

An Empirical Study on Code Comment Completion

Antonio Mastropaolo, Emad Aghajani, Luca Pascarella, Gabriele Bavota
SEART @ Software Institute, Università della Svizzera italiana (USI), Switzerland

Abstract—Code comments play a prominent role in program comprehension activities. However, source code is not always documented and code and comments not always co-evolve. To deal with these issues, researchers have proposed techniques to automatically generate comments documenting a given code at hand. The most recent works in the area applied deep learning (DL) techniques to support such a task. Despite the achieved advances, the empirical evaluations of these approaches show that they are still far from a performance level that would make them valuable for developers. We tackle a simpler and related problem: Code comment completion. Instead of generating a comment for a given code from scratch, we investigate the extent to which state-of-the-art techniques can help developers in writing comments faster. We present a large-scale study in which we empirically assess how a simple n -gram model and the recently proposed Text-To-Text Transfer Transformer (T5) architecture can perform in autocompleting a code comment the developer is typing. The achieved results show the superiority of the T5 model, despite the n -gram model being a competitive solution.

Index Terms—Empirical Study, Code Comments

I. INTRODUCTION

Code comprehension can take up to 58% of developers' time [1]. In such a process, code comments play a pivotal role [2] helping developers in understanding the source code at hand. However, as shown in recent studies, code comments may completely lack [3] or not co-evolve with the related code [4]–[7], even becoming misleading for developers.

To support developers in code comprehension activities, researchers have proposed techniques and tools aimed at automatically documenting a given code at hand. These approaches can be roughly classified into two categories: extractive [8]–[11] and abstractive [10], [12]–[17]. The former create a summary of a code component which includes information extracted from the component being summarized (*e.g.*, a class is summarized by exploiting a set of predefined templates filled in with information extracted from the class code). While simple to implement and efficient in terms of execution time, these approaches fall short in case the code component uses a poor vocabulary. Abstractive approaches try to overcome this limitation by including in the generated summaries information that is not present in the code component to document. Among those techniques, the ones using deep learning (DL) to generate natural language descriptions for a given snippet of code are on the rise [13]–[17], with attempts also made for automatically updating comments given a code change [18].

Despite the substantial improvements brought by DL techniques in addressing the *code comment generation* problem, the findings reported even in the most recent empirical studies show how these techniques are still far from being useful tools for software developers.

For example, in the recent work by Mastropaolo *et al.* [19], the authors show that state-of-the-art techniques are able to generate comments equivalent to those written by humans in only $\sim 10\%$ of cases. For this reason, we believe that completely relying on “*machines*” to write comments is, as of today, a far-fetched goal that, while worth investigating, is unlikely to result in major breakthrough in the short term. In this work, we tackle the simpler problem of *code comment completion*, in which the “*machine*” is in charge of completing a comment that the developer starts writing, similarly to what done for code tokens by code completion techniques [20]–[23].

The code comment completion problem has been firstly tackled by Ciurumelea *et al.* [24] in the context of Python code: They study whether a deep learning model can predict the next word that a developer is likely to type while commenting code. This is, to the best of our knowledge, the only work done in this area. Stemming from their idea, we present in this paper a large-scale study assessing the ability of a simple n -gram model and of the recently proposed Text-To-Text Transfer Transformer (T5) architecture [25] in supporting code comment completion for Java programs. As compared to the work by Ciurumelea *et al.* [24], besides focusing on a different context (*i.e.*, Python vs Java) we: (i) investigate the actual advantages brought by a DL-based model (T5) over a simpler n -gram model that can be trained in a fraction of the time required by T5; (ii) do not limit our study to predicting the single next word the developer is likely to type, but evaluate how the investigated techniques perform when asked to predict longer word sequences (*e.g.*, the next 10 words), providing a more advanced completion support to developers. We also study the complementarity of the two techniques and report qualitative examples of correct and wrong predictions to understand their strengths and limitations.

Our study has been run on a dataset composed by 497,328 Java methods with their related comments. The achieved results can be summarized as follows. First, the T5 model outperforms the n -gram model, achieving superior performance in all the comment completion scenario we tested. Second, despite being more performant, the T5 model exploits as input not only the first part of the comment already written by the developer (also used by the n -gram model), but also a *context* representing the relevant code for the comment to complete. This means that the T5 model, as we tested it, can only be used when the developer writes the comment after the code has been already implemented (the assumption made by approaches for automated code documentation [13]–[17]). Thus, the applicability of the n -gram model is higher (*i.e.*, it can be used also when the code is not yet implemented).

II. T5 TO SUPPORT CODE COMMENT COMPLETION

The T5 model has been introduced by Raffel *et al.* [25] to support multitask learning in Natural Language Processing (NLP). The idea is to reframe NLP tasks in a unified text-to-text format in which the input and output are always text strings. For example, a single model can be trained to translate across languages and to autocomplete sentences. This is possible since both tasks can be represented in a text-to-text format (*e.g.*, in the case of translation, the input is a sentence in a given language, while the output is the translated sentence). The T5 is trained in two phases: *pre-training*, which allows defining a shared knowledge-base useful for a large class of sequence-to-sequence tasks (*e.g.*, guessing masked words in English sentences to learn about the language), and *fine-tuning*, which specializes the model on a specific downstream task (*e.g.*, learning the translation of sentences from English to German). We briefly overview the T5 model and explain how we adapted it for supporting code comment completion. Finally, we describe the hyperparameter tuning of the model and the decoding strategy for generating the predictions.

A. An Overview of T5

The T5 is based on the transformer model architecture that allows handling a variable-sized input using stacks of self-attention layers. When an input sequence is provided, it is mapped into a sequence of embeddings passed into the encoder. Raffel *et al.* [25] propose five variants of T5: *small*, *base*, *large*, *3 Billion*, and *11 Billion*. These variants differ in terms of architectural complexity, with the smaller model having 60M parameters against the 11B of the largest one. While the accuracy of the most complex variants is higher as compared to the less complex models, the training complexity increases with the number of parameters [25]. Considering our computational resources, we decided to use the simplest T5_{small} model and, for this reason, we expect the achieved results to be a lower bound for the performance of a T5-based model in the task of code comment completion.

T5_{small} architectural details. The T5_{small} architecture is characterized by six blocks for encoders and decoders. The feed-forward networks in each block consist of a dense layer with an output dimensionality (d_{ff}) of 2,048. The *key* and *value* matrices of all attention mechanisms have an inner dimensionality (d_{kv}) of 64, and all attention mechanisms have eight heads. All the other sub-layers and embeddings have a dimensionality (d_{model}) of 512.

B. Problem Definition

We instantiate the T5 to the problem of code comment completion in Java. We tackle the problem at method-level granularity, meaning that we expect the model to learn how to autocomplete a code comment used to document a method or part of it. In Java, a method can be documented using a Javadoc comment (that we indicate with C_{JD}) or inner comments (C_I). Each comment is relevant to a specific *context*. For example, the context of a C_{JD} is the entire method, while the context of an C_I can be a single line or a code block.

Given a *context* and an incomplete comment (either a C_{JD} or a C_I), the trained model must predict the tokens needed to complete the comment. This implies that we must build a training dataset in which code comments are linked to the relevant part of the code they document (*i.e.*, the context). While this is trivial for C_{JD} comments, a heuristic is needed for C_I comments, since they are not explicitly linked to certain statements. Section II-C describes how we built such a dataset.

We pre-train the T5 by randomly masking tokens in comments asking the model to guess the masked tokens (Section II-D). This builds the shared knowledge that we then specialize in two fine-tuning tasks, namely *Inner-comment_{task}* and *Javadoc_{task}*, consisting in predicting the missing part of inner and Javadoc comments, respectively (Section II-E).

C. Dataset Preparation

We start from the *CodeSearchNet* dataset [26], providing 6M functions from open-source projects. We only focus on the Java subset, composed of $\sim 1.5M$ methods. We extract the set of instances using the *docstring* (*i.e.*, the method's *Javadoc*) and *function* fields. Then, we run a pre-processing aimed at preparing our dataset.

