# CCRep: Learning Code Change Representations via Pre-Trained Code Model and Query Back

Zhongxin Liu
*Zhejiang University*
Hangzhou, China
liu_zx@zju.edu.cn

Zhijie Tang
*Zhejiang Univeristy*
Hangzhou, China
tangzhijie@zju.edu.cn

Xin Xia*
*Huawei*
China
xin.xia@acm.org

Xiaohu Yang
*Zhejiang University*
Hangzhou, China
yangxh@zju.edu.cn

*Abstract*—Representing code changes as numeric feature vectors, i.e., code change representations, is usually an essential step to automate many software engineering tasks related to code changes, e.g., commit message generation and just-in-time defect prediction. Intuitively, the quality of code change representations is crucial for the effectiveness of automated approaches. Prior work on code changes usually designs and evaluates code change representation approaches for a specific task, and little work has investigated code change encoders that can be used and jointly trained on various tasks. To fill this gap, this work proposes a novel Code Change Representation learning approach named CCRep, which can learn to encode code changes as feature vectors for diverse downstream tasks. Specifically, CCRep regards a code change as the combination of its before-change and after-change code, leverages a pre-trained code model to obtain high-quality contextual embeddings of code, and uses a novel mechanism named query back to extract and encode the changed code fragments and make them explicitly interact with the whole code change. To evaluate CCRep and demonstrate its applicability to diverse code-change-related tasks, we apply it to three tasks: commit message generation, patch correctness assessment, and just-in-time defect prediction. Experimental results show that CCRep outperforms the state-of-the-art techniques on each task.

*Index Terms*—code change, representation learning, commit message generation, patch correctness assessment, just-in-time defect prediction

## I. INTRODUCTION

During software development, developers constantly perform code changes to implement new features, fix bugs, and maintain existing code (e.g., refactoring) [1]. A code repository can be regarded as a sequence of ordered code changes. A code change can be represented as the combination of the code versions before and after the change or as a flat text, such as *diff*.

Analyzing and understanding code changes are important for a bulk of software engineering tasks. For example, commit message generation [2] requires to summarize the content and intent of code changes, and vulnerability fix identification [3] needs to analyze vulnerability-related information in a code change. To automate these tasks, a popular and effective paradigm is first encoding code changes into feature vectors, which are expected to capture the information related to the target task, and then leveraging machine learning or information retrieval (IR) techniques for automation [2], [4]–[7]. Such feature vectors are referred to as code change representations. Intuitively, the more precise the code change representations are, the less challenging the downstream learning or retrieval process will be. In contrast, if the representations miss some key information, it would be very hard, if not impossible, to obtain good results. Therefore, for code-change-related tasks, the quality of code change representations is critical to the effectiveness of automated approaches [8].

Many researchers have investigated code change representation techniques for specific downstream tasks. Some studies converted code changes into numerical vectors based on manually crafted features, such as the sizes of code changes and the syntactic structures being changed [4], [5], [9], [10]. Another line of work leverages neural networks to learn code change representations on downstream tasks in an end-to-end manner [2], [3], [11], [12], i.e., learning-based techniques. Compared to the techniques based on manually crafted features, learning-based approaches automatically learn representations from data and have shown to be more effective in many tasks [3], [13]–[15]. However, many of them adopt task-specific architectures and are trained from scratch, which makes it non-trivial to adapt them to other tasks, especially the tasks with only small datasets. In addition, existing learning-based techniques either only focus on the changed code [3], [8], [16], separately encode the changed code and its context [14], [17], or encode the code change as a whole [2], [13], [18]. Some of them ignore the context or do not highlight the changed code. All of them lack explicit interaction between the changed code and the whole code change. These hinder existing techniques from effectively capturing information from code changes.

Only a few studies focus on general approaches for learning code change representations that can be used in diverse tasks [8], [16]. Yin et al. [16] proposed to learn distributed representations of small code edits by training an auto-encoder to reconstruct edits. Unfortunately, their approach only focuses on small code edits (i.e., a single hunk with no more than 3 changed lines) while many software engineering tasks require encoding code changes with multi hunks [2], [3]. Hoang et al. [8] proposed an approach named CC2Vec to learn code change representations that can be used to boost multiple code-change-related tasks. However, CC2Vec only considers the added and removed code lines and ignores their context. Also,

*Corresponding author.

it requires commit messages, i.e., natural language descriptions of code changes, to supervise the representation learning process. However, commit messages are not always available for code-change-related tasks [6], [14], and it is challenging to collect and identify high-quality commit messages [2], [19], [20].

Considering the lack of general code change representation approaches and to boost existing solutions to code-change-related tasks, this paper proposes a **C**ode **C**hange **Rep**resentation learning approach named **CCRep**, which acts as a general code change encoder and can be used in various downstream tasks. Compared to Yin et al.'s work [16], CCRep targets commit-level code changes, which are employed by plenty of common code-change-related tasks. Different from CC2Vec [8], our approach is jointly trained with task-specific components like classifiers on the target task, not requiring additional labels for supervision. Considering the limitations of existing task-specific learning-based techniques, first, CCRep adopts a pre-trained code model to learn the representations of the code before and after a change. The pre-trained code model can build strong code representations [21], forming a good basis for learning code change representations on diverse downstream tasks. Moreover, a novel mechanism named **query back** is proposed and used in CCRep to highlight the changed code and learn to adaptively select important information from the code change by explicitly interacting the changed code fragments with the whole code change.

Specifically, given a code change, CCRep first splits it into the before-change code and the after-change code, compares the two code versions, and records the alignment information between them. Next, a pre-trained code model is adopted to compute the contextual embeddings of the two code versions, respectively. Then, the query-back mechanism is leveraged to capture the information related to the changed code from the contextual embeddings. In detail, it locates the changed code fragments via the alignment information and extracts a feature vector from them to capture change details. This change-aware feature vector is used as a query to "retrieve" related context information from the before-change and after-change code through attention [22], namely query back, and produce the final code change representation.

