# Deep Just-In-Time Inconsistency Detection Between Comments and Source Code

**Sheena Panthaplackel[1], Junyi Jessy Li[2], Milos Gligoric[3], Raymond J. Mooney[1]**

[1]Department of Computer Science
[2]Department of Linguistics
[3]Department of Electrical and Computer Engineering
The University of Texas at Austin
spantha@cs.utexas.edu, jessy@austin.utexas.edu, gligoric@utexas.edu, mooney@cs.utexas.edu

## Abstract

Natural language comments convey key aspects of source code such as implementation, usage, and pre- and post-conditions. Failure to update comments accordingly when the corresponding code is modified introduces inconsistencies, which is known to lead to confusion and software bugs. In this paper, we aim to detect whether a comment becomes inconsistent as a result of changes to the corresponding body of code, in order to catch potential inconsistencies *just-in-time*, i.e., before they are committed to a code base. To achieve this, we develop a deep-learning approach that learns to correlate a comment with code changes. By evaluating on a large corpus of comment/code pairs spanning various comment types, we show that our model outperforms multiple baselines by significant margins. For extrinsic evaluation, we show the usefulness of our approach by combining it with a comment update model to build a more comprehensive automatic comment maintenance system which can both detect and resolve inconsistent comments based on code changes.

## 1  Introduction

Comments serve as a critical communication medium for developers, facilitating program comprehension and code maintenance tasks (Buse and Weimer 2010; de Souza, Anquetil, and de Oliveira 2005). Code is highly-dynamic in nature, with developers constantly making changes to address bugs and feature requests. Many code changes require reciprocal updates to the accompanying comments to keep them in sync; however, this is not always done in practice (Wen et al. 2019; Fluri et al. 2009; Ratol and Robillard 2017; Jiang and Hassan 2006; Zhou et al. 2017; Tan et al. 2007). Outdated comments which inaccurately portray the code they accompany adversely affect the software development cycle by causing confusion (Wen et al. 2019; Jiang and Hassan 2006; Tan et al. 2007; Zhou et al. 2017) and misguiding developers, hence making code vulnerable to bugs (Jiang and Hassan 2006; Tan et al. 2007; Ibrahim et al. 2012). Therefore, it is desirable to have systems that can automatically detect such inconsistencies and alert developers.

Previous work has explored heuristic-based approaches for automatically detecting specific types of inconsistencies (e.g., identifier naming (Ratol and Robillard 2017), parameter constraints (Zhou et al. 2017), `null` values and exceptions (Tan et al. 2012), locking (Tan et al. 2007),



(a) Inconsistent

(b) Consistent

Figure 1: In the example from the Apache Ignite project shown in Figure 1(a), the existing comment becomes inconsistent upon changes to the corresponding method, and in the example from the Alluxio project shown in Figure 1(b), the existing comment remains consistent after code changes.

interrupts (Tan, Zhou, and Padioleau 2011)). Some have also addressed the notion of coherence between comments and code as a text similarity problem with traditional machine learning models that leverage bag-of-words techniques (Corazza, Maggio, and Scanniello 2018; Cimasa et al. 2019). In contrast, we design an approach that generalizes across types of inconsistencies and captures deeper comment/code relationships. Furthermore, prior research has predominantly focused on detecting inconsistencies that already reside in a software project, within the code repository. We refer to this as *post hoc inconsistency detection* since it occurs potentially many commits *after* the inconsistency has been introduced.

Ideally, these inconsistencies should be detected before they ever enter the repository (e.g., during code review) since they pose a threat to the development cycle and reliability of the software until they are found. Because inconsistent comments generally arise as a consequence of developers failing to update comments immediately following code changes (Wen et al. 2019), we aim to detect whether a comment becomes inconsistent as a result of changes to the accompanying code, *before* these changes are merged into a code base. We refer to this as *just-in-time inconsistency detection*, as it allows catching potential inconsistencies right before they can materialize.

1

Detecting inconsistencies immediately following code changes allows us to utilize information from the version of the code before the changes, for which the comment is consistent. By considering how the changes affect the relationship the comment holds with the code, we can determine whether the comment remains consistent after the changes. For instance, in Figure 1(a), the comment describes the return type of nodeIds() as an array. When the method is modified to return a Set instead of an array, the comment no longer describes the correct return type, making it inconsistent. Such analysis is not possible in post hoc inconsistency detection since the exact code changes that triggered inconsistency cannot be easily pinpointed, making it difficult to align the comment with relevant parts of the code.

Moreover, due to challenges in crafting data extraction rules (Tan et al. 2007; Tan, Zhou, and Padioleau 2011) and annotating substantial amounts of data (Corazza, Maggio, and Scanniello 2018), prior post hoc work relies on a limited set of examples and projects. In contrast, we build a large corpus for just-in-time inconsistency detection by mining commit histories of software projects for code changes with and without corresponding comment updates.

Few approaches exploit code changes for inconsistency detection and these rely on task-specific rules (Sadu 2019), hand-engineered surface features (Liu et al. 2018; Malik et al. 2008), and bag-of-words techniques (Liu et al. 2018). Instead, we *learn* salient characteristics of these inputs through a deep-learning framework that encodes their syntactic structures. Namely, we use recurrent neural networks (RNNs) and gated graph neural networks (GGNNs) (Li et al. 2016) to learn contextualized representations of the comment and code changes and multi-head attention (Vaswani et al. 2017) to relate these representations in order to discern how the code changes affect the comment. We also study how manual features can complement our neural approach.

Furthermore, on its own, an inconsistency detection system can only flag comments that developers failed to update. Actually amending them to reflect code changes requires significant developer effort. Approaches for automatically updating comments based on code changes have been recently proposed (Panthaplackel et al. 2020b; Liu et al. 2020). However, they do not handle cases in which an update is not needed, such as in Figure 1(b). While the type of the key argument is modified, its purpose is unchanged (i.e., it still represents the key to be checked in PROPERTIES). Based on our user study (Panthaplackel et al. 2020b), such cases deteriorated the overall quality of the system. As a form of extrinsic evaluation, we evaluate the utility of our approach by integrating it with this comment update model, to build a more comprehensive automatic comment maintenance system that detects and resolves inconsistencies.

To summarize, our main contributions are as follows: (1) We develop a deep learning approach for just-in-time inconsistency detection that correlates a comment with changes in the corresponding body of code and which outperforms the post hoc setting as well as several baselines. (2) For training and evaluation, we construct a large corpus of comments paired with code changes in the corresponding methods, encompassing multiple types of method comments and consisting of 40,688 examples that are extracted from 1,518 open-source Java projects.[1] (3) We demonstrate the value of inconsistency detection in a comprehensive automatic comment maintenance system, and we show how our approach can support such a system.

## 2 Task

Our task is to determine whether a comment is inconsistent, or semantically out of sync with the corresponding method. Most inconsistencies result from developers making code changes without properly updating the accompanying comments. Suppose $M_{old}$ from the consistent comment/method pair $(C, M_{old})$ is modified to $M$. If $C$ is not in sync with $M$ and is not updated, it will become inconsistent once $M$ is committed. We frame this problem in two distinct settings, with the task being constant across both: determine whether $C$ is inconsistent with $M$.

- **Post hoc:** Here, only the existing version of the comment/method pair is available; the code changes that triggered the inconsistency are unknown.

- **Just-in-time:** Here, the goal is to catch inconsistencies before they are committed. Unlike the post hoc setting, $M_{old}$ is available, allowing us to analyze the changes between $M_{old}$ and $M$.

In line with most prior work in inconsistency detection (Corazza, Maggio, and Scanniello 2018; Tan et al. 2007, 2012; Khamis, Witte, and Rilling 2010), we focus on identifying inconsistencies in comments comprising API documentation for Java methods. API documentation consists of a main description and a set of tag comments (Oracle 2020). While some have considered treating the full documentation as a single comment (Corazza, Maggio, and Scanniello 2018), we choose to perform inconsistency detection at a more fine-grained level, analyzing individual comment types within this documentation. Furthermore, in contrast to previous studies tailored to a specific tag (Zhou et al. 2017; Tan et al. 2012) or specific keywords and templates (Tan et al. 2007; Tan, Zhou, and Padioleau 2011), we simultaneously consider multiple comment types with diverse characteristics. Namely, we address inconsistencies in the @return tag comment, which describes a method's return type, and the @param tag comment, which describes an argument of the method. Additionally, we examine inconsistencies in the less-structured summary comment, derived from the first sentence of the main description.

## 3 Architecture

We aim to determine whether $C$ is inconsistent by understanding its semantics and how it relates to $M$ (or changes between $M_{old}$ and $M$). We show an overview of our approach in Figure 2. First, the comment encoder, a Bi-GRU (Cho et al. 2014), encodes the sequence of tokens in $C$ (Figure 2 (1)). When learning a representation for a given

---

[1]Data and implementation are available at https://github.com/panthap2/deep-jit-inconsistency-detection.
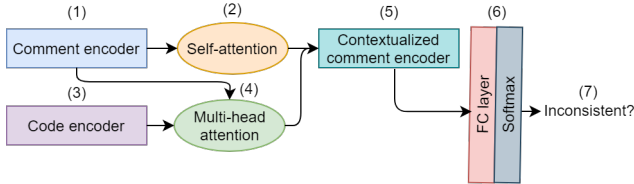
Figure 2: High-level architecture of our approach.

```
<Keep> public static boolean containsKey ( <KeepEnd>
<ReplaceOld> String <ReplaceNew> PropertyKey <ReplaceEnd>
<Keep> key ) { return PROPERTIES . containsKey ( key <KeepEnd>
<Insert> . toString ( ) <InsertEnd>
<Keep> ) ; } <KeepEnd>
```

Figure 3: Sequence-based code edit representation ($M_{edit}$) corresponding to Figure 1(b), with removed tokens in red and added tokens in green.

token, the forward and backward BiGRU passes, in principle, provide context of other tokens in $C$. However, this information can get diluted, especially when there are long-range dependencies, and the relevant context can also vary across tokens. To address this, we update these representations from the comment encoder with more context about how they relate to the other tokens through multi-head self-attention (Vaswani et al. 2017) (Figure 2 (2)). Next, we learn code representations with a code encoder (Figure 2 (3)), which can be a sequence encoder (cf. §3.1) or an abstract syntax tree (AST) encoder (cf. §3.2).

Since the essence of the task comes down to whether $C$ accurately reflects $M$, we must capture the relationship between $C$ and $M$ (or changes between $M_{old}$ and $M$). Prior work does this by computing comment/code similarity through lexical overlap rules (Ratol and Robillard 2017; Sadu 2019), which do not work well when different terms have similar meanings, and cosine similarity between vector representations, which have been found to perform poorly on their own (Liu et al. 2018; Cimasa et al. 2019). Furthermore, this notion of similarity is only appropriate for the summary comment which provides an overview of the corresponding method as a whole. More specialized comment types like @return and @param describe only specific parts of the method. Therefore, their representations may not be very similar to the representation of the full method. In contrast, we learn the relationship between comments and code by computing multi-head attention between each hidden state of the comment encoder and the hidden states of the code encoder (Figure 2 (4)).