First, we discarded instances having $\#tokens \geq 256$, where $\#tokens = \#function_tokens + \#docstring_tokens$ (*i.e.*, $\#tokens$ is the total number of tokens used to represent both the method and the comments associated to it). Such a filter is needed to limit the computational expense of training the model, and removed $\sim 9\%$ of instances from the dataset. Then, we excluded instances containing non ASCII characters as well as comments composed by less than three tokens (words), since unlikely to represent an interesting scenario for code comment completion. We also excluded comments representing instances of self-admitted technical debt (SATD) [27] (*i.e.*, comments documenting temporary workaround). Such a choice was dictated by the fact that we are interested in training our model to complete comments describing a method (or part of it) rather than comments used to document technical debt and very likely to be project-specific. We adopted a simple heuristic to discard SATD comments, excluding all comments starting with *TOFIX*, *TODO*, and *FIXME*. We are aware that more complex state-of-the-art techniques for SATD detection could be used (see *e.g.*, [28]), but we preferred a simpler unsupervised heuristic for our pre-processing pipeline.

We discarded commented code statements using the *code-type* library [29]. Then, we cleaned comments by replacing links with a special *_LINK_* token (using the *urlextract* library [30]), dates with a *_NUM_* token (using the *datefinder* library [31]), and references to code components with a *_REF_* token. The latter are only handled in Javadoc comments exploiting the *@link* tag used to reference code elements. We further clean Javadoc comments by stripping *HTML/XML* tags using the *BeautifulSoup* library [32].

Then, we removed from each method all C_I (inner comments) that are “orphans”, *i.e.*, C_I followed **and** preceded by at least one blank line. As previously explained, to train the T5, we need to link each comment to its *context*.

Javadoc comments are linked to the whole method, while for inner comments we adopt a heuristic that would not work with orphan comments (*i.e.*, we cannot know what lines of code they likely document) — details in Section II-E. As a last step in the processing of inner comments, we merge in a single C_I subsequent inline comments that are not interleaved by empty lines or code statements. This is done since they likely represent a single multi-line comment.

Type of instance	#Instances Pre-training	#Instances Fine-tuning
Inner comment(s) only (D1)	45,764	33,590
Javadoc comment only (D2)	232,121	115,904
Javadoc & Inner comments (D3)	53,667	16,282
Total	331,552	165,776

TABLE I: Instances used for pre-training and fine-tuning.

Finally, we removed duplicates, obtaining the study dataset composed of 497,328 instances, with each instance being a method with its related Javadoc and/or inner comments.

We randomly split this dataset using 2/3 of it for pre-training and 1/3 for fine-tuning. Table I shows the number of instances in each of the two datasets, distinguishing between instances only containing inner comments (D1), only containing Javadoc (D2), and featuring both (D3).

D. Pre-training of T5

In the *pre-training* phase, we use a self-supervised task similar to the one used by Raffel *et al.* [25], consisting of masking tokens in natural language sentences and asking the model to guess the masked tokens. Since we want the model to learn how to generate comments given a certain context, we randomly mask 15% of tokens appearing in the comment-related part of each instance (Javadoc or inner comments). Tokens representing the method code were not masked. The pre-training has been performed for 200k steps.

We also created a new *SentencePiece* [33] (*i.e.*, a tokenizer for neural text processing) model by training it on the entire pre-training dataset, in such a way that it can handle both code and comments. We set its length to 32k wordpieces.

E. Fine-tuning of T5

Once pre-trained, we fine-tune the T5 model in a multi-task setting, in which the two tasks are represented in our case by the automatic completion of (i) Javadoc comments and (ii) inner comments. Having a single model fine-tuned on these two strongly related tasks could result in an effective transfer learning, in which knowledge gained by the model while learning a task (*e.g.*, $Javadoc_{task}$) can be transferred to other similar tasks (*e.g.*, $Inner-comment_{task}$).

1) *Preparing the Dataset for the Model Fine-Tuning:* We further process the 165,776 instances selected for the fine-tuning (see Table I) through the following steps.

Processing Javadoc instances (datasets D2 and D3 in Table I). For the $Javadoc_{task}$ we assume that the *context* documented by a C_{JD} is represented by the whole method.

Thus, given an instance composed by $\{C_{JD}, context\}$ (*i.e.*, a Javadoc comment followed by the method it documents) we simply swap the position of the two elements and add a special separation token $\langle sep \rangle$ to delimit the comment, obtaining an instance in the form: $\{context\langle sep \rangle C_{JD}\langle sep \rangle\}$. The rationale behind this transformation is to “force” the model to process the context before predicting the missing parts of the comment.

Once this is done, we use the method `sent_tokenize` from the *nlk* library [34] to split the C_{JD} in each instance into the y sentences composing it. Then, we take the first sentence s_1 and remove the remaining $y - 1$. Assuming s_1 is composed by n tokens, we randomly extract five different integers between 1 and $n - 1$, and use them to create five variants of s_1 each one having the last $n - m_i$ tokens masked, where m_i is one of the five random integers.

By training the model on these five masked sentences, the learning is focused on guessing how to finalize an incomplete sentence in a comment the developer is writing. Let us justify and explain this process. First, we remove the $y - 1$ following sentences because we assume that a developer writes the comment linearly, starting from the first to the last sentence. Thus, when the developer is writing the first sentence, the remaining $y - 1$ do not exist yet. Second, at most $n - 1$ tokens can be masked in a sentence composed by n tokens, since at least the first token must be provided by the developer (otherwise, the task would be comment generation rather than completion). Third, the choice of creating five different variants for a given sentence is a tradeoff between experimenting with a different number of masked tokens for each sentence and considering all possible combinations of masked tokens, that would lead to an excessive number of training instances.

Such a process is repeated for all y sentences composing the C_{JD} , hiding the sentences following the one under process while keeping the ones preceding it. Thus, for each instance in our fine-tuning dataset, we create up to $y \cdot 5$ instances (*i.e.*, y sentences with five different sets of masked tokens). Less than five instances are created if C_{JD} has less than six tokens, since it would not be possible to mask five different sets of tokens (excluding the first one). Fig. 1 shows an example of masking performed on a single instance. The original instance is reported on top of the figure, and two sentences compose its C_{JD} . This results in the creation of 10 fine-tuning instances. Due to space constraints, Fig. 1 only shows one of the instances generated for each sentence.

Processing Inner-comment instances (datasets D1 and D3 in Table I). For the $Inner-comment_{task}$ the first step to perform when preparing the fine-tuning dataset is the identification of the *context* relevant for a given inner method. We define the following heuristic to identify, given an C_I in a specific instance, the context that can help the model in understanding how to support C_I completion. We use Fig. 2 as a running example to explain the heuristic, showing an example of instance having a single C_I . Starting from the line l_{C_I} containing it, we expand the context both above and below it until specific conditions are met.

Javadoc masking

```
Scroll the screen to the left.
The underlying application should have at least one scroll view belonging
to the class 'android.widget.ScrollView'.
public void scrollLeft(){
    logger.entering();
    WebElement webElement = this.findElement(
        By.className(SCROLLVIEW_CLASS)
    );
    swipeLeft(webElement);
    Logger.exiting();
}

S1: Scroll the screen <MASK>
...
S2: Scroll the screen to the left. The underlying application <MASK>
...
```

Fig. 1: Example of the Javadoc masking process

Identifying the relevant context

```
public void addMBeanServers(Set<MBeanServerConnection> servers){
    // Example of inner comment within a method
    if (!isJBoss()) {
        InitialContext ctx;
        try{
            ctx = new InitialContext();
            MBeanServer server =
                (MBeanServer) ctx.lookup("java:comp/env/jmx/runtime");
            if (server != null) {
                servers.add(server);
            }
        } catch (NamingException e) {
            ...
        }
    }
}
```

Fig. 2: Identification of relevant context of an inner comment

In particular, while expanding above and below l_{C_i} , we stop when we find one of the following: (i) an empty line, (ii) a closing curly brace, or (iii) another code comment. In both cases, we do not expand the context out of the method (e.g., if none of the above conditions is met while expanding above l_{C_i} , we stop at the method signature). In the example shown in Fig. 2, there is no above context since we immediately hit an empty line, while the context below is stopped when we find the first closed curly brace. It is important to clarify that the *context* we identify is not necessarily the part of code documented by C_i . However, our interest is to provide the model with the relevant code surrounding C_i , to allow it exploiting useful information to support the comment completion. Once linked each C_i to its context, we perform the same processing previously described for the C_{JD} instances (i.e., splitting the comment into sentences and creating five variants of each sentence, each having a different number of tokens masked at the end of it).

Data sources	Train	Eval	Test
<i>Javadoc_{task}</i>	1,398,135	174,624	175,084
<i>Inner-comment_{task}</i>	272,944	34,705	34,138
Total	1,671,079	209,329	209,222

TABLE II: Instances used for the fine-tuning

2) *Dataset Splitting*: Table II shows the fine-tuning dataset we obtained from the above-described process. We split it into 80%-10%-10% for train, test and validation, respectively. The dataset for the *Javadoc_{task}* dominates, in terms of size, the one for the *Inner-comment_{task}*.