To show the effectiveness and the generalization of CCRep, we evaluate CCRep on three code-change-related tasks: 1) commit message generation (CMG), 2) automated patch correctness assessment (APCA), 3) just-in-time defect prediction (JIT-DP). Experimental results show that: on CMG, CCRep improves the state-of-the-art approaches by 11.8% and 12.8% in terms of BLEU on two datasets. For APCA, CCRep improves the best baseline by 5.0% and 10.2% in terms of F1 and AUC, respectively. On JIT-DP, CCRep also outperforms the best-performing baseline by 2.1%-10.7% in terms of AUC on five projects. We also conduct ablation studies, which show that both the pre-trained code model and the query-back mechanism are helpful for each task.

The contributions of this work can be summarized as follows:

- We propose the novel query-back mechanism for encoding code changes, which can highlight the changed code and learn to adaptively select important information from a code change.
- We propose a novel code change representation approach named CCRep, which consists of a pre-trained code model and the query-back mechanism. CCRep is plug-and-play and can be used in diverse code-change-related tasks.
- We comprehensively evaluate CCRep on three downstream tasks. Experimental results show that CCRep outperforms the state-of-the-art baselines on each task.
- We provide empirical evidence of the generalizability of pre-trained code models on diverse code-changed-related tasks.
- We release our replication package[1], including our source code and used datasets, for follow-up works.

The remainder of this paper is organized as follows: Section 2 introduces the problem, pre-trained code models and the motivation of the query-back mechanism. Section 3 elaborates on our approach. Section 4 presents the procedures and results of our evaluation on three tasks. We discuss the variants and the limitations of our approach and the threats to validity in Section 5. After briefly reviewing the related work in Section 6, we conclude and point out future work in Section 7.

## II. PROBLEM AND PRELIMINARY

This section describes the problem we aim to solve, briefly introduces pre-trained code models and presents the motivation of the query-back mechanism.

### A. Problem

This work focuses on proposing a learning-based code change representation approach, or in other words a code change encoder, that can be used in various code-change-related tasks. A code change $T$ consists of the code versions before and after the change, i.e., $T^b$ and $T^a$. Both $T^b$ and $T^a$ consist of a sequence of tokens, i.e., $T^b = [t_1^b, t_2^b, ..., t_{|T^b|}^b]$ and $T^a = [t_1^a, t_2^a, ..., t_{|T^a|}^a]$, where $|x|$ refers to the length of $x$. A code change encoder can be viewed as a function $f$ which converts $T$ into a numeric vector $v$, i.e., $v = f(T)$. Because different tasks may value different properties of code changes, in this work, we expect the proposed approach to be jointly trained with task-specific components on the target task.

### B. Pre-Trained Code Model

Model pre-training is widely used in the natural language processing (NLP) community and the produced pre-trained models have shown to be effective in various NLP tasks [23]–[25]. The rationales behind the impressive effectiveness of pre-trained models include: 1) pre-trained models learn high-quality language representations from huge corpora, 2) pre-trained models provide good parameter initializations for downstream tasks [26], and 3) pre-trained models are usually

---

[1]https://github.com/ZJU-CTAG/CCRep

```
1@@ -277,8 +277,8 @@ public class EditPost extends
      Activity {
2 } else {
3
4    if (extras != null) {
5-       id = WordPress.currentBlog.getId();
6       try {
7+          id = WordPress.currentBlog.getId();
8          blog = new Blog(id, this);
9       } catch (Exception e) {
10          Toast.makeText(this,
```

**Commit Message**: Moved post id creation to try catch block to help EditPost activity recover if there's no valid currentBlog

large and can avoid overfitting on the small datasets of downstream tasks. [26]. Recently, researchers also applied model pre-training to code and released a number of pre-trained code models, such as CodeBERT [21], GPT-C [27], PLBART [28] and CodeT5 [29]. These models use Transformer-based architectures, can be used to encode and/or generate code or code-related texts, and have been shown to significantly boost code understanding and generation tasks [21], [29]. The impressive performance and generality of these models inspire us to investigate their feasibility in code change representation learning.

### C. Motivation of Query-Back Mechanism

A code change contains both the changed code and its context. Table I presents a code change with its commit message collected from the WordPress-Android project [30]. We can see from this example that: 1) The changed code is the core of a code change. For this example, by inspecting line 5 and line 7, we can know that the developer "moved post id creation". 2) The context may provide important information for understanding the code change. For instance, based on the context in Table I, we can know that the code change is related to "try catch block" and "EditPost". 3) Not all the context is useful. For this example, line 2-4, line 8 and line 10 are unrelated to the content and intent of this code change. As discussed in Section I, existing code change representation approaches either ignore the context [3], [8], [16], do not highlight the changed code [2], [13], [18], or consider all the context without adaptive information selection [14], [17]. These hinder their effectiveness and generality, and motivate us to propose the query-back mechanism to explicitly highlight the changed code and learn to adaptively capture information from the code change.

### III. APPROACH

The framework of CCRep is presented in Figure 1. CCRep takes a code change as input and generates its vector representation. It consists of three parts:
**Code Change Preprocessing**. Given a code change $T$, this part prepares the token sequences of the before-change and after-change code, i.e., $T_b$ and $T_a$, identifies and aligns the

modified code fragments in them, and stores their alignment information.
**Contextual Code Embedding**. This part leverages a pretrained code model to obtain the contextual embedding of each token in $T_b$ and $T_a$, respectively.
**Query-Back Mechanism**. This part takes as input the contextual embeddings of $T_b$ and $T_a$ and the alignment information of the changed code fragments, leveraging the query-back mechanism to produce the vector representation of $T$.

