Using Large Language Models to Document Code: A First Quantitative and Qualitative Assessment

Ian Guelman ianguelman@dcc.ufmg.br Federal University of Minas Gerais (UFMG) Belo Horizonte, Minas Gerais, Brazil

Laerte Xavier

laertexavier@pucminas.br Pontifical Catholic University of Minas Gerais (PUCMG) Belo Horizonte, Minas Gerais, Brazil

ABSTRACT

Code documentation is vital for software development, improving readability and comprehension. However, it's often skipped due to its labor-intensive nature. AI Language Models present an opportunity to automate the generation of code documentation, easing the burden on developers. While recent studies have explored the use of such models for code documentation, most rely on quantitative metrics like BLEU to assess the quality of the generated comments. Yet, the applicability and accuracy of these metrics on this scenario remain uncertain. In this paper, we leveraged OpenAI GPT-3.5 to regenerate the Javadoc of 23,850 code snippets with methods and classes. We conducted both quantitative and qualitative assessments, employing BLEU alongside human evaluation, to assess the quality of the generated comments. Our key findings reveal that: (i) in our qualitative analyses, when the documents generated by GPT were compared with the original ones, 69.7% were considered equivalent (45.7%) or required minor changes to be equivalent (24.0%); (ii) indeed, 22.4% of the comments were rated as having superior quality than the original ones; (iii) the use of quantitative metrics is susceptible to inconsistencies, for example, comments perceived as having higher quality were unjustly penalized by the BLEU metric.

CCS CONCEPTS

• Software and its engineering \rightarrow Documentation; • Computing methodologies \rightarrow Natural language generation.

KEYWORDS

Code Documentation, Large Language Models, GPT, BLEU Score

1 INTRODUCTION

Documentation plays an important role on effective software development [6], with comments standing out as the principal resource to document code. Code comments are, in short, explanations added to the source code to briefly describes a particular code snippet, improving the readability and comprehension of the code [14]. Writing meaningful comments can be labor-intensive and time-consuming, making developers skip this important step of software development [6]. On the other hand, it is clear that practitioners expect code snippets to be commented and are willing to use automated code comment generation tools [4].

Arthur Gregório Leal arthurgregorioleal@gmail.com Pontifical Catholic University of Minas Gerais (PUCMG) Belo Horizonte, Minas Gerais, Brazil

Marco Tulio Valente mtov@dcc.ufmg.br Federal University of Minas Gerais (UFMG)

Belo Horizonte, Minas Gerais, Brazil

Recently, the advancements of AI Language Models present an opportunity to automate the generation of code documentation, making the process of writing comments seamless, easing the burden on developers. However, there are few studies on the quality of the comments generated by Generative Pre-training Transformers (GPT). An exception is a study by Khan and Uddin [6], in which they use CodeSearchNet [5] dataset to evaluate the use of GPT-3 for automatic code documentation. However, in this study, the authors only conduct a quantitative evaluation, centered on a metric called Bilingual Evaluation Understudy (BLEU) [10], which measures the similarity between texts. Nonetheless, it is challenging to generate conclusive and practical results from this study, due to the absence of reference values to interpret the presented BLEU results.

Thus, in this short paper, we present the first results of a quantitative and qualitative study on the use of Large Language Models (LLMs), specifically OpenAI GPT-3.5 Turbo, for code documentation. To the best of our knowledge, we are the first to evaluate these comments qualitatively. We used a dataset containing 2,357 class-level and 21,493 method-level comments from three well-known Java projects and compared the comments generated by GPT with those already existing in the code. This was first performed quantitatively, using BLEU, and then qualitatively. For this, a sample of 709 comments was carefully and independently analyzed by two of the paper's authors and classified on a scale of 1 to 4.

Our ongoing research already provided two contributions that we viewed as relevant. First, we show that the BLEU metric has significant shortcomings. For example, it does not yield consistent results when comparing comments at the class level. And, in the case of methods, it fails to properly distinguish GPT-generated comments that require minor and major changes. Second, and due to also conducting a qualitative analysis, we were able to show that the comments generated by GPT are of high quality. For example, 69.7% were considered equivalent to the original comments (45.7%) or required minor changes to be equivalent (24.0%). Additionally, 22.4% of the comments were rated as having superior quality than the original ones.