This could result in an unbalanced effectiveness of the model for the two tasks. In other words, the model could perform better on the *Javadoc_{task}* and being less effective on *Inner-comment_{task}*. However, as pointed out by Arivazhagan *et al.* [35], there is no free lunch in choosing the balancing strategy when training a multi-task model, with each strategy having its pros and cons (e.g., oversampling of less represented datasets negatively impacts the performance of the most representative task). For this reason, we decided not to perform any particular adaptation of our training set but to follow the true data distribution when creating batches.

F. Decoding Strategy

The given output layer's values allow different possible decoding strategies to generate the output token streams. We adopt a *greedy decoding* strategy since it we believe it is better suited for the problem we are tackling. Indeed, it provides the developer with the most likely completion (rather than with a set of completions as, for example, a *beam search* would do. Indeed, with multiple options a developer would likely spend more time selecting among different tool suggestions rather than writing the comment themselves.

G. Hyperparameter Tuning

We rely on the configurations used by Mastropaolo *et al.* [19]. Concerning the pre-training, we do not tune the hyperparameters of the T5 model because the pre-training step is task-agnostic, and this would provide limited benefits. Instead, we experiment with four different learning rates schedule for the fine-tuning phase, using the configurations reported in Table III.

Learning Rate Type	Parameters
Constant (C-LR)	$LR = 0.001$
Slanted Triangular (ST-LR)	$LR_{starting} = 0.001$ $LR_{max} = 0.01$ $Ratio = 32$ $Cut = 0.1$
Inverse Square Root (ISQ-LR)	$LR_{starting} = 0.01$ $Warmup = 10,000$
Polynomial Decay (PD-LR)	$LR_{starting} = 0.1$ $LR_{end} = 0.01$ $Power = 0.5$

TABLE III: Learning-rates tested for hyperparameter tuning

We fine-tune the model for 100k steps for each configuration; then, we compute the percentage of perfect predictions (i.e., cases in which the model can correctly predict all masked tokens in the comment) achieved on both tasks in the evaluation set. The achieved results reported in Table IV showed a slight superiority of the *slanted triangular* (column 2) that we use in our study.

Task	C-LR	ST-LR	ISQ-LR	PD-LR
<i>Javadoc_{task}</i>	30.77%	33.60%	31.73%	30.65%
<i>Inner-comment_{task}</i>	9.38%	10.55%	9.70%	9.51%

TABLE IV: Hyperparameter tuning results

III. STUDY DESIGN

The *goal* of this study is to experiment the extent to which a T5 model and an n -gram model can help developers in writing comments faster by supporting code comment completion.

In particular, we answer the following research question: *To what extent can T5 and n -gram models be leveraged to support the auto-completion of code comments?*

We answer this research question by using the 209,222 test set instances in Table II as a context for our study. This means that the two trained models (*i.e.*, T5 and n -gram model) are run on the same test set instances to predict the masked parts of code comments. The T5 model has been trained as described in Section II, while in the following we explain how we implemented and trained the n -gram models. The code and data used in our study are publicly available [36].

A. N -Gram Model

An n -gram model can predict a single token following the $n - 1$ tokens preceding it. We publicly release the implementation of our n -gram model in Python [36]. We trained the n -gram model by using the same instances used to fine-tune the T5 without, however, the masked tokens. We experimented with three different values for n (*i.e.*, $n = 3$, $n = 5$ and $n = 7$). Even though the n -gram model is meant to predict a single token given the $n - 1$ preceding tokens, we designed a fair comparison for the scenario in which we mask more than one token. In particular, we use the n -gram model in the following way: Let us assume that we are predicting, using a 3-gram model, how to complete a sentence having five tokens T , of which the last two are masked (M): $\langle T_1, T_2, T_3, M_4, M_5 \rangle$. We provide as input to the model T_2 and T_3 to predict M_4 , obtaining the model prediction P_4 . Then, we use T_3 and P_4 to predict M_5 , thus obtaining the predicted sentence $\langle T_1, T_2, T_3, P_4, P_5 \rangle$. Basically, all predictions are joined to predict multiple contiguous tokens.

We experimented with the different values of n by running the models on the evaluation set, taking the best one (*i.e.*, 5-gram) and comparing its performance with that of the T5 model. We report the results achieved for the 3-gram and 7-gram models in our replication package [36].

B. Evaluation Metrics and Data Analysis

We compare the T5 and 5-gram models using six metrics.

Perfect predictions: This metric measures the percentage of cases (*i.e.*, instances in the test set) in which the sequence predicted by the model equals the oracle sequence. Since we want to investigate the extent to which the experimented techniques can actually support comment completion, we compute the perfect predictions when the technique is only required to guess the first masked token, the first two, the first three, etc. For example, if we assume that in the previous prediction $\langle T_1, T_2, T_3, P_4, P_5 \rangle$ P_4 is correct while P_5 is wrong, we will consider this as a perfect prediction only when looking at the first token to predict, and as a wrong one when looking at the first two.

We compute the percentage of perfect predictions when trying to predict the first k masked tokens, with k going from 1 to 10 at steps of 1 (*i.e.*, 1, 2, 3, etc.) and for the most challenging scenario in which $k > 10$ (*i.e.*, more than 10 masked tokens must be correctly predicted to consider this prediction as a perfect one).

BLEU score [37]: This metric measures how similar the candidate (predicted) and reference (oracle) texts are. Given a size n , the candidate and reference texts are broken into n -grams, and the algorithm determines how many n -grams of the candidate text appear in the reference text. The BLEU score ranges between 0 (the sequences are completely different) and 1 (the sequences are identical). For both the tasks, we compute the BLEU- $\{1, 2, 3, 4\}$ and their geometric mean (*i.e.*, BLEU-A). Due to the way in which the BLEU- X is computed (*i.e.*, at least X tokens must be part of the prediction task) we only compute the BLEU-A metric when the number of tokens for a given prediction is at least 4. As done for perfect predictions, we report results for different values of k .

Levenshtein distance [38]: To understand the effort needed by developers to convert a prediction generated by the model into a correct comment, we compute the Levenshtein distance at word-level, *i.e.*, the minimum number of word edits (insertions, deletions or substitutions) needed to convert the predicted comment into the reference one. Also in this case, we report results for different values of k .

Overlap metrics: We also compute the complementarity between the T5 and the n -gram model. Let PP_{T5_t} and PP_{NG_t} be the sets of perfect predictions achieved by T5 and the n -gram model, where $t \in \{Inner-comment_{task}, Javadoc_{task}\}$. We compute the following metrics:

$$Shared_t = \frac{|PP_{T5_t} \cap PP_{NG_t}|}{|PP_{T5_t} \cup PP_{NG_t}|}$$

$$OnlyT5_t = \frac{|PP_{T5_t} \setminus PP_{NG_t}|}{|PP_{T5_t} \cup PP_{NG_t}|} \quad OnlyNG_t = \frac{|PP_{NG_t} \setminus PP_{T5_t}|}{|PP_{T5_t} \cup PP_{NG_t}|}$$

$Shared_t$ measures the percentage of perfect predictions shared between the two compared approaches, while $OnlyT5_t$ and $OnlyNG_t$ measure the percentage of cases in which the perfect prediction is only achieved by T5 or the n -gram model, respectively, on the task t .

Confidence Analysis: Both models provide, together with the generated prediction, a score between 0 and 1 indicating the confidence of the prediction, with 1 being the maximum confidence. We check whether the *confidence* of the predictions can be used as a reliable proxy for their “quality”. If this is the case then, a possible implementation of these models into a tool could take advantage of this proxy to automatically filter out low-confidence predictions. For each model, we compute the confidence level for the two sets of “perfect” and wrong predictions comparing their average confidence.

Qualitative analysis of the predictions: To better understand the strengths and weaknesses of the models, we analyze more closely the generated predictions. First, we check what type of words on two models are able to correctly predict.

We do this by performing a *Part-of-Speech (POS) TAG* analysis. For each test set instance we check the POS category of each masked word using 12 POS categories [39] (e.g., adjective, adverb, noun). Then, for each POS category, we compute for both models the percentage of times they were able to correctly predict it. Such an analysis is useful to understand whether the words correctly predicted by the models are mostly trivial ones (e.g., determiners such as “the”, “a”) or also more challenging words representing, for example, nouns. On top of that, we discuss examples of predictions made by the two models.

A statistical comparison between the T5 and the n -gram model is performed using the McNemar’s test [40] and Odds Ratios (ORs) on the perfect predictions they can generate.