### A. Code Change Preprocessing

Given a code change $T$, we first use a text diff tool, such as Python difflib [31], to convert it as a code diff. A code diff contains one or more hunks and each hunk includes the lines that are deleted from the before-change code (the deleted lines), the lines that are added in the after-change code (the added lines), and the unchanged lines before and after them (the context). We can extract the changed code fragments at different granularity levels, e.g., changed code tokens or changed code lines, from a code diff. Then, we split the code diff as a sequence of hunks and preprocess these hunks as follows:

**Line Aligning**. As demonstrated in Table II, for each hunk, we first align the before-change and after-change code line-by-line using Python difflib and obtain multiple aligned line pairs. A newly added/deleted line is regarded to be aligned with an empty line. We label the lines in each pair with a line index $l_i$ starting from 1. Specifically, for each line pair: 1) If it refers to a line change, i.e., addition, deletion or replacement, we label the lines in the pair with $l_i$. 2) If its two lines are the same, they are both labeled with the default index 0. After processing one aligned line pair, no matter whether this line pair refers to a line change or not, $l_i$ is increased by 1. For example, in Table II, the first, second, and fourth line pairs refer to line replacement, deletion, and addition, respectively. The lines in them are labeled with 1, 2 and 4, respectively. The lines in the third/fifth line pairs are the same, so they are labeled with 0. When we finish aligning one hunk, the current line index $l_i$ is passed to the next hunk as its initial line index. After this step, every line in both the before-change and after-change code will have a line index, denoted as $l_i^b$ and $l_i^a$.

**Tokenizing**. The embedding layer in the pre-trained code model is tightly bound to the vocabulary of the model's tokenizer. Therefore, to correctly make use of the pre-trained embedding layer, we use the tokenizer provided by the pre-trained code model to tokenize each code line into a token sequence. Such tokenizer is usually based on subwords, e.g., BPE [32], and needs to build a subword vocabulary from a corpus before pre-training. It can split a text into subwords and avoid the out-of-vocabulary problem [33]. Besides, for each token, we store the index $l_i$ of the line it belongs to.

**Flattening**. To prepare the flat token sequences that can be processed by the pre-trained code model, for each of the before-change and after-change code, we independently collect its tokens from all hunks in the diff and concatenate the tokens into a sequence. The token sequences of the before-change
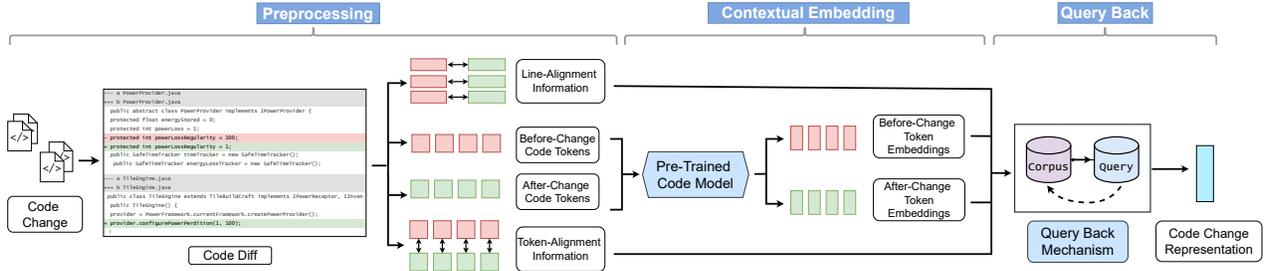
Fig. 1. The overall framework of CCRep.

TABLE II
AN ILLUSTRATIVE EXAMPLE OF LINE ALIGNING

| Line Pair Index | Before-Change Code Line | After-Change Code Line | Line Change Type | $l_i^b$ | $l_i^a$ |
|---|---|---|---|---|---|
| 1 | if (cursor != null) { | if (cursor != null && cursor.moveToFirst()) { | Replace | 1 | 1 |
| 2 | cursor.moveToFirst(); | - | Delete | 2 | N/A |
| 3 | int idx = cursor.getColumnIndex(); | int idx = cursor.getColumnIndex(); | Keep | 0 | 0 |
| 4 | - | if (idx != -1) | Add | N/A | 4 |
| 5 | result = cursor.getString(idx);} | result = cursor.getString(idx);} | Keep | 0 | 0 |

"-" refers to an empty line.

and after-change code are denoted as $T^b = [t_1^b, t_2^b, ..., t_{|T^b|}^b]$ and $T^a = [t_1^a, t_2^a, ..., t_{|T^a|}^a]$, respectively.

**Token Aligning**. We align $T^b$ and $T^a$ token by token using Python difflib to identify the changed tokens. After aligning, every token in $T^b$ ($T^a$) will get a token change flag $m_i^b$ ($m_i^a$), which is 1 for changed tokens and 0 for unchanged ones.

### B. Contextual Code Embedding

This part takes as input the token sequences of the before-change and after-change code, i.e., $T^b$ and $T^a$, aiming to independently encode them as two sequences of contextual embeddings, i.e., $H^b = [h_1^b, h_2^b, ..., h_{|T^b|}^b]$ and $H^a = [h_1^a, h_2^a, ..., h_{|T^a|}^a]$. $H^b$ and $H^a$ are expected to capture the syntactic and semantic information of the code before and after the change. CCRep leverages a pre-trained code model as the code encoder. Because pre-trained code models are shown to be able to produce high-quality code representations and can be applied to datasets of different sizes [21], [27], [29]. Specifically, our implementation of CCRep uses CodeBERT, since it is widely used and has been shown to perform well on multiple code-related tasks [3], [12], [21]. Given $T^b$ or $T^a$, CodeBERT uses a multi-layer Transformer [22] to make code tokens aggregate context information from each other and outputs their contextual embeddings $H^b$ or $H^a$. Please note that CCRep is agnostic to pre-trained code models and CodeBERT can be substituted with other pre-trained models that can be used as a code encoder.

### C. Query-Back Mechanism

This part aims to produce the final representation $v$ of the input code change $T$. To help CCRep effectively capture important information, we propose a novel mechanism named **query back** for this part. Its main idea is to encode the changed code fragments as a change-aware query $q$, use such

query to "retrieve" important information from the before-change and after-change code (the "corpus"), and produce the final code change representation $v$ based on the "retrieved" information. It can be seen that this "retrieval" process makes the changed code explicitly interact with the whole code change, so that we can adaptively extract important information from the "corpus".