The remainder of this paper is organized as follows. Section 2 describes the study design. Section 3 describes the quantitative evaluation via BLEU and Section 4 elaborates on the qualitative evaluation through human assessment. Section 5 correlates the

results of these evaluations. Section 6 addresses threats to validity. Section 7 explores similar research papers. Finally, Section 8 concludes.

2 STUDY DESIGN

2.1 Dataset

For our research, we focused on three prominent GitHub repositories primarily written in Java: spring-boot, spring-framework, and guava. They were chosen for their relevance in the Java open-source ecosystem and their ongoing activity, evidenced by consistent contributions from developers.

2.2 Ground-Truth

We extracted all Javadocs and their corresponding code snippets from the selected repositories, yielding 60,689 comments. Following Khan and Uddin [6], we also removed comments that: (i) are not in English; (ii) contained special tokens (e.g., and https://); or (iii) mentioned author details (i.e., using @author tag). Of the remaining comments, we selected those referring to classes and methods, resulting in 23,850 instances, as can be seen in Table 1.

Table 1: Javadocs instances

Repository	Classes	Methods	TOTAL
spring-framework	986	14,321	15,307
guava	591	4,577	5,168
spring-boot	780	2,595	3,375
TOTAL	2,357	21,493	23,850

2.3 Large Language Model

Model. We used OpenAI GPT-3.5 Turbo model, designated as gpt-3.5-turbo-0125 and released on January 25, 2024.

Prompt. We employed one-shot learning—providing a single example with task-specific instructions—to guide GPT into generating Javadocs from a code snippet. We used the prompt presented in Figure 1. We also randomly selected a comment from each granularity (classes and methods) to be used as an example in this prompt. These comments were then removed from our ground-truth.

```
Context: Suppose you are a Java developer who needs to document some [methods/classes] using source code comments following the Javadoc format.

For example, for the following [method/class]:

'java
[CODE_EXAMPLE]

You should generate the following Javadoc comment:
[JAVADOC_EXAMPLE]

Task: Generate a single comment (in Javadoc format) for the following [method/class]

''java
[EXTRACTED_CODE]

In your answers, only include the suggested comment.
```

Figure 1: Prompt structure

Integration. A Python script was employed to automate the integration with the OpenAI GPT API, facilitating the generation of Javadocs. This script submits prompts to the API, and parses the responses to extract only the Javadoc content. The code used in this process is included in the replication package.

3 QUANTITATIVE ASSESSMENT

We first perform a quantitative assessment of the Javadocs generated by GPT. For that, we compare the generated instances against the ones collected as ground-truth (Section 2.2). For each pair of comments (original and GPT-generated), we calculated the BLEU Score. BLEU is a popular metric to evaluate the quality of machine generated texts [3, 13]. In short, it compares the number of matching n-grams (sequences of n words) between a candidate and the reference text. The resulting score ranges from 0 to 1, with higher values indicating greater similarity. Finally, to address the limited overlap of longer n-gram sequences in short Javadocs, we adopted the use of a smoothed BLEU Score [6], ensuring that the evaluation remains sensitive to partial matches, which provides a more accurate comparison.

Figure 2 presents the obtained results. Considering all comments, GPT achieved a score of 27.04%, in the median. In this context, the first and third quartiles are 19.68% and 35.72%, respectively. At class level, GPT reached a median score of 20.93% (first quartile: 16.37%; third quartile: 25.60%). Finally, for comments generated at method level, the results are: 28.20%, in the median; 20.34% and 36.60% in the first and third quartiles, respectively.

Thus, GPT achieved better results for method level comments (statistically significant according to a Mann-Whitney test, assuming a significance level of 95%, which we will also use in the remaining tests reported in this paper). These results might be due to the more specific scope of methods, when compared to classes, facilitating the generation of more accurate documentation.