We combine the context vectors resulting from both attention modules to form enhanced representations of the tokens in $C$, which carry context from other parts of $C$ as well as the code. These are then passed through another BiGRU encoder (Figure 2 (5)). We take the final state of this encoder to be the vector representation of the full comment, and we feed it through fully-connected and softmax layers (Figure 2 (6)). This leads to the final prediction (Figure 2 (7)).
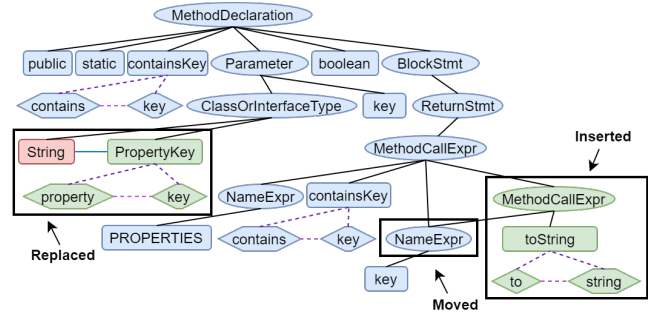


Figure 4: AST-based code edit representation ($M_{edit}$) corresponding to Figure 1(b), with removed nodes in red and added nodes in green.

## 3.1 Sequence Code Encoder

In the just-in-time setting, we represent the changes between $M_{old}$ and $M$ with an edit action sequence, $M_{edit}$. We have previously shown that explicitly defining edits in such a way outperforms having the model implicitly learn them (Panthaplackel et al. 2020b). Each action consists of an action type (Insert, Delete, Keep, ReplaceOld, ReplaceNew) that applies to a span of tokens, as shown in Figure 3. We encode $M_{edit}$ with a BiGRU. Because $M_{old}$ is unavailable in the post hoc setting, we cannot construct an edit action sequence. So, we encode the sequence of tokens in $M$.

## 3.2 AST Code Encoder

To better exploit the syntactic structure of code, we leverage its abstract syntax tree (AST). Following prior work in other tasks (Fernandes, Allamanis, and Brockschmidt 2019; Yin et al. 2019), we encode ASTs and AST edits using gated graph neural networks (GGNNs) (Li et al. 2016). For the post hoc setting, we encode $T$, an AST-based representation corresponding to $M$. In the just-in-time setting, we instead encode $T_{edit}$, an AST-based edit representation. We use GumTree (Falleri et al. 2014), to compute AST node edits between $T_{old}$ (corresponding to $M_{old}$) and $T$, identifying inserted, deleted, kept, replaced, and moved nodes. We merge the two, forming a unified representation, by consolidating identical nodes, as shown in Figure 4.

GGNN encoders for $T$ and $T_{edit}$ use *parent* (public → MethodDeclaration) and *child* (MethodDeclaration → public) edges. Like prior work (Fernandes, Allamanis, and Brockschmidt 2019), we add "subtoken nodes" for identifier leaf nodes to better handle previously unseen identifier names. To integrate these new nodes, we add *subnode* (toString → to), *supernode* (to → toString), *next subnode* (to → string), and *previous subnode* (string → to) edges. When encoding $T_{edit}$, we also include an *aligned* edge type between nodes in the two trees that correspond to an update (String and PropertyKey). Additionally, we learn *edit* embeddings for each action type. To identify how a node is edited (or not edited), we concatenate the corresponding edit embedding to its initial representation that is fed to the GGNN.

| | Train | Valid | Test | Total |
|---|---|---|---|---|
| @return | 15,950 | 1,790 | 1,840 | 19,580 |
| @param | 8,640 | 932 | 1,038 | 10,610 |
| Summary | 8,398 | 1,034 | 1,066 | 10,498 |
| Full | 32,988 | 3,756 | 3,944 | 40,688 |
| Projects | 829 | 332 | 357 | 1,518 |

Table 1: Data partitions.

## 4  Data

By detecting inconsistencies at the time of code change, we can extract automatic supervision from commit histories of open-source Java projects. Namely, we compare consecutive commits, collecting instances in which a method is modified. We extract the comment/method pairs from each version: $(C_1, M_1)$, $(C_2, M_2)$. In prior work, we isolate comment updates made based on code changes through cases in which $C_1 \neq C_2$ (Panthaplackel et al. 2020b). By assuming that the developer updated the comment because it would have otherwise become inconsistent as a result of code changes, we take $C_1$ to be inconsistent with $M_2$, consequently leading to a *positive example*, with $C=C_1$, $M_{old}=M_1$, and $M=M_2$. For *negative examples*, we additionally examine cases in which $C_1=C_2$ and assume that if the existing comment would have become inconsistent, the developer would have updated it. Following this process, we collect @return, @param, and summary comment examples. We additionally incorporate 7,239 positive @return examples from our prior work (Panthaplackel et al. 2020b) which studies @return comment updates.

While convenient for data collection, the assumptions we make do not always hold in practice. For instance, if $C_1$ is refactored without altering its meaning, we would assign a positive label because $C_1 \neq C_2$, despite it actually being consistent. Because such cases of *comment improvement* are not within the scope of our work, we adopt previously proposed heuristics (Panthaplackel et al. 2020b) to reduce the number of instances in which the comment and code changes are unrelated. The negative label is also noisy since $C_1=C_2$ when a developer fails to update comments in accordance with code changes, pointing to the problem we are addressing in this paper. We minimize such cases by limiting to popular, well-maintained projects (Jarczyk et al. 2014). For more reliable evaluation, we curate a clean sample of 300 examples (corresponding to 101 projects) from the test set, consisting of 50 positive and 50 negative examples of each comment type.

In line with prior work (Ren et al. 2019; Movshovitz-Attias and Cohen 2013), we consider a cross-project setting with no overlap between the projects from which examples are extracted in training/validation/test sets. From our data collection procedure, we obtain substantially more negative examples than positive ones, which is not surprising because many changes do not require comment updates (Wen et al. 2019). We downsample negative examples, for each partition and comment type, to construct a balanced dataset. Statistics of our final dataset are shown in Table 1.

Comments are tokenized based on space and punctuation. We parse methods into sequences using javalang (Thunes 2020). Comment and code sequences are subtokenized (e.g.,

camelCase → camel, case; snake_case → snake, case), as done in prior work (Alon et al. 2019; Fernandes, Allamanis, and Brockschmidt 2019), to capitalize on composability and better address the open vocabulary problem in learning from source code (Cvitkovic, Singh, and Anandkumar 2019). Details on data statistics, filtering, and annotation procedures are given in Appendix A.

## 5  Models

We outline baseline, post hoc, and just-in-time inconsistency detection models.

### 5.1  Baselines

**Lexical overlap:** A comment often has lexical overlap with the corresponding method. We include a rule-based just-in-time baseline, OVERLAP($C$, deleted), which classifies $C$ as inconsistent if at least one of its tokens matches a code token belonging to a `Delete` or `ReplaceOld` span in $M_{edit}$.

**Corazza, Maggio, and Scanniello (2018):** This post hoc bag-of-words approach classifies whether a comment is coherent with the method that it accompanies using an SVM with TF-IDF vectors corresponding to the comment and method. We simplify the original data pre-processing, but validate that the performance matches the reported numbers.

**CodeBERT BOW:** We develop a more sophisticated bag-of-words baseline that leverages CodeBERT (Feng et al. 2020) embeddings. These embeddings were pretrained on a large corpus of natural language/code pairs. In the post hoc setting, we consider CodeBERT BOW($C$, $M$), which computes the average embedding vectors of $C$ and $M$. These vectors are concatenated and fed through a feedforward network. In the just-in-time setting, we compute the average embedding vector of $M_{edit}$ rather than $M$, and we refer to this baseline as CodeBERT BOW($C$, $M_{edit}$).

**Liu et al. (2018):** This is a just-in-time approach for detecting whether a block/line comment becomes inconsistent upon changes to the corresponding code snippet. Their task is slightly different as block/line comments describe low-level implementation details and generally pertain to only a limited number of lines of code, relative to API comments. However, we consider it as a baseline since it is closely related. They propose a random forest classifier which leverages features which capture aspects of the code changes (e.g., whether there is a change to a `while` statement), the comment (e.g., number of tokens), and the relationship between the comment and code (e.g., cosine similarity between representations in a shared vector space). We re-implemented this approach based on specifications in the paper, as their code was not publicly available. We disregard 9 (of 64) features that are not applicable in our setting. Details about our re-implementation are given in Appendix B.

### 5.2  Our Models

**Post hoc:** We consider three models, with different ways of encoding the method. SEQ($C$, $M$) encodes $M$ with a GRU, GRAPH($C$, $T$) encodes $T$ with a GGNN, and HYBRID($C$, $M$, $T$) uses both. Multi-head attention in HYBRID($C$, $M$,

| | | Cleaned Test Sample | | | | Full Test Set | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Model | P | R | F1 | Acc | P | R | F1 | Acc |
| Baselines | OVERLAP($C$, deleted) | 77.7 | 72.0 | 74.7 | 75.7 | 74.1 | 62.8 | 68.0 | 70.4 |
| | Corazza, Maggio, and Scanniello (2018) | 65.1 | 46.0 | 53.9 | 60.7 | 63.7 | 47.8 | 54.6 | 60.3 |
| | CodeBERT BOW($C$, $M$) | 66.2 | 70.4 | 67.9 | 66.9 | 68.9 | 73.2 | 70.7 | 69.8 |
| | CodeBERT BOW($C$, $M_{edit}$) | 65.5 | 80.9 | 72.3 | 69.0 | 67.4 | 76.8 | 71.6 | 69.6 |
| | Liu et al. (2018) | 77.6 | 74.0 | 75.8 | 76.3 | 77.5 | 63.8 | 70.0 | 72.6 |
| Post hoc | SEQ($C$, $M$) | 58.9 | 68.0 | 63.0 | 60.3 | 60.6 | 73.4 | 66.3 | 62.8 |
| | GRAPH($C$, $T$) | 60.6 | 70.2 | 65.0 | 62.2 | 62.6 | 72.6 | 67.2 | 64.6 |
| | HYBRID($C$, $M$, $T$) | 53.7 | 77.3 | 63.3 | 55.2 | 56.3 | 80.8 | 66.3 | 58.9 |
| Just-In-Time | SEQ($C$, $M_{edit}$) | 83.8 | 79.3 | 81.5 | 82.0 | 80.7 | 73.8 | 77.1 | 78.0 |
| | GRAPH($C$, $T_{edit}$) | 84.7 | 78.4 | 81.4 | 82.0 | 79.8 | 74.4 | 76.9 | 77.6 |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) | 87.1 | 79.6 | 83.1 | 83.8 | 80.9 | 74.7 | 77.7 | 78.5 |
| Just-In-Time + features | SEQ($C$, $M_{edit}$) + features | 91.3 | 82.0 | 86.4 | 87.1 | 88.4 | 73.2 | 80.0 | **81.8** |
| | GRAPH($C$, $T_{edit}$) + features | 85.8 | **87.1** | 86.4 | 86.3 | 83.8 | **78.3** | 80.9 | 81.5 |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | **92.3** | 82.4 | **87.1** | **87.8** | **88.6** | 72.4 | 79.6 | 81.5 |

Table 2: Results for baselines, post hoc, and just-in-time models. Differences in F1 and Acc between just-in-time vs. baseline models, just-in-time vs. post hoc models, and just-in-time + features vs. just-in-time models are statistically significant.