IV. RESULTS DISCUSSION

Fig. 3 depicts the results achieved by the T5 and 5-gram model in terms of perfect predictions, BLEU-A score, and Levenshtein distance computed for predictions of different lengths (k). The results for the other metrics (e.g., BLEU-1 to BLEU-4) can be found in our replication package [36]. The middle and bottom parts of Fig. 3 show the results achieved in the *Javadoc_{task}* and *Inner-comment_{task}*, respectively, while the top part aggregates the results of the two datasets.

In both the evaluated tasks (i.e., *Javadoc_{task}* and *Inner-comment_{task}*) as well as overall, T5 outperforms the 5-gram model by a significant span for all metrics considered in our study. This is especially true when the two models are required to predict a limited number of tokens ($k \leq 6$) following the ones already written by the developer in the code comment. For example, when only the subsequent word must be predicted (i.e., $k = 1$), T5 can achieve more than 50% of perfect predictions in the *Javadoc_{task}* and more than 25% in the *Inner-comment_{task}*. The 5-gram model, in both scenarios, achieves less than 16% of perfect predictions.

The difference in performance is particularly remarkable for the *Javadoc_{task}*, in which the T5 model achieves better results compared to the *Inner-comment_{task}*. Such a finding might be due to two important factors. First, as explained in Section II-E2, the fine-tuning dataset used for the *Javadoc_{task}* is larger than the one used for the *Inner-comment_{task}*, thus likely providing more knowledge to the model about the vocabulary usually adopted by developers in Javadoc comments and, more in general, about this specific task. Second, the Javadoc has, by its nature, a more regular structure making use of tags (e.g., @param) that could help the model in better predicting the comment, especially given the fact that the T5 model exploits the relevant code context during the prediction. Still, even on the *Inner-comment_{task}* the T5 model achieves three times more perfect predictions of the 5-gram when predicting up to seven tokens (bottom-left corner of Fig. 3).

When the number of tokens to predict increases, the gap in performance between the two approaches gets thinner. In the most complex scenario in which the models predict more than 10 tokens in the comment, the T5 achieves, overall, 13% of perfect predictions against the 3% of the 5-gram model.

Task (d)	Shared _t	OnlyT5 _t	OnlyNG _t
<i>Javadoc_{task}</i>	17.06%	75.87%	7.07%
<i>Inner-comment_{task}</i>	19.70%	67.78%	12.52%

TABLE V: Perfect predictions overlap between T5 and 5-gram

The difference is smaller for the *Inner-comment_{task}*, with the best approach (T5) achieving only 1% of perfect predictions.

The McNemar’s test always indicates significant differences in terms of perfect predictions ensured by the T5 and the 5-gram model, with ORs ranging between 8.04 for the *Inner-comment_{task}* and 17.56 for the *Javadoc_{task}* (OR=16.79 for the overall dataset). The better performance of T5 is confirmed by the other evaluation metrics we adopted, namely the BLEU-A and the Levenshtein distance. The BLEU-A gap is up to five times in favor of the T5 over the 5-gram, confirming a substantial difference in performance between the two models. On the *Javadoc_{task}*, T5 always achieves a BLEU-A score higher than at least $\sim 16\%$ as compared to that achieved by the 5-gram, independently of the number of tokens the two models are asked to predict. As already observed for the perfect predictions, the BLEU-A gap is much smaller in the *Inner-comment_{task}*, however still showing a plus $\sim 5\%$ in favor of T5 up to six tokens. The difference in performance tends to decrease while increasing the number of tokens to predict.

By focusing on the Levenshtein distance (right part of Fig. 3, the lower the better), we can observe that, as expected, the number of token-level edits needed to convert a prediction into the reference one tends to increase for both models when more tokens are predicted. An analysis of both tasks (i.e., *Javadoc_{task}* and *Inner-comment_{task}*) points out that T5 requires a developer intervention in a lower number of cases than 5-gram. However, a clear conclusion can be drawn by looking at the three graphs in the right part of Fig. 3: When the two models are not able to generate a perfect prediction, the effort required by developers to convert the prediction into the comment they actually want to write might be too high. For example, when predicting the next five tokens the developer is likely to type, the T5 requires, on average, to changes to 3.2 of the predicted tokens.

An important point to discuss in the comparison between T5 and 5-gram is the different datasets used for their training. Indeed, T5 benefited of a pre-training phase in which it exploited additional code that was not made available to the 5-gram during training. Thus, we performed an ablation study on the T5 model by removing its pre-training step and checking to what extent its superior performances are due to the performed pre-training. While details can be found in our replication package [36], we can summarize our findings as follows: The pre-training phase increases the performance of the T5 in terms of perfect predictions in a range going from $\sim 0.5\%$ to $\sim 2\%$, depending on the task and on the k value. Thus, while pre-training is beneficial, the performance of the T5 are still better than those of the 5-gram model even when both models are only trained on the fine-tuning dataset.

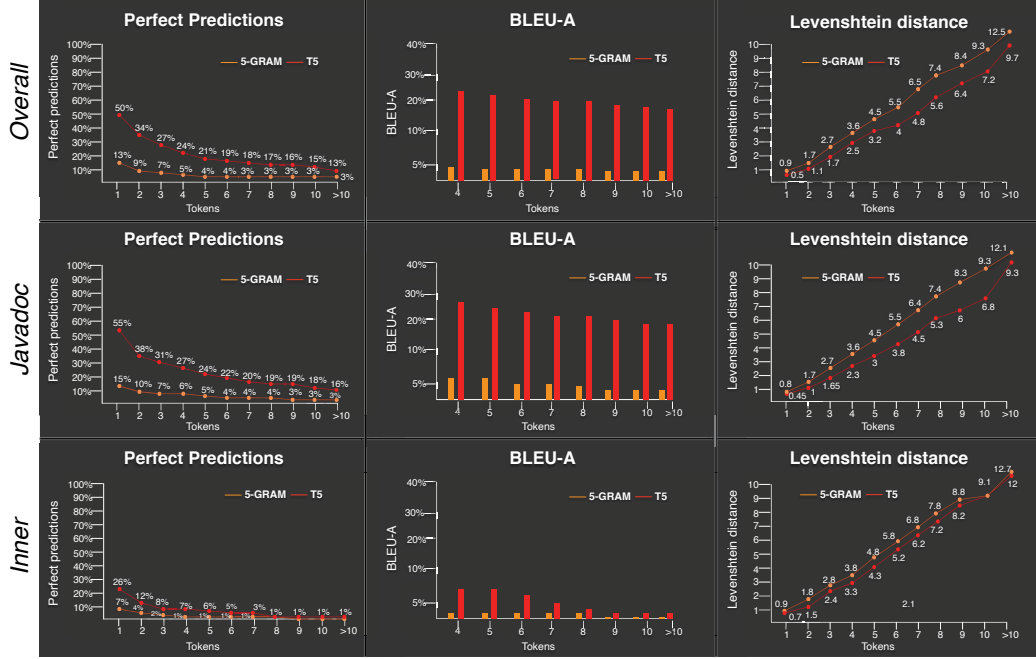


Fig. 3: Performance of the T5 model against the 5-gram model

Table V reports the results of the overlap metrics we computed (see Section III-B). For the *Javadoc_{task}*, only 17.06% of the perfect (*i.e.*, correct) predictions are shared among the two models, while 75.87% are correctly generated only by the T5 model. The 5-gram is responsible for the remaining 7.07% of perfect predictions, that are missed by the T5. This shows, at least for the *Javadoc_{task}*, a limited (but existing) complementarity between the models. Similar results are achieved for the *Inner-comment_{task}*. The two models share 19.70% of perfect predictions, with 67.78% of them correctly predicted by T5 only. The 5-gram model contributes the remaining 12.52% again showing some complementarity between the models.

Fig. 4 shows the POS analysis (*i.e.*, percentage of correct predictions of each POS type) that confirms the superior performance of T5 across all investigated POS categories: Independently from the type of word to predict, the T5 outperforms the 5-gram model. Also, the performances on the *Javadoc_{task}* are, as expected, superior. Not surprisingly, determiners are the ones having the highest percentage of correct predictions. Nevertheless, POS types like nouns, adjective, and verbs which are certainly more challenging to predict still exhibit a good percentage of correct predictions. This analysis, combined with the previous one showing the performance of the model at different values of k , shows that the perfect predictions obtained by the two models (and in particular by the T5) are not only the result of trivial single word ($k = 1$) predictions involving simple POS types, but also include more challenging prediction scenarios. Fig. 5 reports five qualitative examples of predictions performed by both models: the first two are successful predictions made only by the T5 for the *Javadoc_{task}* and the *Inner-comment_{task}*, respectively.