In assessing BLEU Scores, it is important to note that a score of 30% is generally considered the minimum for acceptable results, while scores above 50% are typically viewed as indicating good quality [8]. Considering the median values, GPT fails to generate a Javadoc that, when compared to the original one, achieves a BLEU Score higher than the 30% threshold.

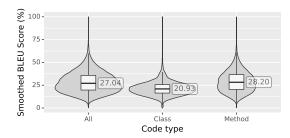


Figure 2: Smoothed BLEU Score

We also compared the size of each pair of comments. GPT-generated Javadocs have a median size that is comparable to the original documentation. Specifically, the median increase in Javadoc length for methods was +30 characters, for classes it was -19 characters, and overall, there was an increase of +25 characters. In terms of lines of code (LOC), there was a median increase of one line for methods, no change for classes, and an overall increase of one line. These differences, both in chars and LOC, are statically significant, as indicated by a Mann-Whitney test.

In comparison to state-of-the-art models explored in the literature, we also concluded that GPT-3.5 Turbo achieved better results. For example, Khan and Uddin [6] previously assessed the performance of various Transformer based models in generating documentation from source code. For Java, they reported the following BLEU Scores for each model: REDCODER [11] with a score of 22.95%, Codex [6] at 22.81%, and CodeBERT [2] achieving 17.65%.

Finding #1: GPT-3.5 achieved, on median, a smoothed BLEU Score of 20.93% for classes and 28.20% for methods, when comparing the generated Javadoc with the original ones. This score is slightly better than benchmarks set in similar studies. For instance, Khan and Uddin [6] reported a score of 22.81% when evaluating a GPT-3 model on documenting Java code.

4 QUALITATIVE ASSESSMENT

Several authors debate whether automated similarity metrics (e.g., BLEU) are effective for evaluating the quality of machine generated texts [3, 9, 12]. While there is a consensus that certain metrics outperform others depending on contextual factors, the overarching agreement is that none can supplant human assessment [1, 14].

Since human evaluation is time-consuming and costly [14], we used random sampling approach to select a representative sample of our dataset, comprising both GPT-generated and original Javadoc comments. With a 95% confidence level and a margin of error of 5%, we randomly selected 331 class-level and 378 method-level Javadoc for human evaluation. Each pair of comments in this sample was carefully evaluated by the first and the second author of this paper, independently. Both are experienced developers with four and six years of experience, respectively. Particularly, they compared the generated Javadocs with the original ones and ranked them with a score ranging from 1 to 4 in a scale inspired by Tran et al. [15].

We refer to this score as **Human Assessment Score**. Its scale is as follows:

- Four: A score of 4 represents that the GPT-generated comment is definitely equivalent and contains the same information as the original one;
- Three: A score of 3 means that minor changes must be performed on the generated comment for it to become equivalent with the original one;
- Two: A score of 2 means that major changes must be performed on the generated comment for it to become equivalent with the original one;
- One: A score of 1 represents that the GPT-generated comment is definitely not equivalent to the original one;

To reduce the bias, after their individual classification, both authors discussed the disagreements to decide a final score for every pair. An overall Cohen's Kappa value of 0.55 was achieved between both authors' classifications, indicating a moderate level of agreement [7]. Additionally, they ranked 64.6% of the comments in the same category, suggesting a reasonable consistency between them. This score reached in consensus is the one analyzed in the remaining of this paper. The consensus process was conducted in two phases: first, after reviewing 10% of the sample to potentially

reduce further disagreements, and after assessing the remaining comments.

Even though the obtained Cohen's kappa coefficient might seem modest at first glance, we claim it is rather adequate, given the subjective nature of comment assessment due to personal interpretation. Notably, as can be seen in Figure 3, 80% of the disagreements occurred within neighboring ranks, with nearly half (45%) happening specifically between ranks 3 and 4. This pattern suggests that while the authors occasionally differed in their assessments, their judgments predominantly fluctuate between neighboring ranks, reflecting a fundamental agreement on the overall similarity of the Javadocs.