$T$) is computed with the hidden states of the two encoders separately and then combined.

**Just-In-Time:** To allow fair comparison with the post hoc setting, these models are identical in structure to the models described above except that $M_{edit}$ is used instead of $M$.

**Just-In-Time + features:** Because injecting explicit knowledge can boost the performance of neural models (Chen et al. 2017; Xuan, Hieu, and Le 2018), we investigate adding comment and code features to our approach. These are computed at the token/node-level and concatenated with embeddings before being passed to encoders. Features are derived from prior work on comments and code (Panthaplackel et al. 2020a,b), including linguistic (e.g., POS tags) and lexical (e.g., comment/code overlap) features.

### 5.3 Model Training

Models are trained to minimize negative log likelihood. We use 2-layer BiGRU encoders (hidden dimension 64). GGNN encoders (hidden dimension 64) are rolled out for 8 message-passing steps, also use hidden dimension 64. We initialize comment and code embeddings, of dimension 64, with pretrained ones (Panthaplackel et al. 2020b). Edit embeddings are of dimension 8. Attention modules use 4 attention heads. We use a dropout rate of 0.6. Training ends if the validation F1 does not improve for 10 epochs.

## 6 Intrinsic Evaluation

We report common classification metrics: precision (P), recall (R), and F1 (w.r.t. the positive label) and accuracy (Acc), averaged across 3 random restarts. We also perform significance testing (Berg-Kirkpatrick, Burkett, and Klein 2012).

In Table 2, we report results for baselines, post hoc and just-in-time inconsistency detection models. In the post hoc setting, we find that our three models can achieve higher F1 scores than the bag-of-words approach proposed by Corazza, Maggio, and Scanniello (2018); however, they underperform the CodeBERT BOW($C$, $M$) baseline and significantly underperform all just-in-time models, including the simple rule-based baseline. This demonstrates the benefit

of performing inconsistency detection in the just-in-time setting, in which the code changes that trigger inconsistency are available. Additionally, by encoding the syntactic structures of the comment and code changes, our just-in-time models outperform this rule-based baseline as well as all other baselines and post hoc approaches. While the HYBRID($C$, $M_{edit}$, $T_{edit}$) model achieves slightly higher scores (on the basis of F1 and accuracy) than SEQ($C$, $M_{edit}$) and GRAPH($C$, $T_{edit}$), the differences are not statistically significant.

Our just-in-time models outperform the rule-based and feature-based baselines, without any hand-engineered rules or features. However, by incorporating surface features into our just-in-time models, we can further boost performance (by statistically significant margins). This suggests that our approach can be used in conjunction with task-specific rules (Tan et al. 2007; Tan, Zhou, and Padioleau 2011; Tan et al. 2012; Ratol and Robillard 2017) and feature sets (Liu et al. 2018) to build improved systems for specific domains.

Furthermore, in Table 3, we analyze the performance of the three just-in-time + features models with respect to individual comment types. While these models are trained on all comment types together without explicitly tailoring it in any way to handle them differently, they are still able to achieve reasonable performance across types. We provide further analysis of individual comment types and compare to comment-specific baselines in Appendix C.

## 7 Extrinsic Evaluation

We further evaluate how our approach could be used to build a comprehensive *just-in-time comment maintenance system* which first determines whether a comment, $C$, has become inconsistent upon code changes to the corresponding method ($M_{old} \rightarrow M$), and then automatically suggests an update if this is the case. To do this, we combine the inconsistency detection approach with our previously proposed comment update model (Panthaplackel et al. 2020b) which updates comments based on code changes. For training and evaluating this combined system, we have two sets of comment/method pairs from consecutive commits for each example in our cor-

| | Model | Cleaned Test Sample | | | | Full Test Set | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | P | R | F1 | Acc | P | R | F1 | Acc |
| @return | SEQ($C$, $M_{edit}$) + features | 88.5* | 72.0* | **79.4*** | **81.3*** | **87.6*** | 73.3* | 79.8* | **81.4*** |
| | GRAPH($C$, $T_{edit}$) + features | 81.2 | **77.3** | 79.1* | 79.7 | 82.2 | **79.3** | **80.6** | 80.9* |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | **88.7*** | 72.0* | **79.4*** | **81.3*** | 87.3* | 73.7* | 79.8* | **81.4*** |
| @param | SEQ($C$, $M_{edit}$) + features | 90.0 | **95.3** | 92.5 | 92.3† | 92.2 | 88.3† | 90.2 | 90.4 |
| | GRAPH($C$, $T_{edit}$) + features | **96.5** | 92.0 | **94.2** | **94.3** | **94.5** | **89.0†** | **91.7** | **91.9** |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 94.6 | 89.3 | 91.8 | 92.0† | 93.3 | 85.9 | 89.4 | 89.9 |
| Summary | SEQ($C$, $M_{edit}$) + features | **96.0** | 78.7 | 86.5§ | 87.7 | 84.7§ | 58.3 | 69.0 | **73.9§** |
| | GRAPH($C$, $T_{edit}$) + features | 80.8 | **92.0** | 86.0§ | 85.0 | 76.0 | **66.4** | **70.6** | 72.5 |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 93.7 | 86.0 | **89.5** | **90.0** | **85.0§** | 57.0 | 68.1 | 73.5§ |

Table 3: Evaluating performance with respect to different types of comments. Scores are averaged across 3 random restarts, and scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

pus. Recall from our data collection procedure that we extracted pairs of the form $(C_1, M_1)$, $(C_2, M_2)$, where $C=C_1$, $M_{old}=M_1$, and $M=M_2$. We now introduce $C_{new}=C_2$, the gold comment for $M$. If $C$ is consistent with $M$, $C=C_{new}$.

## 7.1 Evaluation Method

The GRU-based SEQ2SEQ update model encodes $C$ and a sequential representation of the code changes ($M_{edit}$). Using attention (Luong, Pham, and Manning 2015) and a pointer network (Vinyals, Fortunato, and Jaitly 2015) over learned representations of the inputs, a sequence of edit actions ($C_{edit}$) is generated, identifying how $C$ should be edited to form the updated comment ($C_{new}$). This model also employs the same linguistic and lexical features as the just-in-time + features models. The model is trained on only cases in which $C$ has to be updated and is not designed to ever copy the existing comment. We consider three different configurations for adding inconsistency detection in this model:

**Update w/ implicit detection:** We augment training of the update model with negative examples (i.e., $C$ does not need to be updated). The model implicitly does inconsistency detection by learning to copy $C$ for such cases. Inconsistency detection is evaluated based on whether it predicts $C_{new}=C$.

**Pretrained update + detection:** The update model is Panthaplackel et al. (2020b), trained on only positive examples. At test time, if the detection model classifies $C$ as inconsistent, we take the prediction of the update model. Otherwise, we copy $C$, making $C_{new}=C$. We consider three of the pretrained just-in-time detection models.

**Jointly trained update + detection:** We jointly train the inconsistency detection and update models on the full dataset (including positive and negative examples). We consider three of our just-in-time detection techniques. The update model and detection model share embeddings and the comment encoder for all three, and for the sequence-based and hybrid models, the code sequence encoder is also shared. During training, loss is computed as the sum of the update and detection components. For negative examples, we mask the loss of the update component since it does not have to learn to copy $C$. At test time, if the detection component predicts a negative label, we directly copy $C$ and otherwise take the prediction of the update model.

## 7.2 Results

We report precision, recall, F1, and accuracy for detection. As we have done previously (Panthaplackel et al. 2020b), we evaluate update through exact match (xMatch) as well as metrics used to evaluate text generation (BLEU-4 (Papineni et al. 2002) and METEOR (Banerjee and Lavie 2005)) and text editing tasks (SARI (Xu et al. 2016) and GLEU (Napoles et al. 2015)). In Table 4, we compare performances of combined inconsistency detection and update systems on the cleaned test sample. As reference points, we also provide scores for a system which never updates (i.e., always copies $C$ as $C_{new}$) and Panthaplackel et al. (2020b), which is designed to always update (and only copy $C$ if an invalid edit action sequence is generated). For completeness, we also provide results on the full dataset (which are analogous) in Appendix D.

Since our dataset is balanced, we can get 50% exact match by simply copying $C$ (i.e., never updating). In fact, this can even beat Panthaplackel et al. (2020b) on xMatch, METEOR, BLEU-4, SARI, and GLEU. This underlines the importance of first determining whether a comment needs to be updated, which can be addressed with our inconsistency detection approach. On the majority of the update metrics, both of these underperform the other three approaches (Update w/ implicit detection, Pretrained update + detection, and Jointly trained update + detection). SARI is calculated by averaging N-gram F1 scores for edit operations (add, delete, and keep). So, it is not surprising that the *Update w/ implicit detection* baseline, which learns to copy, performs fewer edits, consequently underperforming on this metric. Because Panthaplackel et al. (2020b) is designed to *always* edit, it can perform well on this metric; however, the majority of the pretrained and jointly trained systems can beat this.

The *Update w/ implicit detection* baseline, which does not include an explicit inconsistency detection component, performs relatively well with respect to the update metrics, but it performs poorly on detection metrics. Here, we use generating $C$ as the prediction for $C_{new}$ as a proxy for detecting inconsistency. It achieves high precision, but it frequently copies $C$ in cases in which it is inconsistent and should be updated, hence underperforming on recall. The pretrained and jointly trained approaches outperform this model by wide statistically significant margins across the majority of metrics, demonstrating the need for inconsistency detection.