Note that for the first one, the 5-gram does not generate any prediction likely due to the fact that the 4-gram provided as input was never encountered in the training set. The T5, instead, correctly guesses the subsequent 12 tokens in the comment out of the 14 we masked (thus, this is a perfect prediction when considering $k = 12$). The third qualitative example is a prediction failed by both techniques, with the comment generated by the T5 model being more similar to the reference one. Finally, the two predictions at the bottom represent cases in which the 5-gram model correctly completes the comment while the T5 comment fails. Also in this case the first of the two examples is taken from the *Javadoc_{task}*, while the second represents an *Inner-comment_{task}*. The complete list of predictions is publicly available [36].

Overall, our quantitative analysis showed the superiority of the T5 model in the code comment completion task. However, it is important to mention that the T5 model exploits during the prediction the *context* we provide as input (*i.e.*, the code likely to be relevant for the specific comment). This means that, from a practical point of view, such a model can be exploited for code comment completion only assuming that the developer first writes the code and, then, the comment to document it. Clearly, this is not always the case and limits the applicability of the T5 model we experimented with. Such a problem is instead not present in the n -gram model, that can always be applied as long as $n-1$ tokens are present before triggering the prediction of the n^{th} token. A tool integrated into an IDE to support code comment completion could exploit both models: The n -gram model could be triggered if no context can be captured for the comment being written, while the T5 can perform the prediction when relevant code is present.

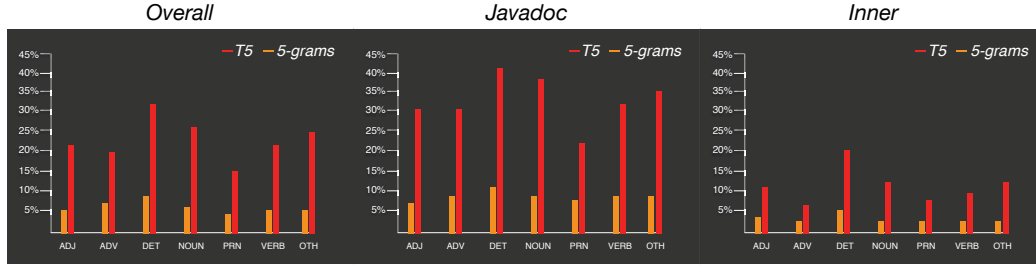


Fig. 4: Part-of-speech tag analysis: ADJ=adjective, ADV=adverb, DET=determiner, PRN=pronoun, OTH=others which includes conjunction, number, particle, adposition, and X, where X are words that cannot be assigned a part-of-speech category

Qualitative examples

```
Attempts to normalize a vector. If this fails, the <MASK>
public static Vector4f normalizeSafe(Vector4f v){
    try{
        return v.normalize();
    } catch (ArithmeticException ex){
        return Vector4f.ZERO;
    }
}

MASK: method catches the exception and return a zero vector of the same
dimension involved
T5: method catches the exception and return a zero vector of the same
protected void onSizeChanged(int w, int h, int oldw, int oldh) {
    super.onSizeChanged(w, h, oldw, oldh);

    // Make sure <MASK>
    if (h != oldh) {
        recomputeScrollPosition(h, oldh, mPageMargin, mPageMargin);
    }
}

MASK: scroll position is set correctly
T5: scroll position is set correctly
5-Gram: there is capacity for at
public void write(ChannelHandlerContext ctx, Object message,
    ChannelPromise promise) {
    throws Exception {
        checkState(
            state == State.PROTECTED,
            "Cannot write ... %s", state);
    }
}

// Just <MASK>
pendingUnprotectedWrites.add(msg, promise);
}

MASK: add the message to the pending queue.
T5: add the first one and try again.
5-Gram: update the new credentials.
@Return an account limit structure that <MASK>
public java.util.List<AccountLimit> getAccountLimits() {
    if (accountLimits == null) {
        accountLimits =
            new com.amazonaws.internal.SdkInternalList<AccountLimit>();
    }
    return accountLimits;
}

MASK: contain a list of AWS CloudFormation account limits and their
values.
5-Gram: contain a list of AWS CloudFormation account limits and their
values.
T5: describes the layer's limits.
public ByteBuffer encode(OffHeapValueHolder<V> valueHolder);
ByteBuffer serialized;
-
byteBuffer.putLong(valueHolder.expirationTime());
//represent the hits on previous versions. It is kept for
//compatibility reasons with <MASK>
byteBuffer.putLong(0L);
-
return byteBuffer;
}

MASK: previously saved data
5-Gram: previously saved data
T5: write/write
```

Fig. 5: Qualitative examples.

Fig. 6 depicts the average *confidence* level for perfect (continues lines) and wrong (dashed) predictions made by the T5 (red lines) and the 5-gram (orange) model. As it can be seen, the confidence of both models is a good proxy for the quality of the predictions. This is particularly true for the T5, for which we can see as the correct predictions tend to have a confidence greater than 0.5 for all k values, while the wrong predictions have confidence approaching 0.0 when k increases. Thus, a threshold based on confidence could be used in IDE tools to avoid recommending predictions likely to be wrong.

Finally, it is worth mentioning that, as compared to the task of code comment generation (*i.e.*, generating a comment from scratch given a code as input), the performance achieved in terms of perfect prediction is substantially higher. Indeed, when looking at a similar model (T5) experimented in the context of Javadoc generation, it achieved $\sim 10\%$ perfect predictions [19]. The results in Fig. 3 show that, depending on the length of the prediction on the *Javadoc*_{task}, for code comment completion, the same model can achieve $\sim 16\text{--}55\%$ perfect predictions. Thus, it might be more feasible for such a (simpler) problem to develop tools that can already support developers in their everyday coding activities.

V. THREATS TO VALIDITY

Construct validity. While building our dataset, we made three important assumptions. First, by providing the code context as input to the T5, we assume that the comment is written by developers when the code is already implemented, which is not always the case. Still, this does not invalidate the reported results, but it means, from a practical point of view, that comment completion recommendations could be triggered by the T5 only when comments are added after the code is written. Second, when fine-tuning the T5 model, we hid from the comment under analysis the sentences following the ones in which we masked tokens. Thus, we are assuming that the comment is written linearly (*i.e.*, one sentence after the other) which, again, is not always the case. Third, in our study we use the original comment written by developers as oracle, assuming that it represents a valuable reference for the experimented model. Also this assumption, while done in many previous works [13]–[17], might be wrong.

Internal validity. An important factor that influences DL performance is hyperparameters tuning. For the pre-training phase, we used the default T5 parameters selected in the original paper [25]. For the fine-tuning, we did not change the model architecture (*e.g.*, number of layers), but we experimented with different learning rates. We are aware that a more extensive calibration would likely produce better results.

External validity. Our study involved $\sim 500\text{k}$ Java methods, thus ensuring a good generalizability at least for Java code. Our results cannot be generalized to other languages.

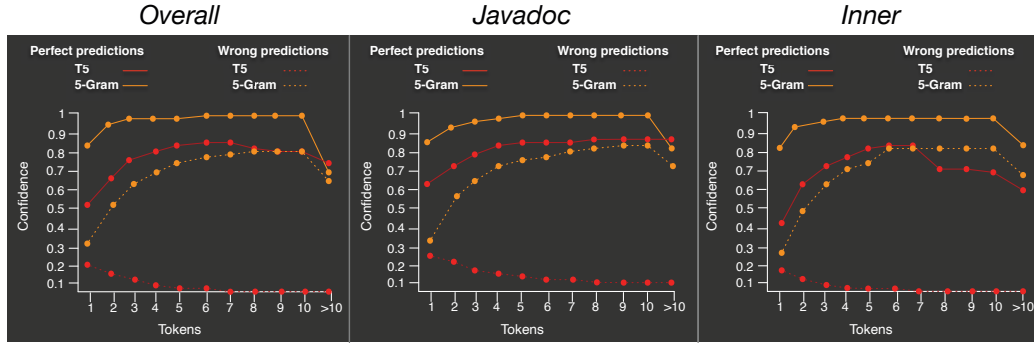


Fig. 6: Confidence level in relation to the length of the predicted tokens

VI. RELATED WORK

Our work is related to approaches (semi)automatically documenting source code. They can be roughly classified into approaches relying on summarization techniques (e.g., [41], [42]), mining crowd knowledge (e.g., [43], [44]) and using supervised learning techniques (e.g., [17], [45]). We briefly discuss the first two categories, while we focus more on the latter being it closer to our research.