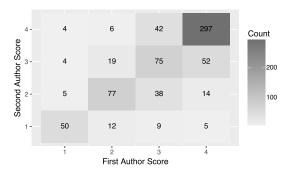


Figure 3: Human Assessment Score Agreement Matrix

As presented in Figure 4, almost the majority of the generated comments (45.7%) received a Score of 4, indicating equivalence to the original Javadoc. This was followed by Scores 3 (24.0%), 2 (21.2%), and 1 (9.2%).

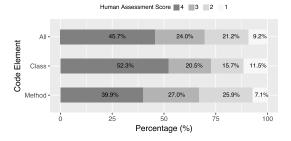


Figure 4: Human Assessment Score Distribution

Interestingly, we also identified automatically generated comments that not only achieved equivalence with the original Javadoc but also augmented the information encapsulated within them. Those are distinguished with a 4+ Score, representing 49% of the Javadocs rated as 4, and 22.4% of the whole sample.

Finding #2: The performance of GPT was highly positive in our qualitative evaluation. For instance, 69.7% of the comments generated by GPT received Scores of 3 (requiring minor changes) or 4 (definitely equivalent). Additionally, 22.4% of the comments received a Score of 4+, indicating that they have a higher quality than the original one.

4.1 Examples

To illustrate and for a better understanding of our results, this section presents examples of comments for each score and discuss why they received their respective classification.

4.1.1 Score 1 (Definitely not Equivalent). The original Javadoc presented in Figure 5 refers to a method that does nothing except by throwing an exception when called. Although GPT also states that an exception will be thrown, it has a clear misunderstanding of what happens with the map, saying that all the values will be replaced, when in reality they will remain unmodified. Thus, it is appropriate to classify it with Score 1, once it is clearly not equivalent with the original comment. Generally, the other generated Javadocs ranked with Score 1 also show a lack of understanding of a key concept of the given code element.

```
Score 1
Original
/**

Guaranteed to throw an exception and leave the map unmodified.

* @throws UnsupportedOperationException always

* @deprecated Unsupported operation.

* GPT-Generated
/**

* Replaces all key-value pairs in the map by applying the specified function, always throwing UnsupportedOperationException.

* @param function the function to apply (not used as the method always throws an exception)

* @deprecated This method is deprecated and should not be called

* @DoNotCall(""Always throws UnsupportedOperationException"")

* @throws UnsupportedOperationException always thrown as this operation is not supported

*/
```

Figure 5: Human Assessment Score-Example Score 1

- 4.1.2 Score 2 (Major Changes). Both original and generated Javadocs presented in Figure 6 refer to a waiting time before shutting down. However, the original Javadoc is far more detailed, also mentioning what happens if a task is submitted during this period, the default waiting time and how to overwrite it. Since many important pieces of information are missing and the generated comment requires major additions to become equivalent, it was classified with Score 2. Typically, the other Javadocs classified as 2 also present a significant gap in information compared to the original comments.
- 4.1.3 Score 3 (Minor Changes). The original and the generated Javadocs in Figure 7 clearly mentions the functionality of the documented method: reversing the elements in an array. While the generated one goes straight to the point, the original one cites an equivalent method and discusses its efficiency. Although the generated comment lacks some information, those are minor and are not relevant nor required to the understanding of the method. Hence, it was appropriately classified with Score 3. Commonly, Javadocs that received Score 3 provide the same information about the snippet functionality as the original comment, but with less detail on the reasons behind certain choices.
- 4.1.4 Score 4 (Definitely Equivalent). The original and the generated Javadocs presented in Figure 8 are almost identical, they provide the same information, the only difference being subtle changes in the phrasing within the documentation. Since both are definitely

```
Original
* Configure the amount of time we'll wait before shutting down
 If a task is submitted during the {@code shutdownQuietPeriod},
it is guaranteed
* to be accepted and the {@code shutdownQuietPeriod} will start
over.
 By default, this is set to
* {@link LoopResources#DEFAULT_SHUTDOWN_QUIET_PERIOD} which is 2
seconds but
  can also be overridden with the system property
* {@link reactor.netty5.ReactorNetty#SHUTDOWN_QUIET_PERIOD
  ReactorNetty . SHUTDOWN_QUIET_PERIOD } .
 @see #setShutdownTimeout(Duration)
GPT-Generated
* Set the duration of the quiet period before shutdown.
 @param shutdownQuietPeriod the duration of the quiet period
```