6

| | Update Metrics | | | | | Detection Metrics | | | |
|---|---|---|---|---|---|---|---|---|---|
| | xMatch | METEOR | BLEU-4 | SARI | GLEU | P | R | F1 | Acc |
| Never Update | 50.0 | 67.4 | 72.1 | 24.9 | 68.2 | 0.0 | 0.0 | 0.0 | 50.0 |
| Panthaplackel et al. (2020b) | 25.9 | 60.0 | 68.7 | 42.0* | 67.4 | 54.0 | **95.6** | 69.0 | 57.1 |
| Update w/ implicit detection | 58.0 | 72.0 | 74.7 | 31.5 | 72.7 | **100.0** | 23.3 | 37.7 | 61.7 |
| Pretrained update + detection | | | | | | | | | |
| $\text{SEQ}(C, M_{edit})$ + features | **62.3**$^\dagger$ | 75.6* | 77.0* | 42.0* | 76.2 | 91.3* | 82.0$^\S$ | 86.4* | 87.1$^{\S\P}$ |
| $\text{GRAPH}(C, T_{edit})$ + features | 59.4 | 74.9$^\S$ | 76.6$^\dagger$ | **42.5**$^\parallel$ | 75.8*$^\dagger$ | 85.8 | 87.1 | 86.4* | 86.3$^\dagger$ |
| $\text{HYBRID}(C, M_{edit}, T_{edit})$ + features | **62.3**$^\dagger$ | 75.8$^{\dagger\parallel}$ | **77.2** | 42.3$^\dagger$ | **76.4** | 92.3 | 82.4$^\S$ | 87.1$^\dagger$ | 87.8*$^\parallel$ |
| Jointly trained update + detection | | | | | | | | | |
| $\text{SEQ}(C, M_{edit})$ + features | 61.4* | **75.9**$^\parallel$ | 76.6$^\dagger$ | 42.4$^{\dagger\parallel}$ | 75.6$^\dagger$ | 88.3$^\dagger$ | 86.2 | 87.2$^\dagger$ | 87.3$^{\S\parallel}$ |
| $\text{GRAPH}(C, T_{edit})$ + features | 60.8 | 75.1$^\S$ | 76.6$^\dagger$ | 41.8* | 75.8* | 88.3$^\dagger$ | 84.7* | 86.4* | 86.7$^{\dagger\P}$ |
| $\text{HYBRID}(C, M_{edit}, T_{edit})$ + features | 61.6* | 75.6*$^\dagger$ | 76.9* | 42.3$^\dagger$ | 75.9* | 90.9* | 84.9* | **87.8** | **88.2*** |

Table 4: Results on joint inconsistency detection and update on the cleaned test sample. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

We do not observe a significant difference between the pretrained and jointly trained systems. The pretrained models achieve slightly higher scores on most update metrics and the jointly trained models achieve slightly higher scores on the detection metrics; however, these differences are small and often statistically insignificant. Overall, we find that our approach can be useful for building a real-time comment maintenance system. Since this is not the focus of our paper but rather merely a potential use case, we leave it to future work for developing more intricate joint systems.

## 8 Related Work

**Code/Comment Inconsistencies:** Prior work analyze how inconsistencies emerge (Fluri et al. 2009; Jiang and Hassan 2006; Ibrahim et al. 2012; Fluri, Wursch, and Gall 2007) and the various types of inconsistencies (Wen et al. 2019); but, they do not propose techniques for addressing the problem.

**Post Hoc Inconsistency Detection:** Prior work propose rule-based approaches for detecting pre-existing inconsistencies in specific domains, including locks (Tan et al. 2007), interrupts (Tan, Zhou, and Padioleau 2011), `null` exceptions for method parameters (Zhou et al. 2017; Tan et al. 2012), and renamed identifiers (Ratol and Robillard 2017). The comments they consider are consequently constrained to certain templates relevant to their respective domains. We instead develop a general-purpose, machine learning approach that is not catered towards any specific types of inconsistencies or comments. Corazza, Maggio, and Scanniello (2018) and Cimasa et al. (2019) address a broader notion of coherence between comments and code through text-similarity techniques, and Khamis, Witte, and Rilling (2010) determine whether comments, specifically `@return` and `@param` comments, conform to particular format. We instead capture deeper code/comment relationships by learning their syntactic and semantic structures. Rabbi and Siddik (2020) propose a siamese network for correlating comment/code representations. In contrast, we aim to correlate comments and code through an attention mechanism.

**Just-In-Time Inconsistency Detection:** Liu et al. (2018) detect inconsistencies in a block/line comment upon changes to the corresponding code snippet using a random forest classifier with hand-engineered features. Our approach does not require such extensive feature engineering. Although their task is slightly different, we consider their approach as a baseline. Stulova et al. (2020) concurrently present a preliminary study of an approach which maps a comment to the AST nodes of the method signature (before the code change) using BOW-based similarity metrics. This mapping is used to determine whether the code changes have triggered a comment inconsistency. Our model instead leverages the full method context and also learns to map the comment directly to the code changes. Malik et al. (2008) predict whether a comment will be updated using a random forest classifier utilizing surface features that capture aspects of the method that is changed, the change itself, and ownership. They do not consider the existing comment since their focus is not inconsistency detection; instead, they aim to understand the rationale behind comment updating practices by analyzing useful features. Sadu (2019) develops at approach which locates inconsistent identifiers upon code changes through lexical matching rules. While we find such a rule-based approach (represented by our $\text{OVERLAP}(C, \text{deleted})$ baseline) to be effective, a learned model performs significantly better. Svensson (2015) builds a system to mitigate the damage of inconsistent comments by prompting developers to validate a comment upon code changes. Comments that are not validated are identified, indicating that they may be out of date and unreliable. Nie et al. (2019) present a framework for maintaining consistency between code and todo comments by performing actions described in such comments when code changes trigger the specified conditions to be satisfied.

## 9 Conclusion

We developed a deep learning approach for just-in-time inconsistency detection between code and comments by learning to relate comments and code changes. Based on evaluation on a large corpus consisting of multiple types of comments, we showed that our model substantially outperforms various baselines as well as post hoc models that do not consider code changes. We further conducted an extrinsic evaluation in which we demonstrated that our approach can be used to build a comprehensive comment maintenance system that can detect and update inconsistent comments.

7

## Acknowledgments

## Ethics Statement

Through this work, we aim to reduce time-consuming confusion and vulnerability to software bugs by keeping developers informed with up-to-date-documentation, in order to consequently help improve developers' productivity and software quality. Buggy software and incorrect API usage can result in significant malfunctions in many everyday operations. Maintaining comment/code consistency can help prevent such negative-impact events.

However, over-reliance on such a system could result in developers giving up identifying and resolving inconsistent comments themselves. By presuming that the system detects all inconsistencies and all of these are properly addressed, developers may also take the available comments for granted, without carefully analyzing their validity. Because the system may not catch all types of inconsistencies, this could potentially exacerbate rather than resolve the problem of inconsistent comments. Our system is not intended to serve as an infallible safety net for poor software engineering practices but rather as a tool that complements good ones, working alongside developers to help deliver reliable, well-documented software in a timely manner.

## References

Allamanis, M. 2019. The Adverse Effects of Code Duplication in Machine Learning Models of Code. In *SPLASH*, Onward!, 143–153.

Alon, U.; Brody, S.; Levy, O.; and Yahav, E. 2019. code2seq: Generating Sequences from Structured Representations of Code. In *International Conference on Learning Representations*.

Banerjee, S.; and Lavie, A. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 65–72.

Berg-Kirkpatrick, T.; Burkett, D.; and Klein, D. 2012. An Empirical Investigation of Statistical Significance in NLP. In *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 995–1005.

Buse, R. P. L.; and Weimer, W. R. 2010. Learning a Metric for Code Readability. *IEEE Transactions on Software Engineering* 36(4): 546–558.

Chen, R.-C.; Yulianti, E.; Sanderson, M.; and Croft, W. B. 2017. On the Benefit of Incorporating External Features in a Neural Architecture for Answer Sentence Selection. In *SIGIR Conference on Research and Development in Information Retrieval*, 1017–1020.

Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; and Bengio, Y. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Conference on Empirical Methods in Natural Language Processing*, 1724–1734.

Cimasa, A.; Corazza, A.; Coviello, C.; and Scanniello, G. 2019. Word Embeddings for Comment Coherence. In *Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, 244–251.

Corazza, A.; Maggio, V.; and Scanniello, G. 2018. Coherence of Comments and Method Implementations: A Dataset and an Empirical Investigation. *Software Quality Journal* 26(2): 751–777.

Cvitkovic, M.; Singh, B.; and Anandkumar, A. 2019. Open Vocabulary Learning on Source Code with a Graph-Structured Cache. In *International Conference on Machine Learning*, 1475–1485.

de Souza, S. C. B.; Anquetil, N.; and de Oliveira, K. M. 2005. A Study of the Documentation Essential to Software Maintenance. In *International Conference on Design of Communication: Documenting & Designing for Pervasive Information*, 68–75.

Falleri, J.-R.; Morandat, F.; Blanc, X.; Martinez, M.; and Monperrus, M. 2014. Fine-Grained and Accurate Source Code Differencing. In *International Conference on Automated Software Engineering*, 313–324.

Feng, Z.; Guo, D.; Tang, D.; Duan, N.; Feng, X.; Gong, M.; Shou, L.; Qin, B.; Liu, T.; Jiang, D.; and Zhou, M. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. *ArXiv* abs/2002.08155.

Fernandes, P.; Allamanis, M.; and Brockschmidt, M. 2019. Structured Neural Summarization. In *International Conference on Learning Representations*.

Fluri, B.; Wursch, M.; and Gall, H. C. 2007. Do Code and Comments Co-Evolve? On the Relation Between Source Code and Comment changes. In *Working Conference on Reverse Engineering*, 70–79.

Fluri, B.; Würsch, M.; Giger, E.; and Gall, H. C. 2009. Analyzing the Co-Evolution of Comments and Source Code. *Software Quality Journal* 17(4): 367–394.

Ibrahim, W. M.; Bettenburg, N.; Adams, B.; and Hassan, A. E. 2012. On the Relationship between Comment Update Practices and Software Bugs. *Journal of Systems and Software* 85(10): 2293–2304.

Jarczyk, O.; Gruszka, B.; Jaroszewicz, S.; Bukowski, L.; and Wierzbicki, A. 2014. GitHub Projects. Quality Analysis of Open-Source Software. In *International Conference on Social Informatics*, 80–94.

Jiang, Z. M.; and Hassan, A. E. 2006. Examining the Evolution of Code Comments in PostgreSQL. In *International Workshop on Mining Software Repositories*, 179–180.

Khamis, N.; Witte, R.; and Rilling, J. 2010. Automatic Quality Assessment of Source Code Comments: The JavadocMiner. In *Natural Language Processing and Information Systems*, 68–79.

Li, Y.; Tarlow, D.; Brockschmidt, M.; and Zemel, R. S. 2016. Gated Graph Sequence Neural Networks. In *International Conference on Learning Representations*.