Specifically, this part takes as input the contextual embeddings $H^b$ and $H^a$ produced by the pre-trained code model and the alignment information extracted during preprocessing, and outputs the code change representation $v$. Considering that different tasks may focus on the changed code fragments of different granularity, we propose three variants of the query-back mechanism, namely the token-level, the line-level, and the hybrid query-back mechanisms. Figure 2 shows their architectures.

*1) Token-Level Query-Back Mechanism:* This variant constructs the change-aware query $q^t$ based on the changed code tokens, which is beneficial if the changed code is fine-grained, e.g., only changes the name of a method or renames an identifier. As shown in the upper part of Figure 2, it works as follows:

**Changed Token Selection**. During preprocessing, each token in $T_b$ and $T_a$ is assigned a token change flag $m_i^b$ and $m_i^a$. Based on these flags, we first construct a flag sequence $M^b = [m_1^b, m_2^b, ..., m_{|T_b|}^b]$ ($M^a = [m_1^a, m_2^a, ..., m_{|T_a|}^a]$) for $T_b$ ($T_a$). Then, the contextual embeddings of the changed code tokens are picked out from $H^b$ and $H^a$ based on $M^b$ and $M^a$, and are further concatenated as a new embedding sequence
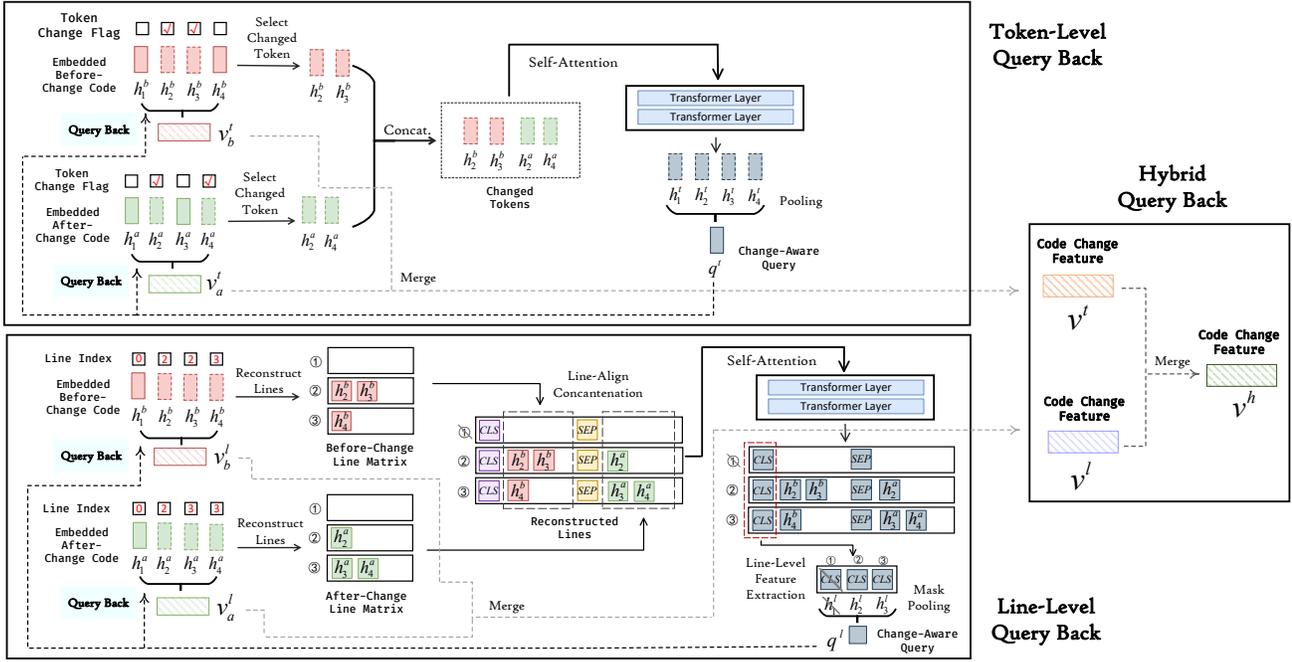
Fig. 2. The architectures of the token-level, the line-level, and the hybrid query-back mechanisms. The upper and lower parts depict the token-level and the line-level query-back mechanisms, respectively. They both take as input the contextual embeddings of the before-change and after-change code (the "corpus"), together with the alignment information, i.e., token change flags and line indices, extracted during preprocessing. Token-level query-back mechanism selects and combines changed tokens based on token change flags, constructs a query $q^t$ using a Transformer and the average pooling, and adaptively retrieves important information from the "corpus" to represent the code change as a feature vector $v^t$. Line-level query-back follows a similar procedure, but it reconstructs the structure of changed line pairs based on line indices, explicitly masks unchanged lines when constructing the query $q^l$, and outputs a feature vector $v^l$ as the code change representation. Hybrid query-back combines $v^t$ and $v^l$ to construct the code change representation $v^h$.

$H'$, as follows:

$$= [h_1^b, h_2^b, ..., h_{|T_b|}^b] \otimes [m_1^b, m_2^b, ..., m_{|T_b|}^b],$$
$$[h_1^{a'}, h_2^{a'}, ..., h_{n_a}^{a'}] = [h_1^a, h_2^a, ..., h_{|T_a|}^a] \otimes [m_1^a, m_2^a, ..., m_{|T_a|}^a],$$
$$H' = [h_1^{b'}, h_2^{b'}, ..., h_{n_b}^{b'}, h_1^{a'}, h_2^{a'}, ..., h_{n_a}^{a'}],$$
$$\text{(1)}$$

where $n_b$ and $n_a$ refer to the numbers of the changed tokens in $T_b$ and $T_a$, $\otimes$ denotes the masked-select operation which selects elements from the left operand based on the flags provided by the right operand. This step localizes the fine-grained token changes and can help the model concentrate on them during feature extraction.