A. Code Summarization

These techniques can generate code summaries in the form of (i) bag of words representing the main responsibilities of the code [42], [46]–[48]; (ii) distilled code, generated by hiding lines not considered fundamental for the code comprehension [49]; and (iii) natural language text, trying to describe the code functionality as humans would do [50], [51]. The first two categories are considered *extractive* approaches, since they synthesize the summary by extracting the most important elements from the original input, while the latter is representative of *abstractive* approaches, that can include in the summary information not present in the code to document.

Both extractive and abstractive techniques have been used to document code components at different granularity levels, such as method (e.g., [41], [42], [50], [52], [53]), method parameters (e.g., [54]), method usages (e.g., [55], [56]) class (e.g., [11], [49], [57]), unit tests [58], and code snippets (e.g., [59], [60]).

The abstractive approaches working at method/code snippet level are the most relevant for our research, since we aim at auto completing code comments that are (i) related to a single method or to a part of it; and (ii) written in natural language. Despite these similarities the tackled problem is different.

B. Mining Crowd Documentation

Another category of approaches relies on mining and processing online resources to aid program comprehension. These approaches have been used to recommend API usage examples [61]–[77], fragments of textual/video tutorials relevant for a coding task at hand [78]–[80], discussions on developers' communications, Q&A websites/forums relevant for a given code [68], [73], [81], and relevant pieces of information in API reference documentation [69].

Closer to our work are the studies leveraging crowd knowledge (e.g., [43], [44]), or pre-existing documentation (e.g., [82], [83]) to document source code. For instance, Rahman *et al.* [43] implemented *CodeInsight* a tool to mine insightful comments from *StackOverflow* discussing bugs or improvement tips for a code snippet at hand. Wong *et al.* [82], [84] leveraged existing code comments (within a set of input projects) to automatically generate comments for a target project with an approach called *ColCom*. In a nutshell, their approach reuses code comments extracted from input projects for type-1 and type-2 clones found in the target project. More recently, Aghajani *et al.* [83] presented *ADANA*, an approach to automatically inject code comments describing a given piece of Android code at the granularity level intended by the developer. *ADANA* reuses the descriptions of semantically similar code snippets retrieved and processed from *GitHub Gist* and *StackOverflow*.

These techniques could be used to complement and refine the recommendations generated by the models we experimented with. For example, once identified code comments in which the T5 model fails to suggest correct completions, crowdsourced approaches could be used to retrieve similar comments wrote by developers in a similar context (i.e., similar code to document) to strengthen the T5 training set on that specific scenario.

C. Machine Learning for Comment Generation & Completion

We focus on approaches exploiting machine learning to summarize and document source code, not discussing works done in related but different areas (e.g., summarizing code changes [85], suggest descriptive method names [15]).

Nazar *et al.* [86] recruited four students with development experience asking them to extract, from a corpus of 127 code snippets, the code lines needed to summarize each snippet. Then, ten developers were asked to inspect the lines selected as relevant for the summary and extract features characterizing them. A total of 21 features have been defined and used to train a SVM classifier able to achieve a precision of 82% in selecting relevant code lines for a code snippet.

Ying and Robilliard [87] presented an approach using supervised learning to summarize code examples.

The supervised algorithm exploits features extracted from the statements composing the code examples. The authors compute the accuracy of their approach by contrasting the generated summaries against those in a manually built oracle, reporting a 71% precision. The same authors also empirically investigated developers write code summaries [88].

DL models are on the rise for the automatic documentation of code. Iyer *et al.* [14] presented *CODE-NN*, a model producing natural language summaries describing C# code and SQL queries. Zheng *et al.* [89] proposed *Code Attention*, an attention module to translate code to comments. It exploits domain features of code snippets to infer their structure.

Liang *et al.* [90] used a Recursive Neural Network named Code-RNN to describe the structural information of source code and produce natural language comments. Their Code-RNN takes advantage of a novel GRU cell (Code-GRU) designed for code comments generation. The authors reported the quality of the generated comments by computing the ROUGE score [91] and showing the superiority of the proposed approach as compared to competitive techniques.

Hu *et al.* [13] use a Deep Neural Network (DNN) to automatically generate comments for a given Java method. To train the DNN, the authors mine $\sim 9k$ Java projects hosted on GitHub by collecting pairs of (method, comment), where “comment” is the first sentence of the Javadoc linked to the method. To assess the effectiveness of their technique, the authors computed the BLEU-4 score [37], showing the superiority of their approach with respect to the competitive technique presented in [14].

LeClair *et al.* [45] presented a neural model combining the AST source code structure and words from code to generate coherent summaries of Java methods. The approach, tested on 2.1M methods, showed its superiority as compared to the previous works by Hu *et al.* [13] and Iyer *et al.* [14].

Haque *et al.* [17] presented an approach aimed at documenting Java methods through an encoder-decoder architecture and representing an improvement of the work by LeClair *et al.* [45]. Their model leverages multiple information about the method to document, and in particular: (i) the source code of the method, as a flattened sequence of tokens representing the method, (ii) its AST representation, and (iii) the “file context”, meaning the code of every other method in the same file. The authors show that adding the contextual information as one of the inputs substantially improves the BLEU score obtained by deep learning techniques. Finally, in a recent work, Mastropaolo *et al.* [19] showed that a T5 Model [25] properly pre-trained and fine-tuned can achieve better performance than the technique presented in [17], generating comments “as humans would do” in $\sim 10\%$ of cases.

To summarize, even the most positive results in the literature show major limitations for automatically generating code comments. For this reason, we tackled the “simpler” comment completion problem proposed by Ciurumelea *et al.* [24].

The authors mine the top-1000 most starred GitHub python projects in November 2018. Then, they pre-process the extracted data and train three different models.

The first is a *Sequential Model* taking as input a sequence of Python docstring tokens and trained to predict the single token following that sequence. The second model takes as input a sequence of the Python docstring concatenated with structural information (*i.e.*, the method signature). Also in this case, it is trained to predict the next token following the docstring provided as input. Finally, a *Context Model (full body)* is trained for the same task, with the structural information provided to the method represented in this case by the complete method body. The authors show that the third model is the one providing the best results.

Our study stems from the work by Ciurumelea *et al.* [24]. However, it is performed in a different context, investigating the actual need for a DL-based model over a simpler n -gram model while looking at more challenging prediction tasks.

VII. CONCLUSIONS AND FUTURE WORK

We presented an empirical study comparing two different techniques, namely the T5 and the n -gram model, in the task of code comment completion (*i.e.*, autocomplete a code comment the developer started writing). The two models are different in nature, with the T5, based on deep learning, exploiting as information to support the completion of a comment C not only the first tokens typed by the developer while writing C but also a “context” representing code relevant for C . The n -gram model, instead, does only consider the $n - 1$ preceding tokens to predict the n^{th} token.

Our results showed the superiority of the T5, achieving significantly better prediction performance as compared to the n -gram model. However, the simplicity of the latter and its wider applicability (it does not require a code context as input), make the two models potentially complementary in the implementation of a code comment completion tool.

Our future work aim at perform additional investigations needed to better understand the strengths and weaknesses of the two models. This include studying: (i) the role played by the “code context” in the performance of the T5 (*i.e.*, what are the performance of the T5 when providing as input *only* the comment tokens “typed” by the developer?); (ii) the performance of the experimented models for different types of code comments (*e.g.*, Pascarella *et al.* [92] presented a taxonomy of code comment types and an approach to automatically classify them); (iii) how the T5 performance varies when using larger models [25]; and (iv) whether the performance of the models improves when training and testing them on code comments from a specific domain (*e.g.*, only comments extracted from Android apps).

Finally, we are working on the integration of the two models in an IDE plugin to assess their usefulness in studies with developers.

ACKNOWLEDGMENT

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 851720).