Figure 6: Human Assessment Score-Example Score 2

```
Score 3
Original
/**

Reverses the elements of {@code array}. This is equivalent to {
@code

* Collections.reverse(Longs.asList(array))}, but is likely to be more efficient.

*

* @since 23.1

*/
GPT-Generated
/**

* Reverse the elements in the given array.

* @param array the array to be reversed

*/
```

Figure 7: Human Assessment Score-Example Score 3

equivalent, we ranked the GPT-generated comment with a Score 4. All the generated Javadocs classified with Score 4 share this characteristic of being nearly identical to their original counterparts, with only subtle linguistic variations.

```
Score 4
Original
/**

* Register a {@link DeferredResultProcessingInterceptor} under the given key.

* @param key the key

* @param interceptor the interceptor to register

*/
GPT-Generated
/**

* Register a {@link DeferredResultProcessingInterceptor} with a given key.

*

* @param key the key to associate with the interceptor

* @param interceptor the interceptor to be registered

*/
```

Figure 8: Human Assessment Score-Example Score 4

4.1.5 Score 4+ (Definitely Equivalent, Containing Additional Information). Although both the original and the generated Javadocs presented in Figure 9 mentions that the given class is a 'simple property descriptor', the comment generated using GPT is far more comprehensive, providing additional information on the class's parent, properties, and methods. This information is really useful for

understanding a method by only reading the documentation. Thus, we decided to classify the generated Javadoc with the distinguished Score of 4+. Every Javadoc in our sample distinguished with Score 4+ share this same trait of providing useful additional information of the snippet compared to the original documentation.

```
Score 4+
Original
/**

* A simple {@link PropertyDescriptor}.

*/
GPT-Generated
/**

* A class representing a simple property descriptor that includes information about the read method, write method,

* property type, and property editor class. This class extends
PropertyDescriptor and provides methods to

* get and set the read and write methods, property type, and
property editor class. It also includes methods

* to determine equality, calculate hash code, and generate a
string representation of the object.

*/
```

Figure 9: Human Assessment Score-Example Score 4+

5 CORRELATING THE RESULTS

Figure 10 presents the relationship between the Human Assessment and BLEU Scores on the level of methods. As expected, both metrics presents an aligned upward trend, except for the 4+ category. This is expected, as the 4+ Javadocs are characterized by their comprehensive content, surpassing the original documentation in detail. Unfortunately, this fact leads to a decrease in textual similarity, impacting the BLEU Scores.

Indeed, for methods, the correlation between Human Assessment Score and BLEU is weakly positive (Spearman: 0.25) and statistically significant when the 4+ category is omitted. However, considering the whole sample yields a further weakened correlation (Spearman: 0.08), although not statistically significant.

In Figure 10, we can also see that BLEU does not distinguish between the comments ranked with Score 2 (major changes required), from the ones with Score 3 (minor changes required to make the generated comment equivalent) neither from the ones with Score 4 (definitely equivalent). Indeed, the BLEU of the comments placed in these categories have similar medians (27.62%, 30.32%, and 32.72%, respectively). Moreover, BLEU fails to distinguish between the two most extreme categories: Score 1 (definitely not equivalent) and Score 4+ (definitely equivalent and containing additional information). Their medians are, respectively, 18.11% and 20.17%. A Mann-Whitney test supports these observations, stating that there are no statistically significant differences in BLEU Scores between the members of each mentioned group.

Finally, when analyzing the results for classes, as presented in Figure 11, we found that BLEU cannot differentiate comments ranked with Score 1 (definitely not equivalent) from those with Score 2 (major changes required) and from the ones with Score 3 (minor changes required to make the generated comment equivalent). Their BLEU Score medians are notably similar—20.82%, 21.59%, and 21.47%, respectively. On the other hand, comments ranked with Scores 4 or 4+ are indeed statistically different from any other rank. These findings are also supported by a Mann-Whitney test.