Liu, Z.; Chen, H.; Chen, X.; Luo, X.; and Zhou, F. 2018. Automatic Detection of Outdated Comments During Code Changes. In *Annual Computer Software and Applications Conference*, 154–163.

Liu, Z.; Xia, X.; Yan, M.; and Li, S. 2020. Automating Just-In-Time Comment Updating. In *International Conference on Automated Software Engineering*, To Appear.

Luong, T.; Pham, H.; and Manning, C. D. 2015. Effective Approaches to Attention-based Neural Machine Translation. In *Conference on Empirical Methods in Natural Language Processing*, 1412–1421.

Malik, H.; Chowdhury, I.; Tsou, H.-M.; Jiang, Z. M.; and Hassan, A. E. 2008. Understanding the Rationale for Updating a Function's

Comment. In *International Conference on Software Maintenance*, 167–176.

Movshovitz-Attias, D.; and Cohen, W. W. 2013. Natural language models for predicting programming comments. In *Annual Meeting of the Association for Computational Linguistics*, 35–40.

Napoles, C.; Sakaguchi, K.; Post, M.; and Tetreault, J. 2015. Ground Truth for Grammatical Error Correction Metrics. In *Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 588–593.

Nie, P.; Rai, R.; Li, J. J.; Khurshid, S.; Mooney, R. J.; and Gligoric, M. 2019. A Framework for Writing Trigger-Action Todo Comments in Executable Format. In *Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 385–396.

Oracle. 2020. Javadoc. `https://docs.oracle.com/javase/8/docs/technotes/tools/windows/javadoc.html`.

Panthaplackel, S.; Gligoric, M.; Mooney, R. J.; and Li, J. J. 2020a. Associating Natural Language Comment and Source Code Entities. In *AAAI Conference on Artificial Intelligence*.

Panthaplackel, S.; Nie, P.; Gligoric, M.; Li, J. J.; and Mooney, R. J. 2020b. Learning to Update Natural Language Comments Based on Code Changes. In *Annual Meeting of the Association for Computational Linguistics*, 1853–1868.

Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *Annual Meeting of the Association for Computational Linguistics*, 311–318.

Rabbi, F.; and Siddik, M. S. 2020. Detecting Code Comment Inconsistency Using Siamese Recurrent Network. In *International Conference on Program Comprehension - Early Research Achievements*, 371–375.

Ratol, I. K.; and Robillard, M. P. 2017. Detecting Fragile Comments. *International Conference on Automated Software Engineering* 112–122.

Ren, X.; Xing, Z.; Xia, X.; Lo, D.; Wang, X.; and Grundy, J. 2019. Neural Network-Based Detection of Self-Admitted Technical Debt: From Performance to Explainability. *Transactions on Software Engineering and Methodology* 28: 1–45.

Sadu, A. 2019. *Automatic Detection of Outdated Comments in Open Source Java Projects*. Master's thesis, Universidad Politécnica de Madrid.

Stulova, N.; Blasi, A.; Gorla, A.; and Nierstrasz, O. 2020. Towards Detecting Inconsistent Comments in Java Source Code Automatically. In *International Working Conference on Source Code Analysis and Manipulation*, 65–69.

Svensson, A. 2015. *Reducing outdated and inconsistent code comments during software development: The comment validator program*. Master's thesis, Uppsala University.

Tan, L.; Yuan, D.; Krishna, G.; and Zhou, Y. 2007. /*iComment: Bugs or Bad Comments?*/. In *Symposium on Operating Systems Principles*, 145–158.

Tan, L.; Zhou, Y.; and Padioleau, Y. 2011. aComment: Mining Annotations from Comments and Code to Detect Interrupt Related Concurrency Bugs. In *International Conference on Software Engineering*, 11–20.

Tan, S. H.; Marinov, D.; Tan, L.; and Leavens, G. T. 2012. @tComment: Testing Javadoc Comments to Detect Comment-Code Inconsistencies. In *International Conference on Software Testing, Verification and Validation*, 260–269.

Thunes, C. 2020. Javalang. `https://pypi.org/project/javalang/`.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is All you Need. In *Conference on Neural Information Processing Systems*, 5998–6008.

Vinyals, O.; Fortunato, M.; and Jaitly, N. 2015. Pointer Networks. In *Advances in Neural Information Processing Systems*, 2692–2700.

Wen, F.; Nagy, C.; Bavota, G.; and Lanza, M. 2019. A Large-Scale Empirical Study on Code-Comment Inconsistencies. In *International Conference on Program Comprehension*, 53–64.

Xu, W.; Napoles, C.; Pavlick, E.; Chen, Q.; and Callison-Burch, C. 2016. Optimizing Statistical Machine Translation for Text Simplification. *Transactions of the Association for Computational Linguistics* 4: 401–415.

Xuan, H. N. T.; Hieu, V. C.; and Le, A.-C. 2018. Adding External Features to Convolutional Neural Network for Aspect-based Sentiment Analysis. In *Conference on Information and Computer Science*, 53–59.

Yin, P.; Neubig, G.; Allamanis, M.; Brockschmidt, M.; and Gaunt, A. L. 2019. Learning to Represent Edits. In *International Conference on Learning Representations*.

Zhou, Y.; Ruihang, G.; Taolue, C.; Zhiqiu, H.; Sebastiano, P.; and Harald, G. 2017. Analyzing APIs Documentation and Code to Detect Directive Defects. In *International Conference on Software Engineering*, 27–37.

| | Positive | Negative | Total |
|---|---|---|---|
| @return | 9,807 | 72,826 | 82,633 |
| @param | 5,507 | 19,007 | 24,514 |
| Summary | 5,904 | 69,650 | 75,554 |
| Full | 21,218 | 161,483 | 182,701 |

Table 5: Dataset sizes before downsampling

# A  Additional Data Details

We provide additional information about our procedures for filtering the dataset and curating a sample of the test set for evaluation. We also include various statistics about the examples that comprise our dataset.

## A.1  Filtering

Recall from our data collection procedure (§4) that we extracted comment/method pairs from consecutive commits, of the form $(C_1, M_1)$, $(C_2, M_2)$, where $C=C_1$, $M_{old}=M_1$, and $M=M_2$. We apply heuristics to reduce the number of cases in which there are unrelated comment and code changes. We filter out positive examples in which the differences between $C_1$ and $C_2$ entail minor cosmetic edits (e.g., reformatting, spelling corrections). Similar to prior work (Panthaplackel et al. 2020b), for @return examples, we require there to be a code change to at least one return statement or the return type of the method. We discard all @return examples (positive and negative) that do not satisfy this condition. This is because @return comments describe aspects of the return value of a method, which is typically related to the method return type and return statements. We apply the same constraint to summary comment examples, since they often describe aspects of the output (e.g., Figure 1(a)). Because @param comments generally pertain to the method's arguments, we only use examples in which an argument name or type is changed within the method. Next, we reduce the number of noisy negative examples in which a developer fails to update comments in accordance with code changes by limiting to the top 1,000 starred and forked projects, which are considered popular and well-maintained (Jarczyk et al. 2014). Furthermore, because we consider AST representations, we remove all examples consisting of a method which cannot be parsed into an AST. Additionally, we remove duplicate examples, as they have been found to negatively affect training machine learning models for source code (Allamanis 2019).

## A.2  Downsampling Negative Class

In our data collection procedure, we obtained substantially more negative examples. Because a naïve classifier trained to always predict the negative label can achieve high accuracy in such a setting, we downsample the negative class to obtain a balanced dataset. We provide the sizes of the positive and negative classes before downsampling in Table 5. Note that we also discard some examples to ensure no overlap between the projects in training, validation, and test.

## A.3  Curating Test Sample

We construct a clean test set by randomly sampling without replacement from the full test set in such a way that we at-

| | @return | @param | Summary | Full |
|---|---|---|---|---|
| $C$ | 9.7 | 8.4 | 13.3 | 10.3 |
| $M_{old}$ | 131.1 | 186.9 | 137.0 | 147.2 |
| $M$ | 131.9 | 187.7 | 135.4 | 147.3 |
| $M_{edit}$ | 179.4 | 240.9 | 186.6 | 197.3 |
| $T_{old}$ | 127.2 | 184.1 | 130.5 | 142.9 |
| $T$ | 128.1 | 184.5 | 129.5 | 143.2 |
| $T_{edit}$ | 154.3 | 213.7 | 159.1 | 171.1 |

Table 6: Statistics on the average lengths of comment and code representations.

tain a balanced sample, in terms of both labels (i.e., positive, negative) and comment type (i.e., @return, @param, summary). Rather than assigning new labels to mislabeled examples, we choose to remove such examples altogether from this sample. Note that for examples that are incorrectly labeled as negative and should actually be updated, we do not have ground truth updated comments (i.e., $C_{new}$) for evaluation. So, by removing mislabeled examples, we can use the same set of examples to evaluate the combined inconsistency detection and update system. This cleaning procedure was done by one of the authors of this paper who has 8+ years of Java experience.

We remove 11% of examples for having the incorrect label, 3% for being uncertain about the correct label due to the limited context provided in the method, and 6% from being poor examples (e.g., comments like "document me" or code changes that simply comment out the entire method). Therefore, we find 17-20% noise. For the individual comment types, the percent of noise is 6-15% for @return, 14-16% for @param, and 26-28% for summary comments. For individual labels, it is 26-28% for positive and 8-12% for negative.

## A.4  Data Statistics

In Table 6, we show the average lengths of comment and code representations for the various types of comments in our dataset. The lengths for $C$ and sequential code representations (i.e., $M_{old}$, $M$, $M_{edit}$) are computed based on the subtokenized sequences that are used by our model. Note that the $M_{edit}$ representation also includes edit keywords. We report the sizes of the AST representations ($T_{old}$, $T$, $T_{edit}$) in terms of number of nodes. This also includes the added *subnodes*.