REFERENCES

- [1] X. Xia, L. Bao, D. Lo, Z. Xing, A. E. Hassan, and S. Li, "Measuring program comprehension: A large-scale field study with professionals," *IEEE Transactions on Software Engineering*, pp. 951–976, 2018.
- [2] S. C. B. de Souza, N. Anquetil, and K. M. de Oliveira, "A study of the documentation essential to software maintenance," in *International Conference on Design of Communication*, 2005, pp. 68–75.
- [3] D. Spinellis, "Code documentation," *IEEE Software*, vol. 27, no. 4, pp. 18–19, July 2010.
- [4] B. Fluri, M. Wursch, and H. C. Gall, "Do code and comments co-evolve? on the relation between source code and comment changes," in *Proceedings of the 14th Working Conference on Reverse Engineering*, 2007, pp. 70–79.
- [5] B. Fluri, M. Wursch, E. Giger, and H. C. Gall, "Analyzing the co-evolution of comments and source code," *Software Quality Journal*, vol. 17, no. 4, pp. 367–394, 2009.
- [6] M. Linares-Vásquez, B. Li, C. Vendome, and D. Poshyanyk, "How do developers document database usages in source code?" in *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, Nov 2015, pp. 36–41.
- [7] F. Wen, C. Nagy, G. Bavota, and M. Lanza, "A large-scale empirical study on code-comment inconsistencies," in *Proceedings of the 27th International Conference on Program Comprehension, ICPC 2019*, pp. 53–64.
- [8] P. Rodeghero, S. Jiang, A. Armaly, and C. McMillan, "Detecting user story information in developer-client conversations to generate extractive summaries," ser. ICSE 2017, 2017, pp. 49–59.
- [9] S. Haiduc, J. Aponte, L. Moreno, and A. Marcus, "On the use of automated text summarization techniques for summarizing source code," in *2010 17th Working Conference on Reverse Engineering*, 2010, pp. 35–44.
- [10] G. Sridhara, L. Pollock, and K. Vijay-Shanker, "Automatically detecting and describing high level actions within methods," in *2011 33rd International Conference on Software Engineering (ICSE)*, 2011, pp. 101–110.
- [11] L. Moreno, J. Aponte, G. Sridhara, A. Marcus, L. Pollock, and K. Vijay-Shanker, "Automatic generation of natural language summaries for java classes," in *ICPC 2013*, pp. 23–32.
- [12] P. W. McBurney and C. McMillan, "Automatic source code summarization of context for java methods," *IEEE Transactions on Software Engineering*, vol. 42, no. 2, pp. 103–119, 2016.
- [13] X. Hu, G. Li, X. Xia, D. Lo, and Z. Jin, "Deep code comment generation," ser. ICPC '18. Association for Computing Machinery, 2018, pp. 200–210.
- [14] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer, "Summarizing source code using a neural attention model," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp. 2073–2083.
- [15] M. Allamanis, H. Peng, and C. A. Sutton, "A convolutional attention network for extreme summarization of source code," *CoRR*, vol. abs/1602.03001, 2016.
- [16] X. Hu, G. Li, X. Xia, D. Lo, and Z. Jin, "Deep code comment generation with hybrid lexical and syntactical information," *Springer Empirical Software Engineering*, vol. 25, pp. 2179–2217, 2020.
- [17] S. Haque, A. LeClair, L. Wu, and C. McMillan, "Improved automatic summarization of subroutines via attention to file context," ser. MSR '20, 2020, pp. 300–310.
- [18] Z. Liu, X. Xia, M. Yan, and S. Li, "Automating just-in-time comment updating," in *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2020, pp. 585–597.
- [19] A. Mastropaolo, S. Scalabrino, N. Cooper, D. N. Palacio, D. Poshyanyk, R. Oliveto, and G. Bavota, "Studying the usage of text-to-text transfer transformer to support code-related tasks," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 336–347.
- [20] M. Bruch, M. Monperrus, and M. Mezini, "From examples to improve code completion systems," ser. ESEC/FSE 2009, 2009, pp. 213–222.
- [21] V. Raychev, M. Vechev, and E. Yahav, "Code completion with statistical language models," in *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '14. ACM, 2014, pp. 419–428.
- [22] A. Svyatkovskiy, S. Lee, A. Hadjitofi, M. Riechert, J. Franco, and M. Allamanis, "Fast and memory-efficient neural code completion," 2020.
- [23] S. Brody, U. Alon, and E. Yahav, "Neural edit completion," *arXiv preprint arXiv:2005.13209*, 2020.
- [24] A. Ciurumelea, S. Proksch, and H. C. Gall, "Suggesting comment completions for python using neural language models," in *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, 2020, pp. 456–467.
- [25] 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," 2019.
- [26] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "Codesearchnet challenge: Evaluating the state of semantic code search," *arXiv preprint arXiv:1909.09436*, 2019.
- [27] A. Potdar and E. Shihab, "An exploratory study on self-admitted technical debt," in *2014 IEEE International Conference on Software Maintenance and Evolution*, 2014, pp. 91–100.
- [28] X. Ren, Z. Xing, X. Xia, D. Lo, X. Wang, and J. Grundy, "Neural network-based detection of self-admitted technical debt: From performance to explainability," *ACM Trans. Softw. Eng. Methodol.*, vol. 28, no. 3, p. 15, 2019.
- [29] <https://github.com/jdkato/codetype>.
- [30] <https://pypi.org/project/urlextract/>.
- [31] <https://github.com/akoumjian/datefinder>.
- [32] <https://www.crummy.com/software/BeautifulSoup/>.
- [33] T. Kudo and J. Richardson, "Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," *CoRR*, vol. abs/1808.06226, 2018.
- [34] <https://www.nltk.org>.
- [35] N. Arivazhagan, A. Bapna, O. Firat, D. Lepikhin, M. Johnson, M. Krikun, M. X. Chen, Y. Cao, G. F. Foster, C. Cherry, W. Macherey, Z. Chen, and Y. Wu, "Massively multilingual neural machine translation in the wild: Findings and challenges," *CoRR*, vol. abs/1907.05019, 2019.
- [36] "Replication package <https://github.com/antonio-mastropaolo/ICSME2021-Completion>."
- [37] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: A method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ser. ACL '02, 2002, pp. 311–318.
- [38] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals," in *Soviet physics doklady*, vol. 10, no. 8, 1966, pp. 707–710.
- [39] S. Petrov, D. Das, and R. McDonald, "A universal part-of-speech tagset," in *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*. Istanbul, Turkey: European Language Resources Association (ELRA), May 2012, pp. 2089–2096. [Online]. Available: http://www.lrec-conf.org/proceedings/lrec2012/pdf/274_Paper.pdf
- [40] Q. McNemar, "Note on the sampling error of the difference between correlated proportions or percentages," *Psychometrika*, vol. 12, no. 2, pp. 153–157, 1947.
- [41] S. Haiduc, J. Aponte, L. Moreno, and A. Marcus, "On the use of automated text summarization techniques for summarizing source code," in *2010 17th Working Conference on Reverse Engineering*, 2010, pp. 35–44.
- [42] P. Rodeghero, C. McMillan, P. W. McBurney, N. Bosch, and S. D'Mello, "Improving automated source code summarization via an eye-tracking study of programmers," in *Proceedings of the 36th International Conference on Software Engineering*, ser. ICSE 2014, 2014, pp. 390–401.
- [43] M. M. Rahman, C. K. Roy, and I. Keivanloo, "Recommending insightful comments for source code using crowdsourced knowledge," in *2015 IEEE 15th International Working Conference on Source Code Analysis and Manipulation (SCAM)*, 2015, pp. 81–90.
- [44] C. Vassallo, S. Panichella, M. Di Penta, and G. Canfora, "CODES: Mining source code descriptions from developers discussions," ser. ICPC 2014, pp. 106–109.
- [45] A. LeClair, S. Jiang, and C. McMillan, "A neural model for generating natural language summaries of program subroutines," ser. ICSE '19, 2019, pp. 795–806.
- [46] B. P. Eddy, J. A. Robinson, N. A. Kraft, and J. C. Carver, "Evaluating source code summarization techniques: Replication and expansion," in *International Conference on Program Comprehension (ICPC)*, 2013, pp. 13–22.
- [47] S. Haiduc, J. Aponte, and A. Marcus, "Supporting program comprehension with source code summarization," in *International Conference on*