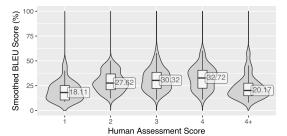


Figure 10: BLEU by Human Assessment Score on Methods

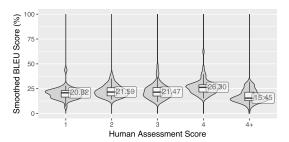


Figure 11: BLEU by Human Assessment Score on Classes

5.1 Examples of BLEU Scores

To better illustrate the relationship between the Human Assessment and BLEU Score, Table 2 presents the values obtained for each metric across the examples discussed in Section 4.1. Notably, the example in Figure 8 stands out with the only BLEU Score that surpassed the 50% threshold [8], at 60.26%. Similarly, the metrics for the examples in Figures 5 and 6 are in agreement, as both are notably low.

Table 2: Human Assessment Score vs. BLEU

Figure	Human Assessment Score	BLEU Score
Fig. 5	1	14.55%
Fig. 6	2	26.78%
Fig. 7	3	25.67%
Fig. 8	4	60.26%
Fig. 9	4+	7.21%

On the other hand, while the example in Figure 7 has a relative good Human Assessment Score (Score 3)—suggesting the GPT-generated documentation requires only minor changes to be equivalent to the original—its BLEU Score is low, at 25.67%, thus failing to reach the 30% threshold [8]. This discrepancy is more pronounced for the example in Figure 9. Despite achieving the best Human Assessment Score (Score 4+), it has the lowest BLEU Score at 7.21%. Such inconsistencies, while somewhat expected due to BLEU's reliance on n-gram matching, highlight BLEU's limitations in contexts demanding a deeper understanding of the underlying meaning, where semantic outweighs syntactic similarity.

Finding #3: Our qualitative analysis revealed at least two critical problems related to using a metric like BLEU to evaluate the quality of the comments generated by GPT. First, BLEU was not able to properly distinguish between all of the Human Assessment Score categories. Second, comments having a better quality than the original ones (Score 4+, 22.43% of our sample) were unjustly penalized by the BLEU metric.

6 THREATS TO VALIDITY

There are at least two threats to the validity of our current results. Firstly, GPT may have had prior exposure to the repositories under analysis due to its training on open-source projects. Nevertheless, the model used was trained on data available only up until September 2021, minimizing this threat since it is likely that substantial updates to the repositories have occurred since this date. Lastly, the extracted Javadocs were used as ground truth for comparison. However, there is a chance that they do not exhibit high quality. To address this concern, a sample of 240 comments was initially analyzed, revealing that more than 93% met the appropriate quality standards. Notably, the selected repositories are maintained by reputable companies, further supporting the suitability of the evaluated Javadocs.

7 RELATED WORK

Few previous studies in the literature have explored the adoption of LLMs to automatic generation of code documentation. For example, Khan and Uddin [6] explored the potential of GPT-3 based Codex for this purpose, highlighting its superiority over existing methods. They evaluated Codex's performance across six programming languages, finding that it outperforms existing models with an average BLEU Score of 20.6%. Their findings suggest Codex-generated documentation to be comparable to human-written in terms of readability and informativeness. They emphasized Codex's ability to generate a more comprehensible documentation, also providing additional information. However, their study does not include a qualitative assessment based on human judgment, similar to the one we conducted and described in this paper.

In addition, Hu et al. [4] conducted a comprehensive study to understand the expectations of practitioners regarding automated code comment generation. The authors performed a mixed-methods study—including interviews, survey, and literature review—to gather insights into current commenting practices, issues, and the desirability of automated commenting tools. Their findings indicate that while most practitioners recognize the value of automated comment generation, there is a significant gap between existing tools and practitioners needs. However, their study did not consider recent AI-models, such as LLMs, which were not yet popular at the study time. By contrast, in our study we leveraged OpenAI's GPT model to generate comments for real Java code.