**Query Construction**. In this step, we make the changed tokens aware of each other and encode them as a single vector, namely the change-aware query $q^t$. In detail, we first use a multi-layer Transformer to extract change-aware features $H^t = [h_1^t, h_2^t, ..., h_{n_b+n_a}^t]$ from $H'$. Then, average pooling is applied to $H^t$ to squash it into the query $q^t$, as follows:

$$H^t = \text{Transformer}(H') \quad \text{(2)}$$

$$q^t = \text{Pooling}(h_1^t, h_2^t, ..., h_{n_b+n_a}^t) \quad \text{(3)}$$

$q^t$ is expected to encode the information of token changes.

**Query Back**. The extracted change-aware query $q^t$ is used to "retrieve" relevant information from the before-change and after-change code (the "corpus") through attention. Since $q^t$ is also learned from the "corpus", we call this attention

query-back attention. Specifically, we adopt the multi-head attention [22] to implement the "retrieval" process, where both the keys and the values are set to the contextual embeddings $H^b$ or $H^a$ and $q^t$ is used as the query, as follows:

$$v_b^t = \text{MultiHead}(q^t, H^b, H^b) \quad \text{(4)}$$

$$v_a^t = \text{MultiHead}(q^t, H^a, H^a) \quad \text{(5)}$$

$v_b^t$ and $v_a^t$ refer to the attended feature vectors retrieved from the before-changed and after-change code, respectively.

**Merging**. Eventually, we merge $v_b^t$ and $v_a^t$ by element-wise addition to get the final code change representation $v^t$:

$$v^t = v_b^t + v_a^t. \quad \text{(6)}$$

*2) Line-Level Query-Back Mechanism:* This variant constructs the change-aware query $q^l$ based on the changed code lines, which captures line-level code change features and can be useful if one or more code lines are completely added or deleted. As shown in the lower part of Figure 2, it works as follows:

**Changed Line Selection**. In this step, based on the line index $l_i$ of each token generated during preprocessing, we first select out the tokens in the changed lines and reconstruct the line structures of the before-change and after-change code $T^b$ and $T^a$, respectively. The reconstruction process is identical for $T^b$ and $T^a$, and we leverage the example in Figure 2 to illustrate it: (1) First, we use the line indexes of all the

tokens in $T^b$ ($T^a$) to form a line index sequence $[l_1^b, l_2^b, ..., l_{|T_b|}^b]$ ($[l_1^a, l_2^a, ..., l_{|T_a|}^a]$). In our example, such sequence is $[0, 2, 2, 3]$ ($[0, 2, 3, 3]$). (2) Then, we initialize an empty matrix of shape $L \times W \times d$ for reconstructing line structures, where $L$, $W$ and $d$ refer to the maximum lines of the code (line-dimension), the maximum tokens in a code line and the dimension of a contextual embedding. In our example, $L$=3 and $W$=3. We refer to this matrix as line matrix. (3) After that, the contextual embedding of each token in $T^b$ ($T^a$) is filled into the corresponding row of the line matrix according to its line index. In our example, $h_2^b$ and $h_3^b$ are filled into the 2-nd row, while $h_4^b$ is filled into the 3-rd row. (4) Finally, we order the contextual embeddings in each row by their token indices in $T^b$ ($T^a$). In this way, line structures are reconstructed and stored in the line matrix, of which each row stores one line, as shown by the "Before-Change Line Matrix" and the "After-Change Line Matrix" in Figure 2. We refer to the process mentioned above as the scattering-reshaping operation, and briefly formulate it as:

$$Ls^b = [h_1^b, h_2^b, ..., h_{|T_b|}^b] \odot [l_1^b, l_1^b, ..., l_{|T_b|}^b], \quad Ls^b \in \mathcal{R}^{L \times W \times d}$$
$$Ls^a = [h_1^a, h_2^a, ..., h_{|T_a|}^a] \odot [l_1^a, l_1^a, ..., l_{|T_a|}^a], \quad Ls^a \in \mathcal{R}^{L \times W \times d}$$
(7)

where $\odot$ refers to the scattering-reshaping operation, $Ls^b$ and $Ls^a$ refer to the line matrices of $T^b$ and $T^a$, respectively. Note that the tokens from the unchanged lines (with line index 0) are filled into the 0-th row of the line matrix. Since we only care the changed lines, we simply drop out the 0-th row. The line matrices $Ls^b$ and $Ls^a$ are concatenated along the second dimension (the token dimension) to pair aligned lines and form longer lines, and two special tokens *CLS* and *SEP* are respectively inserted into the head and the tail of the before-change line in each pair, as follows:

$$Ls = [CLS, Ls^b, SEP, Ls^a], \quad Ls \in \mathcal{R}^{L \times (2W+2) \times d}$$
(8)

We refer to such longer line as paired line.

**Query Construction**. We also use a multi-layer Transformer to model each paired line and make its tokens interact with each other. Since the tokens in each paired line are in order, we further adopt the positional encoding [22]. For each paired line, we use the hidden state of its first token, i.e. $CLS$, as its feature vector. The feature vectors of all paired lines are denoted as $H^l = [h_1^l, h_2^l, ..., h_L^l]$. Then, we mask the vectors of the unchanged lines and apply average pooling to the masked $H^l$ for obtaining the change-aware query $q^l$.

$$q^l = \text{MaskPooling}(H^l, mask^l)$$
(9)

where $mask^l$ indicates the changed lines. $q^l$ is expected to capture the information of the changed lines.

**Query Back & Merging**. These two steps are similar to those of the token-level query-back mechanism. We use $q^l$ as the query to "retrieve" information from $H^b$ and $H^a$, and produce the final code change representation $v^l$.