- Software Engineering - Volume 2 (ICSE'10)*, vol. 2. IEEE, 2010, pp. 223–226.
- [48] S. L. Abebe and P. Tonella, “Extraction of domain concepts from the source code,” *Science of Computer Programming*, vol. 98, pp. 680–706, 2015.
 - [49] J. Fowkes, P. Chanthirasegaran, R. Ranca, M. Allamanis, M. Lapata, and C. Sutton, “Autofolding for source code summarization,” *IEEE Transactions on Software Engineering*, vol. 43, no. 12, pp. 1095–1109, 2017.
 - [50] G. Sridhara, E. Hill, D. Muppaneni, L. Pollock, and K. Vijay-Shanker, “Towards automatically generating summary comments for java methods,” ser. ASE '10, 2010, pp. 43–52.
 - [51] S. Rastkar, G. C. Murphy, and A. W. J. Bradley, “Generating natural language summaries for crosscutting source code concerns,” in *International Conference on Software Maintenance (ICSM)*, 2011, pp. 103–112.
 - [52] P. W. McBurney, C. Liu, C. McMillan, and T. Weninger, “Improving topic model source code summarization,” in *ICPC 2014*, pp. 291–294.
 - [53] N. Dragan, M. L. Collard, and J. I. Maletic, “Reverse engineering method stereotypes,” in *International Conference on Software Maintenance*, 2006, pp. 24–34.
 - [54] G. Sridhara, L. Pollock, and K. Vijay-Shanker, “Generating parameter comments and integrating with method summaries,” in *International Conference on Program Comprehension*, 2011, pp. 71–80.
 - [55] P. W. McBurney and C. McMillan, “Automatic documentation generation via source code summarization of method context,” in *International Conference on Program Comprehension*, ser. ICPC 2014, pp. 279–290.
 - [56] —, “Automatic source code summarization of context for java methods,” *IEEE Transactions on Software Engineering*, vol. 42, no. 2, pp. 103–119, 2016.
 - [57] N. Dragan, M. L. Collard, and J. I. Maletic, “Automatic identification of class stereotypes,” in *International Conference on Software Maintenance*, 2010, pp. 1–10.
 - [58] B. Li, C. Vendome, M. Linares-Vásquez, and D. Poshyvanyk, “Aiding comprehension of unit test cases and test suites with stereotype-based tagging,” ser. ICPC '18, pp. 52–63.
 - [59] G. Sridhara, L. Pollock, and K. Vijay-Shanker, “Automatically detecting and describing high level actions within methods,” in *2011 33rd International Conference on Software Engineering (ICSE)*. IEEE, 2011, pp. 101–110.
 - [60] G. Sridhara, “Automatic generation of descriptive summary comments for methods in object-oriented programs,” Ph.D. dissertation, 2012.
 - [61] J. Li, A. Sun, and Z. Xing, “Learning to answer programming questions with software documentation through social context embedding,” *Information Sciences*, vol. 448, pp. 36–52, 2018.
 - [62] R. Holmes and G. C. Murphy, “Using structural context to recommend source code examples,” in *Proceedings of the 27th International Conference on Software Engineering*, ser. ICSE '05, pp. 117–125.
 - [63] J. Stylos and B. A. Myers, “Mica: A web-search tool for finding api components and examples,” in *VLHCC 2006*, pp. 195–202.
 - [64] G. C. Murphy, R. J. Walker, and R. Holmes, “Approximate structural context matching: An approach to recommend relevant examples,” *IEEE Transactions on Software Engineering*, vol. 32, pp. 952–970, 12 2006.
 - [65] L. Moreno, G. Bavota, M. Di Penta, R. Oliveto, and A. Marcus, “How can i use this method?” ser. ICSE '15, pp. 880–890.
 - [66] S. Subramanian, L. Inozemtseva, and R. Holmes, “Live api documentation,” in *ICSE 2014*, pp. 643–652.
 - [67] S. Jiang, A. Armary, C. McMillan, Q. Zhi, and R. Metoyer, “Docio: Documenting api input/output examples,” ser. ICPC '17, pp. 364–367.
 - [68] C. Treude and M. P. Robillard, “Augmenting api documentation with insights from stack overflow,” ser. ICSE '16, pp. 392–403.
 - [69] M. P. Robillard and Y. B. Chhetri, “Recommending reference api documentation,” *Empirical Softw. Engg.*, vol. 20, no. 6, pp. 1558–1586, 2015.
 - [70] R. P. L. Buse and W. Weimer, “Synthesizing api usage examples,” in *Proceedings of the 34th International Conference on Software Engineering*, ser. ICSE '12. IEEE Press, 2012, pp. 782–792.
 - [71] S. K. Bajracharya, J. Ossher, and C. V. Lopes, “Leveraging usage similarity for effective retrieval of examples in code repositories,” in *Proceedings of the Eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering*, ser. FSE '10, pp. 157–166.
 - [72] C. Mcmillan, D. Poshyvanyk, M. Grechanik, Q. Xie, and C. Fu, “Portfolio: Searching for relevant functions and their usages in millions of lines of code,” *ACM Trans. Softw. Eng. Methodol.*, vol. 22, no. 4, pp. 37:1–37:30, 2013.
 - [73] L. Ponzanelli, G. Bavota, M. Di Penta, R. Oliveto, and M. Lanza, “Mining stackoverflow to turn the ide into a self-confident programming prompter,” in *Proceedings of the 11th Working Conference on Mining Software Repositories*, 2014, pp. 102–111.
 - [74] L. Ponzanelli, S. Scalabrino, G. Bavota, A. Mocci, R. Oliveto, M. Di Penta, and M. Lanza, “Supporting software developers with a holistic recommender system,” ser. ICSE '17, 2017, pp. 94–105.
 - [75] T. Xie and J. Pei, “MAPO: Mining API usages from open source repositories,” in *Proc. of the 2006 Int. Workshop on Mining Software Repositories*, 2006, pp. 54–57.
 - [76] H. Zhong, T. Xie, L. Zhang, J. Pei, and H. Mei, “Mapo: Mining and recommending api usage patterns,” in *ECOOP 2009—Object-Oriented Programming*. Springer, 2009, pp. 318–343.
 - [77] J. Wang, Y. Dang, H. Zhang, K. Chen, T. Xie, and D. Zhang, “Mining succinct and high-coverage API usage patterns from source code,” ser. MSR '13, 2013, pp. 319–328.
 - [78] G. Petrosyan, M. P. Robillard, and R. De Mori, “Discovering information explaining API types using text classification,” ser. ICSE '15, vol. 1, 2015, pp. 869–879.
 - [79] L. Ponzanelli, G. Bavota, A. Mocci, M. Di Penta, R. Oliveto, M. Hasan, B. Russo, S. Haiduc, and M. Lanza, “Too long; didn’t watch! extracting relevant fragments from software development video tutorials,” ser. ICSE '16, pp. 261–272.
 - [80] L. Ponzanelli, G. Bavota, A. Mocci, M. Di Penta, R. Oliveto, B. Russo, S. Haiduc, and M. Lanza, “Codetube: extracting relevant fragments from software development video tutorials,” in *International Conference on Software Engineering Companion (ICSE-C)*. IEEE, 2016, pp. 645–648.
 - [81] L. Ponzanelli, A. Bacchelli, and M. Lanza, “Leveraging crowd knowledge for software comprehension and development,” in *Proc. of the 17th European Conf. on Soft. Maint. and Reeng.*, March 2013, pp. 57–66.
 - [82] E. Wong, T. Liu, and L. Tan, “CloCom: Mining existing source code for automatic comment generation,” in *Software Analysis, Evolution and Reengineering (SANER)*, 2015, pp. 380–389.
 - [83] E. Aghajani, G. Bavota, M. Linares-Vásquez, and M. Lanza, “Automated documentation of android apps,” *IEEE Transactions on Software Engineering*, vol. 47, no. 1, pp. 204–220, 2021.
 - [84] E. Wong, J. Yang, and L. Tan, “Autocomment: Mining question and answer sites for automatic comment generation,” in *International Conference on Automated Software Engineering (ASE)*, 2013, pp. 562–567.
 - [85] S. Rastkar and G. C. Murphy, “Why did this code change?” in *International Conference on Software Engineering (ICSE)*, 2013, pp. 1193–1196.
 - [86] N. Nazar, H. Jiang, G. Gao, T. Zhang, X. Li, and Z. Ren, “Source code fragment summarization with small-scale crowdsourcing based features,” *Frontiers of Computer Science*, vol. 10, no. 3, pp. 504–517, 2016.
 - [87] A. T. T. Ying and M. P. Robillard, “Code fragment summarization,” in *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ser. ESEC/FSE 2013, pp. 655–658.
 - [88] —, “Selection and presentation practices for code example summarization,” in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, ser. FSE 2014, pp. 460–471.
 - [89] W. Zheng, H.-Y. Zhou, M. Li, and J. Wu, “Code attention: Translating code to comments by exploiting domain features,” 2017.
 - [90] Y. Liang and K. Q. Zhu, “Automatic generation of text descriptive comments for code blocks,” *CoRR*, vol. abs/1808.06880, 2018.
 - [91] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.
 - [92] L. Pascarella, M. Bruntink, and A. Bacchelli, “Classifying code comments in java software systems,” *Empir. Softw. Eng.*, vol. 24, no. 3, pp. 1499–1537, 2019.