# B  Re-Implementation of Liu et al. (2018)

From the 64 features used in Liu et al. (2018), we disregard 9 which are not compatible with our setting. Namely, since the body of code they consider are not full methods, they have separate features for the code snippet under consideration and the method to which it belongs. On the other hand, we look at only full methods, and so the separate notion of "code snippet" is irrelevant. They have two separate features *comment: ratio of comment lines to the method, comment: ratio of comment lines to the code snippet*, and we discard the latter. Next, they consider whether there is a name change for the method in which the code snippet appears; however, we look at two versions of the same method and

| | Model | Cleaned Test Sample | | | | Full Test | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **P** | **R** | **F1** | **Acc** | **P** | **R** | **F1** | **Acc** |
| @return | OVERLAP($C$, deleted) | 69.2 | 54.0 | 60.7 | 65.0$^\parallel$ | 67.6$^{\parallel\P}$ | 53.3$^\parallel$ | 59.6 | 63.9 |
| | Corazza, Maggio, and Scanniello (2018) | 73.2 | 60.0$^\P$ | 65.9 | 69.0 | 68.9$^\P$ | 61.2$^\S$ | 64.8 | 66.8 |
| | CodeBERT BOW($C$, $M$) | 84.9* | 74.7$^\S$ | 79.4*$^\P$ | 80.7* | 85.6 | 82.7 | **84.1** | **84.3** |
| | CodeBERT BOW($C$, $M_{edit}$) | 62.5 | 74.0$^\S$ | 67.7$^\parallel$ | 64.7$^\parallel$ | 66.8$^\parallel$ | 78.8*$^\dagger$ | 72.2 | 69.7 |
| | Liu et al. (2018) | 76.9 | 62.0$^\P$ | 68.6$^\parallel$ | 71.7 | 76.0 | 63.0$^\S$ | 68.9 | 71.6 |
| | Khamis, Witte, and Rilling (2010) | 52.1 | **98.0** | 68.1$^\parallel$ | 54.0 | 51.6 | **97.3** | 67.4 | 52.9 |
| | GENMATCH | 64.6 | 62.0$^\P$ | 63.3 | 64.0$^\parallel$ | 60.4 | 54.9$^\parallel$ | 57.5 | 59.5 |
| | SEQ($C$, $M_{edit}$) + features | 85.3* | 75.3$^\S$ | 79.9$^\P$ | 81.0* | 87.2$^\S$ | 75.9 | 81.2*$^\S$ | 82.4 |
| | GRAPH($C$, $T_{edit}$) + features | 87.4 | 77.3* | **82.0** | **83.0** | 84.0* | 78.0* | 80.8*$^\S$ | 81.4*$^\dagger$ |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 84.8* | 78.0$^\dagger$ | 81.2 | 82.0$^\P$ | 84.3* | 78.7$^\dagger$ | 81.3$^\S$ | 81.9* |
| Combined | SEQ($C$, $M_{edit}$) + features | 88.5$^\S$ | 72.0$^\parallel$ | 79.4*$^\P$ | 81.3*$^\P$ | **87.6**$^\S$ | 73.3$^\P$ | 79.8$^\P$ | 81.4*$^\dagger$ |
| | GRAPH($C$, $T_{edit}$) + features | 81.2 | 77.3*$^\dagger$ | 79.1* | 79.7 | 82.2 | 79.3$^\dagger$ | 80.6* | 80.9$^\dagger$ |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | **88.7**$^\S$ | 72.0$^\parallel$ | 79.4* | 81.3* | 87.3$^\S$ | 73.7$^\P$ | 79.8$^\P$ | 81.4*$^\dagger$ |

Table 7: Results for @return examples. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

| | Model | Cleaned Test Sample | | | | Full Test | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **P** | **R** | **F1** | **Acc** | **P** | **R** | **F1** | **Acc** |
| @param | OVERLAP($C$, deleted) | 85.7 | **96.0***$^\S$ | 90.6 | 90.0 | 84.0 | **93.3** | 88.4$^\P$ | 87.8$^\P$ |
| | Corazza, Maggio, and Scanniello (2018) | 74.1 | 40.0 | 51.9 | 63.0 | 59.1$^\S$ | 43.9 | 50.4 | 56.7 |
| | CodeBERT BOW($C$, $M$) | 62.8 | 57.3 | 59.9 | 61.7 | 58.9$^\S$ | 64.4 | 61.5 | 59.7 |
| | CodeBERT BOW($C$, $M_{edit}$) | 81.8 | 84.0 | 82.8 | 82.7 | 75.5 | 82.7 | 78.9$^\S$ | 77.9 |
| | Liu et al. (2018) | 90.4$^\S$ | 62.7 | 74.0 | 78.0 | 88.6$^\P$ | 72.3 | 79.6$^\S$ | 81.5 |
| | Khamis, Witte, and Rilling (2010) | **97.8*** | 90.0$^\parallel$ | 93.8 | 94.0$^\dagger$ | 87.7$^\P$ | 89.0*$^\S$ | 88.3$^\P$ | 88.2$^\P$ |
| | SEQ($C$, $M_{edit}$) + features | 95.4 | **96.0*** | 95.7*$^\dagger$ | 95.7* | 91.4 | 89.2$^\S$ | 90.3$^\dagger$ | 90.4$^{\dagger\S}$ |
| | GRAPH($C$, $T_{edit}$) + features | 97.3* | 94.0 | 95.6* | 95.7* | **94.9*** | 90.0 | **92.4** | **92.6** |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 96.6$^\dagger$ | 95.3$^\S$ | **96.0**$^\dagger$ | **96.0** | 94.3$^\dagger$ | 89.3$^\S$ | 91.7* | 91.9* |
| Combined | SEQ($C$, $M_{edit}$) + features | 90.0$^\S$ | 95.3$^\S$ | 92.5 | 92.3$^\S$ | 92.2 | 88.3* | 90.2$^\dagger$ | 90.4$^\dagger$ |
| | GRAPH($C$, $T_{edit}$) + features | 96.5$^\dagger$ | 92.0 | 94.2 | 94.3$^\dagger$ | 94.5*$^\dagger$ | 89.0*$^\S$ | 91.7* | 91.9* |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 94.6 | 89.3$^\parallel$ | 91.8 | 92.0$^\S$ | 93.3 | 85.9 | 89.4 | 89.9$^\S$ |

Table 8: Results for @param examples. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

require the name to remain the same (as part of our data collection procedure). So, we discard the feature *refactoring: rename method*. Additionally, we discard external features that are not extracted from outside the method (e.g., class-related features) as we focus on detecting inconsistencies using only the method-level context. We leave it to future work to study ways to incorporate the external context into our approach. We discard *code: changes on class attribute, code: class attribute related, refactoring: extract method, refactoring: inline method, refactoring: encapsulate field, refactoring: replace exception with test, comment: ratio of comment lines to class*.

## C   Comment-Specific Performance

Since much of the work in inconsistency detection has focused on comment-specific or task-specific settings, we analyze how our approach performs in a similar setting. Namely, we consider training and evaluating our models and baselines on only data pertaining to individual comment types. We additionally study how comment-specific training compares to combined training, in which we train on the full dataset, comprised of multiple comment types.

### C.1   @return Comments

We train the (learned) baselines introduced in Section 5.1 on only the 15,950 examples pertaining to @return comments. We additionally consider two baselines for @return comments. Khamis, Witte, and Rilling (2010) proposed a heuristic for detecting inconsistency in @return comments: the comment must begin with the correct return type of the corresponding method. We implement a baseline based on this heuristic. We also remove articles (e.g., *a*, *the*) from the beginning of the comment before applying this rule, as we found this to improve performance. Furthermore, Panthaplackel et al. (2020b) released an @return comment generation model trained on a large corpus of @return comment/method pairs. We introduce another baseline, GENMATCH, in which we use the pretrained generation model to generate an @return comment for $M_{old}$ and an @return comment for $M$. If the two comments match exactly, we consider the code change to be irrelevant to @return comments and thus the existing @return comment remains consistent. We compare these baselines with our models, trained on only @return comments. We additionally compare with our models, trained on the combined

| | Model | Cleaned Test Sample | | | | Full Test | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **P** | **R** | **F1** | **Acc** | **P** | **R** | **F1** | **Acc** |
| Summary | OVERLAP($C$, deleted) | 75.0* | 66.0 | 70.2§ | 72.0 | 71.4 | 49.7 | 58.6 | 64.9§ |
| | Corazza, Maggio, and Scanniello (2018) | 55.6 | 30.0 | 39.0 | 53.0 | 61.7 | 41.1 | 49.3 | 57.8 |
| | CodeBERT BOW($C$, $M$) | 61.6† | 78.7† | 68.8 | 64.3 | 63.7§ | **75.6*** | 68.9*†§ | 66.0§† |
| | CodeBERT BOW($C$, $M_{edit}$) | 62.1† | 80.0§ | 69.8§ | 65.3 | 64.5§ | 72.4§ | 68.0†§ | 65.9§† |
| | Liu et al. (2018) | 85.2 | 76.7 | 80.7† | 81.7 | 77.1* | 57.0†¶ | 65.5 | 70.0 |
| | SEQ($C$, $M_{edit}$) + features | 72.7 | **92.7** | 81.4† | 78.7 | 67.7 | 74.3* | 70.6‖ | 68.9 |
| | GRAPH($C$, $T_{edit}$) + features | 74.3* | 92.0* | 82.0 | 79.3 | 68.4 | 70.9 | 69.2* | 68.2 |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 70.7 | 90.0 | 79.2 | 76.3 | 64.5§ | 72.9§ | 68.4†§ | 66.3† |
| Combined | SEQ($C$, $M_{edit}$) + features | **96.0** | 78.7†§ | 86.5* | 87.7 | 84.7† | 58.3† | 69.0*† | **73.9*** |
| | GRAPH($C$, $T_{edit}$) + features | 80.8 | 92.0* | 86.0* | 85.0 | 76.0* | 66.4 | **70.6‖** | 72.5 |
| | HYBRID($C$, $M_{edit}$, $T_{edit}$) + features | 93.7 | 86.0 | **89.5** | **90.0** | **85.0†** | 57.0¶ | 68.1§ | 73.5* |

Table 9: Results for summary comment examples. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

training set, as done in the main paper.

In Table 7, we report results on the 100 @return examples in the cleaned test set as well as the 1,840 @return examples in the full test set. While the CodeBERT BOW($C$, $M_{edit}$) baseline performs quite well here, our approach can outperform baselines (w.r.t. F1 and Acc) on the cleaned test sample, when trained on only @return comments. We find that training on the combined dataset slightly deteriorates performance of our models. This is not surprising as in combined training, models must learn to generalize across comment types, not just @return comments. Nonetheless, the difference in performance between training on the comment-specific and combined sets are relatively small.

## C.2 @param Comments

For @param comments, we consider another baseline designed to follow the heuristic proposed by Khamis, Witte, and Rilling (2010) for this comment type: the comment should begin with the name of the parameter being documented. We remove articles from the beginning of the comment and consider whether the first term is one of the arguments of the method. If this is not the case, we classify it as inconsistent. We consider the comment-specific and combined settings, as we do for @return comments.

In Table 8, we report results on the 100 @param examples in the cleaned test set as well as the 1,038 @param examples in the full test set. We find that rule-based baselines can perform very well for @param comments, especially the Khamis, Witte, and Rilling (2010) baseline. This suggests that most @param comments conform to the format they suggested. Nonetheless, our models are able to *learn* this without explicitly specifying this format and can even achieve higher performance (by statistically significant margins) when trained on only @param comments. The combined setting slightly deteriorates performance of our models; however, the GRAPH($C$, $T_{edit}$) + features model can still perform slightly better than Khamis, Witte, and Rilling (2010) w.r.t. F1 on the cleaned test sample.