*3) Hybrid Query-Back Mechanism:* For some tasks concerning both token changes and line additions/deletions, both the token-level and the line-level change information can be
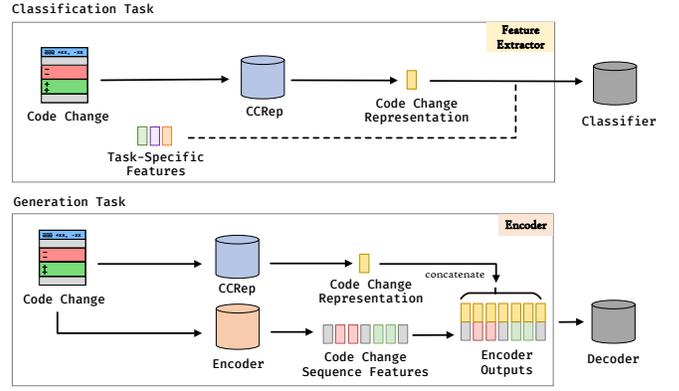


Fig. 3. The usage of CCRep in code-change-related classification and generation tasks.

beneficial. Considering this, we further propose the hybrid query-back mechanism, which fuses the token-level and the line-level query-back mechanisms. Specifically, we first obtain the code change representations produced by the token-level and the line-level query-back mechanisms, i.e., $v^t$ and $v^l$. Then, we linearly project them into the same feature space, normalize them with layer normalization [34], and finally merge them through element-wise addition to produce the final code change representation $v^h$:

$$v^h = \text{LayerNorm}(W_t^T v^t) + \text{LayerNorm}(W_l^T v^l)$$
(10)

where $W_t$ and $W_l$ are learnable parameters of this module.

*D. The Usage of CCRep*

CCRep takes as input a code change and outputs a single numerical vector as the representation of the code change. It acts as a general encoder and can be used to replace or enhance existing code change encoders in both classification and generation tasks, as shown in Figure 3. To use CCRep for classification, the most straightforward way is to combine CCRep and task-specific components (e.g. a commit message encoder), connect them to a classifier, such as a Multi-Layer Perceptron (MLP), and jointly train CCRep, task-specific components and the classifier on the target classification task. Please note that some tasks only take code changes as input, so there may be no task-specific component. For generation tasks, CCRep can be plugged into a neural generation model to enhance its representations of code changes. Specifically, given a code change, neural generation models usually leverage a neural component named encoder to encode it into a sequence of feature vectors, and fed such vectors into another neural component named decoder for generation. We can concatenate the single vector produced by CCRep and each feature vector produced by the encoder of a neural generation model, respectively, to enhance the representations of code changes, and feed those enhanced representations into the decoder for better generation. Also, CCRep and the original encoder and decoder are jointly trained on the target generation task, instead of using a two-stage training scheme as CC2Vec.

## IV. Experiments

To evaluate the effectiveness of CCRep and demonstrate its applicability to diverse downstream tasks, we apply CCRep to three different tasks related to code changes: commit message generation [13], automated patch correctness assessment [14], and just-in-time defect prediction [10]. For each task, we aim to answer two research questions:

**RQ1: How effective is CCRep compared to the state-of-the-art approaches?** We evaluate the three variants of CCRep, which use the token-level, the line-level and the hybrid query-back mechanisms and are referred to as CCRep^token, CCRep^line and CCRep^hybrid, respectively, on the target task, and compare them against the state-of-the-art approaches.

**RQ2: How effective is each component in CCRep?** There are two main components in CCRep, i.e., the pre-trained code model and the query-back mechanism. We conduct ablation studies on each task to investigate their contributions to CCRep's effectiveness.

### A. Commit Message Generation

*1) Background:* When committing a change into a software repository, developers are encouraged to write a textual message to describe the content and intent of this change, namely, commit message. Although commit messages are valuable for program comprehension and software maintenance, developers usually neglect writing high-quality commit messages due to time pressure and lack of direct motivation [19]. To alleviate this problem, researchers propose the Commit Message Generation (CMG) task, which takes as input a commit and outputs a concise description summarizing this commit.

*2) Baselines:* The most recent state-of-the-art CMG approach is FIRA proposed by Dong et.al. [35], which employs Graph Convolution Network (GCN) [36] as the encoder and a Transformer as the decoder. FIRA further proposes a dual copy mechanism to copy both tokens and sub-tokens from the input code change during generation. According to the evaluation results presented in the FIRA paper, we additionally choose the best-performing retrieval-based CMG approach, i.e., NNGen [20], the best-performing learning-based approach, i.e., CODISUM [37], and CoReC [13], which combines retrieval-based and learning-based methods, as the baselines. We also use LogGen [8] as a baseline, since it is proposed in the CC2Vec paper and is the adaption of CC2Vec to CMG.

*3) Our Approach:* Following the description in Section III-D, we plug CCRep into a neural generation model by concatenating the feature vector produced by CCRep with the sequence features of code diff produced by the pre-trained code model (i.e., CodeBERT), as illustrated in the lower part of Figure 3.

*4) Experimental Setting:* We evaluate our approach on CMG with the datasets used in the CoReC paper [13] and the FIRA paper [35]. Both datasets are widely-used for this task. The dataset used by CoReC was initially collected by Jiang et al. [2] and further cleansed by Liu et al. [20], containing 22.0K commits. Since CCRep focuses on code

### TABLE III
### CMG: EVALUATION RESULTS ON THE CoReC DATASET

| Model | BLEU | METEOR | ROUGE-L |
|---|---|---|---|
| LogGen [8] | 9.41 | 5.31 | 12.13 |
| NNGen [20] | 23.29 | 14.26 | 28.68 |
| CODISUM [37] | 13.15 | 7.35 | 16.49 |
| CoReC [13] | 25.27 | 15.34 | 29.73 |
| CCRep^token | **28.24** | **16.99** | **34.23** |
| CCRep^line | 26.65 | 15.74 | 32.01 |
| CCRep^hybrid | 27.25 | 16.58 | 33.30 |