Finally, Evtikhiev et al. [1] conducted an in-depth study to evaluate the applicability of various metrics—including BLEU, ROUGE-L, METEOR, ChrF, CodeBLEU, and RUBY—for assessing the quality of code generation models. Their findings revealed that none of the metrics reliably emulated human judgment, particularly when

differences in model scores were minimal. While their study focuses on code generation, our paper addresses the generation of comments (for the purpose of documenting source code).

8 CONCLUSION AND FUTURE WORK

Our key findings are as follows:

- (1) We concluded that GPT-3.5 Turbo is capable of generating high quality comments. In our qualitative assessment, 45.7% of the generated comments were considered equivalent to the original ones. Indeed, 22.4% of them were assessed as superior in quality than the originals.
- (2) When assessing the efficacy of automated evaluation metrics, we found significant shortcomings with BLEU. At least in our context, BLEU showed relevant inconsistencies. Those highlight the risks of relying solely on automated metrics for evaluating the quality of comments generated by tools such as GPT.

As future work, we plan to extend our dataset to consider projects in other programming languages and comments associated with other code elements, such as statements, declarations, and packages. We plan to correlate the results of the quantitative analysis with other metrics, such as SIDE, METEOR, ChrF, and ROUGE-L. In the qualitative analysis, we intend to extend our sample and also to include evaluations of another expert. Finally, we plan to evaluate GPT's effectiveness in automating documentation for an open-source project, covering all contribution steps, including code reviews, where developers will assess the quality and appropriateness of the generated documentation.

Replication Package: https://doi.org/10.5281/zenodo.11480759

REFERENCES

- [1] Mikhail Evtikhiev, Egor Bogomolov, Yaroslav Sokolov, and Timofey Bryksin. 2023. Out of the BLEU: how should we assess quality of the code generation models? *Journal of Systems and Software* (2023).
- [2] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. Codebert: A pre-trained model for programming and natural languages. arXiv preprint arXiv:2002.08155 (2020).
- [3] Markus Freitag, David Grangier, and Isaac Caswell. 2020. BLEU might be guilty but references are not innocent. arXiv preprint arXiv:2004.06063 (2020).
- [4] Xing Hu, Xin Xia, David Lo, Zhiyuan Wan, Qiuyuan Chen, and Thomas Zimmermann. 2022. Practitioners' expectations on automated code comment generation (ICSE '22). ACM, 1693–1705.
- [5] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Codesearchnet challenge: Evaluating the state of semantic code search. arXiv preprint arXiv:1909.09436 (2019).
- [6] Junaed Younus Khan and Gias Uddin. 2023. Automatic Code Documentation Generation Using GPT-3. In 37th IEEE/ACM International Conference on Automated Software Engineering (ASE).
- [7] J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. biometrics (1977), 159–174.
- [8] Alon Lavie. 2010. Evaluating the output of machine translation systems. In Proceedings of the 9th Conference of the Association for Machine Translation in the Americas: Tutorials.
- [9] Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. arXiv preprint arXiv:2006.06264 (2020).
- [10] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In 40th Annual Meeting of the Association for Computational Linguistics. 311–318.
- [11] Md Rizwan Parvez, Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval augmented code generation and summarization. arXiv preprint arXiv:2108.11601 (2021).

- [12] Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. 2021. Reassessing automatic evaluation metrics for code summarization tasks. In 29th Symposium on the Foundations of Software Engineering (FSE). 1105–1116.
- [13] Ensheng Shi, Yanlin Wang, Lun Du, Junjie Chen, Shi Han, Hongyu Zhang, Dongmei Zhang, and Hongbin Sun. 2022. On the evaluation of neural code summarization. In 4th International Conference on Software Engineering (ICSE). 1597–1608.
- [14] Xiaotao Song, Hailong Sun, Xu Wang, and Jiafei Yan. 2019. A Survey of Automatic Generation of Source Code Comments: Algorithms and Techniques. *IEEE Access* 7 (2019), 111411–111428.
- [15] Ngoc Tran, Hieu Tran, Son Nguyen, Hoan Nguyen, and Tien Nguyen. 2019. Does BLEU Score Work for Code Migration?. In 2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC). 165–176.