## C.3 Summary Comments

Summary comments (e.g., Figure 5) do not have a well-defined structure, and thus we do not have a format-based

```
/**Looks up a field with a given name and if found returns its ordinal.*/

public static int lookupField(
    final RelDataType rowType, String columnName){
    final RelDataTypeField [] fields = rowType.getFields();
    for (int i = 0; i < fields.length; i++) {
        RelDataTypeField field = fields[i];
        if (field.getName().equals(columnName)) {
            return i;
        }
    }
    return -1;
}
```

```
public static RelDataTypeField lookupField(
    final RelDataType rowType, String columnName){
    final RelDataTypeField [] fields = rowType.getFields();
    for (int i = 0; i < fields.length; i++) {
        RelDataTypeField field = fields[i];
        if (field.getName().equals(columnName)) {
            return field;
        }
    }
    return null;
}
```

Figure 5: A summary comment from the Apache Calcite Avatica project becomes inconsistent upon changes to the `lookupField()` method. It should be updated to be *Looks up a field with a given name, returning null if not found.*

baseline as we did for @return and @param comments. We evaluate baselines and our models, trained on comment-specific data, as well as our model trained on the combined training set.

In Table 9, we report results on the 100 summary examples in the cleaned test set as well as the 1,066 summary examples in the full test set. While Liu et al. (2018) is a strong baseline here, we find that we can outperform all baselines in the combined training setting. Unlike the case for @return and @param comments, combined training appears to yield improved performance over comment-specific training for our models. This suggests that the models can extract valuable information from the more structured comments in the training set that pertain to specific parts of the code in order to address the less-structured summary comments.

## D Combined Detection+Update (Full)

We illustrate our approaches for combining the tasks of detection and update in Figure 6. In Table 10, we show results

| | Training | Architecture | Inference |
|---|---|---|---|
| **Update w/ Implicit Detection** | +, − → | Update — Supervision: **c_new** | $P_{inconsistency\ label}$ = output (**Update**) != **C**; $P_{comment}$ = output (**Update**) |
| **Pretrained Update + Detection** | ±, + → ; +, − → | Update — Supervision: **c_new** ; Detection — Supervision: **C != c_new** | $P_{inconsistency\ label}$ = output (**Detection**); if output (**Detection**) *is Inconsistent*: $P_{comment}$ = output (**Update**); else: $P_{comment}$ = **C** |
| **Jointly Trained Update + Detection** | +, − → | Update — Supervision: **c_new** ; Shared: Embeddings, **C** encoder, **M_edit** encoder ; Detection — Supervision: **C != c_new** | |

Figure 6: Overview of the three configurations used for the combined system which makes a prediction for whether $C$ is inconsistent ($P_{inconsistent}$) as well for the updated comment ($P_{comment}$). Modules are trained on either the full training set, including positive/inconsistent (+) and negative/consistent (-) examples, or only positive examples.

| | Update Metrics | | | | | Detection Metrics | | | |
|---|---|---|---|---|---|---|---|---|---|
| | **xMatch** | **METEOR** | **BLEU-4** | **SARI** | **GLEU** | **P** | **R** | **F1** | **Acc** |
| Never Update | 50.0 | 67.7 | 71.6 | 25.1 | 68.3 | 0.0 | 0.0 | 0.0 | 50.0 |
| Panthaplackel et al. (2020b) | 21.5 | 56.2 | 64.7 | 37.6* | 63.4 | 53.1 | 91.8 | 67.2 | 55.3 |
| Update w/ implicit detection | 56.1* | 71.3 | 73.4* | 30.2 | 71.4 | **98.5** | 18.2 | 30.8 | 59.0 |
| Pretrained update + detection | | | | | | | | | |
| $\text{SEQ}(C, M_{edit})$ + features | **57.3**§ | **72.6*** | **73.9**† | 37.8§ | **73.2**§ | 88.4† | 73.2 | 80.0† | 81.8*† |
| $\text{GRAPH}(C, T_{edit})$ + features | 55.2 | 71.8 | 73.5* | 38.0†‖ | 72.8* | 83.8 | **78.3** | 80.9* | 81.5† |
| $\text{HYBRID}(C, M_{edit}, T_{edit})$ + features | **57.3**§ | **72.6*** | **73.9**† | 37.6* | **73.2**†§ | 88.6† | 72.4 | 79.6† | 81.5† |
| Jointly trained update + detection | | | | | | | | | |
| $\text{SEQ}(C, M_{edit})$ + features | 56.5*† | 72.2†§ | 73.5* | 37.9†‖ | 72.9* | 85.7* | 76.7* | 80.9* | 81.9*† |
| $\text{GRAPH}(C, T_{edit})$ + features | 56.2* | 72.0§ | 73.6* | 37.8†§ | 73.0*† | 85.9* | 76.7* | **81.0*** | 82.0*† |
| $\text{HYBRID}(C, M_{edit}, T_{edit})$ + features | 56.8† | 72.4*† | 73.8 | **38.1**‖ | 73.1§ | 86.7 | 75.7 | 80.9* | **82.1*** |

Table 10: Results on joint inconsistency detection and update on the full test set. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

of combined detection+update systems on the full test set. The results are analogous to those presented in Section 7.2 for the cleaned test set. While the differences for the update metrics are less pronounced, the pretrained and jointly trained approaches can again outperform *Update w/ implicit detection* as well as the two reference points: Never Update and Panthaplackel et al. (2020b). The drastic differences in performance with respect to the detection metrics further demonstrate the importance of explicit inconsistency detection in a combined detection+update system. In line with our observations from the cleaned test set, we find the performances of the pretrained and jointly trained systems to be very close.

## E Implicit vs. Explicit Edits

Through $M_{edit}$ and $T_{edit}$, we are *explicitly* defining the code edits between $M_{old}$ and $M$ or between $T_{old}$ and $T$. During preliminary experiments, we also considered having models *implicitly* learn the edits. Namely, instead of providing $M_{edit}$ as the input to the sequence-based code encoder, we encode $M_{old}$ and $M$ separately. Both of these are encoded using the same GRU encoder, but multi-head attention is computed with the two sets of hidden states separately and then combined. We do the same for the graph-based approach (i.e., encode $T_{old}$ and $T$ separately rather than use $T_{edit}$) as well as the hybrid approach. Results for these approaches are shown in Table 11, where $\text{SEQ}(C, M_{old}, M)$, $\text{GRAPH}(C, T_{old}, T)$, and $\text{HYBRID}(C, M_{old}, M, T_{old}, T)$ correspond to the encoding edits implicitly for the sequence-based, graph-based, and hybrid approaches implicitly. We find that implicitly encoding edits leads to performance that is similar (or even worse in some cases) than the post hoc setting. To truly take advantage of the just-in-time setting, we find it necessary to encode edits explicitly, which can boost performance by wide, statistically significant margins.

| | Model | Cleaned Test Sample | | | | Full Test | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **P** | **R** | **F1** | **Acc** | **P** | **R** | **F1** | **Acc** |
| Post hoc | SEQ$(C, M)$ | 58.9 | 68.0 | 63.0 | 60.3 | 60.6 | 73.4 | 66.3 | 62.8 |
| | GRAPH$(C, T)$ | 60.6 | 70.2 | 65.0 | 62.2 | 62.6 | 72.6 | 67.2 | 64.6 |
| | HYBRID$(C, M, T)$ | 53.7 | 77.3 | 63.3 | 55.2 | 56.3 | **80.8** | 66.3 | 58.9 |
| Just-In-Time (implicit) | SEQ$(C, M_{old}, M)$ | 57.8 | 67.1 | 61.6 | 58.3 | 61.5 | 74.0 | 66.9 | 63.4 |
| | GRAPH$(C, T_{old}, T)$ | 58.5 | 67.3 | 62.6 | 59.9 | 61.3 | 71.8 | 66.1 | 63.3 |
| | HYBRID$(C, M_{old}, M, T_{old}, T)$ | 58.8 | 62.9 | 59.9 | 58.7 | 62.9 | 69.3 | 65.1 | 63.1 |
| Just-In-Time | SEQ$(C, M_{edit})$ | 83.8 | 79.3 | 81.5 | 82.0 | 80.7 | 73.8 | 77.1 | 78.0 |
| | GRAPH$(C, T_{edit})$ | 84.7 | 78.4 | 81.4 | 82.0 | 79.8 | 74.4 | 76.9 | 77.6 |
| | HYBRID$(C, M_{edit}, T_{edit})$ | **87.1** | **79.6** | **83.1** | **83.8** | **80.9** | 74.7 | **77.7** | **78.5** |

Table 11: Analyzing implicit code edit representations. Differences in F1 and Acc between just-in-time (explicit) models vs. post hoc models and just-in-time (explicit) vs. just-in-time (implicit) are statistically significant.

| | Update Metrics | | | | | Detection Metrics | | | |
|---|---|---|---|---|---|---|---|---|---|
| | **xMatch** | **METEOR** | **BLEU-4** | **SARI** | **GLEU** | **P** | **R** | **F1** | **Acc** |
| Never Update | 50.0 | 67.4 | 72.1$^\P$ | 24.9 | 68.2 | 0.0 | 0.0 | 0.0 | 50.0 |
| Panthaplackel et al. (2020b) | 25.9 | 60.0 | 68.7 | **42.0** | 67.4$^\S$ | 54.0 | **95.6** | **69.0** | 57.1$^\S$ |
| Update w/ implicit detection | | | | | | | | | |
| $\quad(C, M) \to C_{new}$ | 49.2 | 67.0 | 71.5* | 27.9 | 68.4*$^\dagger$ | 77.4 | 20.0* | 31.4* | 57.3$^\S$ |
| $\quad(C, M_{old}, M) \to C_{new}$ | 48.0 | 66.4 | 71.4* | 25.4* | 67.6$^\S$ | 68.4$^\P$ | 8.7 | 15.4 | 52.3 |
| $\quad(C, M_{edit}) \to C_{new}$ | 50.9* | 68.1 | 72.1$^{\dagger\S\P}$ | 28.5 | 69.2 | 80.1 | 20.2* | 32.2* | 57.6$^\S$ |
| $\quad(C, M) \to C_{edit} \Rightarrow C_{new}$ | 50.6* | 67.7 | 72.3$^\dagger$ | 25.4* | 68.6* | **100.0**$^\S$ | 1.8$^\dagger$ | 3.5$^\dagger$ | 50.9$^\P$ |
| $\quad(C, M_{old}\ M) \to C_{edit} \Rightarrow C_{new}$ | 50.3 | 67.5 | 72.2$^\S$ | 25.3 | 68.4$^\dagger$ | 66.7$^\P$ | 1.6$^\dagger$ | 3.0$^\dagger$ | 50.8$^\P$ |
| $\quad(C, M_{edit}) \to C_{edit} \Rightarrow C_{new}$ | 52.1 | 68.6 | 72.9 | 26.9 | 69.6 | **100.0**$^\S$ | 7.1 | 13.2 | 53.6 |
| $\quad$+features | **58.0** | **72.0** | **74.7** | 31.5 | **72.7** | **100.0**$^\S$ | 23.3 | 37.7 | **61.7** |