### TABLE IV
### CMG: EVALUATION RESULTS ON THE FIRA DATASET

| Model | BLEU | METEOR | ROUGE-L |
|---|---|---|---|
| LogGen [8] | 8.95 | 8.34 | 10.50 |
| NNGen [20] | 9.16 | 9.53 | 11.24 |
| CODISUM [37] | 16.55 | 12.83 | 19.73 |
| CoReC [13] | 13.03 | 12.04 | 15.47 |
| FIRA [35] | 17.67 | 14.93 | 21.58 |
| CCRep^token | **19.93** | **16.27** | **23.81** |
| CCRep^line | 19.79 | 16.06 | 23.60 |
| CCRep^hybrid | 19.70 | 15.84 | 23.41 |

changes, we filter out the commits with non-code changes, such as binary file changes, file creation or file deletion, resulting in 20.5K commits left. We re-train and re-evaluate the baselines on our filtered dataset (hereon, the CoReC dataset) using the replication packages provided by their authors [38]–[41]. FIRA is not evaluated on this dataset since it needs to parse each code change into two ASTs, but the code diff fragments provided by this dataset are not parsable. Following CoReC, our approach uses an LSTM as the decoder on this dataset.

The dataset used by FIRA (hereon, the FIRA dataset) was published by Xu et.al. [37], which contains 75K, 8K and 7.6K commit-message pairs in the training, validation and test sets, respectively. Following FIRA, our approach adopts a Transformer as the decoder and also equips with a copy mechanism and the identifier abstraction used by FIRA. Since Dong et al [35] have evaluated FIRA and other baseline approaches on this dataset, the evaluation results of all the baselines shown in Table IV are directly copied from the FIRA paper.

Following prior work [13], [35], the Adam optimizer [42] is used to minimize the average cross-entropy loss during training, and BLEU [43], METEOR [44] and ROUGE-L [45] are used as evaluation metrics. However, Tao et al. [46] reported that a variant of the original BLEU called B-Norm BLEU is more correlated with human judgment on the quality of generated commit messages. Thus, we use the B-Norm BLEU as the substitution of the original BLEU-4 in our experiments (denoted as BLEU as well).

*5) Results for RQ1:* Experimental results on CMG are shown in Table III and Table IV. CCRep^token, CCRep^line and CCRep^hybrid stand for the three variants of our approach. It is

| Model | BLEU | METEOR | ROUGE-L |
|---|---|---|---|
| CCRep − CodeBERT | 27.90 | 16.56 | 33.36 |
| CCRep − QueryBack | 26.23 | 15.72 | 31.78 |
| CCRep | **28.24** | **16.99** | **34.23** |

TABLE VI
CMG: ABLATION RESULTS ON THE FIRA DATASET

| Model | BLEU | METEOR | ROUGE-L |
|---|---|---|---|
| CCRep − CodeBERT | 18.77 | 15.54 | 22.37 |
| CCRep − QueryBack | 17.88 | 14.78 | 21.33 |
| CCRep | **19.93** | **16.27** | **23.81** |

shown that our approach outperforms all the baselines in terms of all metrics, indicating the effectiveness of CCRep on CMG. Specifically, on the CoReC dataset, our best-performing variant, i.e., CCRep$^{token}$, outperforms the best-performing baseline, i.e., CoReC, by 11.8%, 10.8% and 15.1% in terms of BLEU, METEOR and ROUGE-L, respectively. As for the FIRA dataset, CCRep$^{token}$ outperforms FIRA by 12.8%, 9.0% and 10.3% in terms of the three metrics. We conduct statistical significance tests using paired bootstrap resampling with 1000 resamples following Koehn [47]. All the p-values are less than 0.001, indicating that the performance differences between our approach and the baselines are significant.

*6) Results for RQ2:* To answer RQ2, we compare the best-performing variant of our approach, i.e., CCRep$^{token}$, with two special models, namely CCRep−CodeBERT and CCRep−QueryBack. The former replaces CodeBERT in CCRep$^{token}$ with a RoBERTa-base model [48], which is widely used as a baseline encoder in the NLP community. CCRep−QueryBack removes the query-back mechanism from CCRep$^{token}$ and directly feeds the diff's feature vectors produced by the CodeBERT encoder to the decoder for generation.

Experimental results in Table V and Table VI show that our approach outperforms the two special models in terms of all metrics, indicating the effectiveness of the pre-trained code model and the query-back mechanism. We conduct statistical significance tests like RQ1. The p-values of CCRep compared to CCRep−QueryBack are all less than 0.001, which means that the query-back mechanism brings statistically significant performance improvements. The p-values of CCRep compared to CCRep−CodeBERT are all less than 0.05 except the one calculated on the CoReC dataset in terms of BLEU. This indicates that CodeBERT are also significantly beneficial in most cases.

### B. Automated Patch Correctness Assessment

*1) Background:* Automated Program Repair (APR) [49] aims to automatically generate bug-fixing patches. Many APR approaches follow a generate-and-validate methodology, which examines the generated patches with developer-provided test suits. Because test suits may be inadequate to cover all possible cases, a generated patch that passes all test cases (i.e., plausible patch) may still be incorrect (i.e., overfitting patch). This is known as the *overfitting* problem of APR [50], [51]. To alleviate this problem, many techniques have been proposed to automatically identify correct patches among plausible patches, namely Automated Patch Correctness Assessment (APCA), which is a binary classification task of code changes.

*2) Baselines:* We consider existing APCA approaches that take patches as input as baselines. CACHE, recently proposed by Lin et al. [14], is the state-of-the-art APCA approach. Given a patch, CACHE uses AST paths [52] to extract the features of its deleted code, its added code and their context, and integrates such features as the patch feature for prediction. Another recent work from Tian et al. [6] leverages pre-trained representation learning models, such as BERT [23] and CC2Vec [8], with some feature comparison functions [8] to extract patch features and uses classifiers like Logistic Regression (LR) and Decision Tree (DT) for prediction.

*3) Our Approach:* As illustrated in Section III-D and the upper part of Figure 3, to apply CCRep to APCA, we simply connect CCRep to a two-layer Multi-Layer Perceptron (MLP) classifier for prediction.