Table 12: Results for various configurations of Update w/ implicit detection on the cleaned test sample. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

| | Update Metrics | | | | | Detection Metrics | | | |
|---|---|---|---|---|---|---|---|---|---|
| | **xMatch** | **METEOR** | **BLEU-4** | **SARI** | **GLEU** | **P** | **R** | **F1** | **Acc** |
| Never Update | 50.0$^\dagger$ | 67.7$^\S$ | 71.6$^\dagger$ | 25.1 | 68.3 | 0.0 | 0.0 | 0.0 | 50.0 |
| Panthaplackel et al. (2020b) | 21.5 | 56.2 | 64.7 | **37.6** | 63.4 | 53.1 | 91.8 | **67.2** | 55.3*$^{\dagger\S}$ |
| Update w/ implicit detection | | | | | | | | | |
| $\quad(C, M) \to C_{new}$ | 48.7 | 66.9 | 70.7* | 27.1* | 67.9 | 73.8 | 16.9* | 27.1 | 55.6* |
| $\quad(C, M_{old}, M) \to C_{new}$ | 47.9 | 66.4 | 70.6* | 25.6$^\dagger$ | 67.5 | 65.2 | 8.7 | 15.3$^\dagger$ | 52.0 |
| $\quad(C, M_{edit}) \to C_{new}$ | 50.0*$^\dagger$ | 67.7*$^{\dagger\S}$ | 71.2 | 27.9 | 68.6*$^\dagger$ | 78.7 | **18.7**$^\dagger$ | 30.1* | 56.8$^\S$ |
| $\quad(C, M) \to C_{edit} \Rightarrow C_{new}$ | 50.2* | 67.9* | 71.7 | 25.6$^\dagger$ | 68.5* | **100.0**$^\S$ | 2.3 | 4.6 | 51.2 |
| $\quad(C, M_{old}\ M) \to C_{edit} \Rightarrow C_{new}$ | 50.0$^\dagger$ | 67.8$^\dagger$ | 71.6$^\dagger$ | 25.2 | 68.4$^\dagger$ | 93.3*$^\S$ | 0.8 | 1.5 | 50.4 |
| $\quad(C, M_{edit}) \to C_{edit} \Rightarrow C_{new}$ | 52.0 | 68.9 | 72.2 | 27.0* | 69.4 | 99.6$^\S$ | 7.4 | 13.7$^\dagger$ | 53.7$^\dagger$ |
| $\quad$+features | **56.1** | **71.3** | **73.4** | 30.2 | **71.4** | 98.5* | 18.2*$^\dagger$ | 30.8* | **59.0** |

Table 13: Results for various configurations of Update w/ implicit detection on the full test set. Scores for which the difference in performance is *not* statistically significant are shown with identical symbols.

## F   Update w/ Implicit Detection Configurations

The update components of our combined detection+update systems are based on the architecture proposed by Panthaplackel et al. (2020b) for automatically updating comments based on code changes. As mentioned in Section 7.1, their approach entails encoding $C$ and a sequential code edit representation ($M_{edit}$), and then using attention and a pointer network over these learned representations to decode a sequence of comment edit actions ($C_{edit}$). The edit action sequence is finally parsed into an actual comment ($C_{new}$) as a post-processing step. This approach was initially designed to handle only cases in which a comment has to be updated. Our *Update w/ implicit detection* baseline model applies this approach on both positive (i.e., inconsistent comments that should be updated) and negative (i.e., consistent comments that should not be updated) examples. Since their approach was not designed to support negative examples, we evaluate whether other input/output configurations of their architecture would be better suited for our setting. Because the features proposed in Panthaplackel et al. (2020b) were tailored towards the specific inputs and outputs of their architecture, we disable features for these various configurations. However, we do evaluate how these configurations compare to the full model, which includes features.

We present results for the cleaned test sample and full test set in Tables 12 and 13. Our notation for input/output configurations is as follows: *(inputs) → output*. If the model

|  |  |
|---|---|
| (a) Example from OpenAPI Generator | (b) Example from OWASP ZAP |

Figure 7: Examples in which inconsistencies emerged as a result of developers failing to update comments upon code changes. Predictions of the combined, pre-trained detection+update approach are shown.

generates $C_{edit}$ which is then parsed into $C_{new}$, we use the following notation on the output side: $C_{edit} \Rightarrow C_{new}$. Note that $(C, M_{edit}) \rightarrow C_{edit} \Rightarrow C_{new}$ + *features* corresponds directly to training Panthaplackel et al. (2020b)'s model on the full training set (positive and negative examples), i.e., the *Update w/ implicit detection* model used in the main paper.

Similar to our findings from Appendix E, we observe that explicitly encoding code edits ($M_{edit}$) significantly boosts performance across most metrics, independent of the output configuration. Training the model to generate $C_{edit}$ appears to yield improved performance across most update metrics, but this appears to deteriorate performance w.r.t. the detection metrics. However, none of these configurations that rely on implicit inconsistency detection can outperform Panthaplackel et al. (2020b)'s approach (trained only on positive examples) on detection F1 or Acc. By incorporating features, we see improvements across most metrics, but the SARI and F1 metrics are still substantially higher for Panthaplackel et al. (2020b). Furthermore, the pretrained and jointly trained models presented in tables 4 and 10 can outperform all of these configurations. This study confirms that even under various input/output configurations, without an explicit inconsistency detection component, a combined detection+update system cannot adequately identify inconsistent comments. Even in scenarios in which a possible update cannot be suggested, we argue that flagging inconsistencies and alerting developers would be a critical functionality of such a system. Therefore, implicitly performing inconsistency detection, as done by the various configurations of *Update w/ implicit detection*, is not sufficient.

## G   Detecting Existing Inconsistencies

Recall that in our data collection procedure, we assign the negative (i.e., consistent) label to examples in which the developer did not update a comment following code changes. Based on our inspection of a sample of the full, unannotated test set, we find examples that are mislabeled as negative, and our model can correctly identify some of these cases. For instance, in the example shown in Figure 7(a), the developer failed to amend the comment to indicate that the method no longer returns `enumNumber` but rather its value

or `null` if it is not set. Similarly, in Figure 7(b), the developer failed to update `ZapTextField` to `JPasswordField` in the comment when the return type of the method was modified. The inconsistency in the OWASP ZAP project was fixed after we reported the issue and the inconsistency in the OpenAPI Generator project continues to persist today (at the time of submission). [2]

Our inconsistency detection model correctly predicts the positive label for both of these cases, suggesting that it would have been able to potentially prevent these inconsistencies by alerting developers just-in-time. Furthermore, by combining our pretrained HYBRID($C, M_{edit}, T_{edit}$) + features detection model with a pretrained update model (Panthaplackel et al. 2020b), we can additionally produce suggestions for resolving these inconsistencies. Note that because the update model is trained on subtokenized comments, it is not able to produce tokens like `enumNumber`; however, simple heuristics could be used to conjoin subtokens in the final prediction. While the suggested updates are not perfect, they could serve as starting points to help guide developers in updating comments. Nonetheless, this suggests that our approach can also be applied to detect existing inconsistencies (i.e., post hoc detection) by analyzing the history of changes to a particular comment/method pair.

## H   Hyperparameters

Hyperparameters were tuned on validation data. For hidden dimension size, we considered {64, 128}. For number of attention heads, we considered {1, 4, 8}. For dropout, we considered {0.2, 0.3, 0.6}.

## I   Software and Hardware

We implemented all neural models using PyTorch, and rely on PyTorch's default initialization methods for initializing model weights. We use the *scikit-learn* library to compute evaluation metrics for inconsistency detection. All models were trained on a single GPU, either NVIDIA Titan V GPUs (12 GB) or GeForce GTX Titan Black GPUs (8 GB). The

---

[2]We have reported them as issues in their respective projects. Warning: searching for these issues may reveal authors' identities.

| | | # Epochs | Training time |
|---|---|---|---|
| | SEQ$(C, M)$ | 16.3 | 1h 35.7s |
| | GRAPH$(C, T)$ | 15.3 | 26m 16s |
| | HYBRID$(C, M, T)$ | 28.3 | 2h 28m 10.3s |
| | SEQ$(C, M_{edit})$ | 18.7 | 1h 23m 24.0s |
| Detection | GRAPH$(C, T_{edit})$ | 15.0 | 29m 5.7s |
| | HYBRID$(C, M_{edit}, T_{edit})$ | 14.7 | 1h 31m 42.3s |
| | SEQ$(C, M_{edit})$ + features | 18.7 | 2h 38m 23s |
| | GRAPH$(C, T_{edit})$ + features | 17.3 | 1h 5m 21.7s |
| | HYBRID$(C, M_{edit}, T_{edit})$ + features | 16.0 | 2h 51m 15.3s |
| | Update w/ implicit detection | 33.7 | 43m 20.7s |
| | Pretrained update + detection | | |
| | SEQ$(C, M_{edit})$ + features | 57.0 | 3h 2m 59s |
| |    Detection | 18.7 | 2h 38m 23s |
| |    Update | 38.3 | 24m 36s |
| | GRAPH$(C, T_{edit})$ + features | 55.6 | 1h 29m 57.7s |
| |    Detection | 17.3 | 1h 5m 21.7s |
| Combined Detection+Update |    Update | 38.3 | 24m 36s |
| | HYBRID$(C, M_{edit}, T_{edit})$ + features | 54.3 | 3h 15m 51.3s |
| |    Detection | 16.0 | 2h 51m 15.3s |
| |    Update | 38.3 | 24m 36s |
| | Jointly trained update + detection | | |
| | SEQ$(C, M_{edit})$ + features | 30.7 | 45m 22s |
| | GRAPH$(C, T_{edit})$ + features | 25.7 | 45m 17s |
| | HYBRID$(C, M_{edit}, T_{edit})$ + features | 27.3 | 46m 11s |

Table 14: Average number of training epochs, and training time for inconsistency detection models as well as the combined detection+update models. Note that we used two different types of GPUs for these experiments, and therefore, the times are not necessarily comparable across models. Additionally, following every epoch of training the detection-only models, we compute precision, recall, F1, and Acc (using scikit-learn) on the validation data (as this determines the training termination condition), which adds to the computation time.

average number of training epochs as well as the average training times are provided in Table 14.

## J  Statistical Significance

We use bootstrap statistical significance testing (Berg-Kirkpatrick, Burkett, and Klein 2012) with $p < 0.05$ and 10,000 samples of size 5,000 each.