*4) Experimental Setting:* We use the two datasets constructed by Lin et al. [14] for evaluation and refer to them as CACHE-Small and CACHE-Large. CACHE-Small was constructed by merging and deduplicating the plausible patches collected by Wang et al. [53] and Tian et al. [6], containing 1,183 patches from the Defects4J benchmark. CACHE-Large has 49.7K patches and was built by merging and deduplicating the patches from RepairThemAll [54] and ManySStuBs4J [55]. Both datasets are roughly balanced. Five widely-used classification metrics, including Accuracy, Precision, Recall, F1 and AUC, are used for evaluation. Following Lin et al. [14], we perform 5-fold cross-validation on both datasets, use the Adam optimizer to minimize the binary cross-entropy loss during training, and employ Dropout [56] in the classifier.

Because we use the same datasets and experimental settings as Lin et al. [14], the evaluation results of the baselines are directly borrowed from their paper. For space constraints, we only present the variants of Tian et al.'s approach that: 1) achieve the best F1 or the best AUC, or 2) use CC2Vec and achieve the best F1 or AUC among all the variants using CC2Vec.

*5) Results for RQ1:* The evaluation results are shown in Table VII and Table VIII. It can be seen that on CACHE-Small, our approach outperforms CACHE in terms of all metrics except precision. Our best variant, i.e., CCRep$^{hybrid}$, substantially outperforms CACHE in terms of recall by 12.3 points and improves CACHE in terms of F1 and AUC by 3.9 points and 8.2 points, respectively. For Tian et al.'s approach, our approach outperforms all of its variants, including the CC2Vec-based variants, in terms of F1 and AUC by large margins.

TABLE VII
APCA: EVALUATION RESULTS ON CACHE-SMALL

| Model | Acc. | Pre. | Rec. | F1 | AUC |
|---|---|---|---|---|---|
| LR + CC2Vec [6] | 64.9 | 62.4 | **90.1** | 73.7 | 68.6 |
| LR + code2vec [6] | 66.8 | 68.6 | 72.9 | 70.6 | 70.2 |
| CACHE [14] | 75.4 | **79.5** | 76.5 | 78.0 | 80.3 |
| CCRep$^{token}$ | 80.5 | 73.8 | 89.0 | 80.6 | 88.1 |
| CCRep$^{line}$ | 81.2 | 75.9 | 86.4 | 80.6 | 88.3 |
| CCRep$^{hybrid}$ | **82.2** | 76.3 | 88.8 | **81.9** | **88.5** |

TABLE VIII
APCA: EVALUATION RESULTS ON CACHE-LARGE

| Model | Acc. | Pre. | Rec. | F1 | AUC |
|---|---|---|---|---|---|
| DT + BERT [6] | 95.7 | 93.9 | 97.4 | 95.6 | 95.9 |
| LR + Doc2Vec [6] | 90.4 | 91.9 | 88.0 | 89.9 | 96.1 |
| DT + CC2Vec [6] | 95.6 | 95.4 | 95.7 | 95.5 | 95.7 |
| CACHE [14] | 98.6 | 98.9 | 98.2 | 98.6 | 98.9 |
| CCRep$^{token}$ | 99.4 | 99.6 | 99.2 | 99.4 | 99.97 |
| CCRep$^{line}$ | **99.7** | 99.8 | **99.6** | **99.7** | 99.98 |
| CCRep$^{hybrid}$ | 99.6 | **99.9** | 99.3 | 99.6 | **99.99** |

TABLE IX
APCA: EVALUATION RESULTS OF THE ABLATION STUDIES

| Dataset | Model | Acc. | Pre. | Rec. | F1 | AUC |
|---|---|---|---|---|---|---|
| CACHE-Small | CCRep − CodeBERT | 77.5 | 72.6 | 81.6 | 76.6 | 86.1 |
| | CCRep − QueryBack | 79.3 | 74.4 | 83.2 | 78.4 | 87.0 |
| | CCRep | **82.2** | **76.3** | **88.8** | **81.9** | **88.5** |
| CACHE-Large | CCRep − CodeBERT | 99.6 | 99.7 | 99.6 | 99.6 | **99.99** |
| | CCRep − QueryBack | **99.7** | 99.8 | **99.7** | **99.7** | **99.99** |
| | CCRep | 99.6 | **99.9** | 99.3 | 99.6 | **99.99** |

On CACHE-Large, our approach also outperforms Tian et al.'s approach and CACHE in terms of all metrics. Our best variant, i.e., CCRep$^{hybrid}$, achieves an F1 of 99.6% and an AUC of 99.99%. One possible reason behind the impressive performance of both CACHE and our approach is that CACHE-Large is way larger than CACHE-Small, which may ease the learning. Another possible reason is that CACHE-Large is synthetic to some extent. Specifically, according to Lin et al. [14], in CACHE-Large, almost all correct patches are human-written patches, while all overfitting patches are generated by APR tools. Therefore, an approach can perform well if it is able to tell from human-written patches and APR-generated patches.

Another thing worth mentioning is that CCRep$^{hybrid}$ outperforms the other two variants on the two datasets in terms of AUC, which highlights that both the token-level and the line-level change information is useful for APCA.

*6) Results for RQ2:* We compare the best-performing variant, i.e., CCRep$^{hybrid}$, with CCRep−CodeBERT and CCRep−QueryBack. CCRep−CodeBERT replaces Code-BERT in CCRep$^{hybrid}$ with the RoBERTa-base model [48], which is a widely-used baseline encoder. CCRep−QueryBack directly uses CodeBERT to encode a patch's diff and uses the contextual embedding of a special token *CLS* inserted at

the beginning of the diff as the code change representation. Table IX presents the results, which indicate that both the pre-trained code model and the query back mechanism positively affect the effectiveness of CCRep on CACHE-Small. However, there are no significant performance differences between the two special models and CCRep on CACHE-Large. A possible explanation for this is that distinguishing correct and overfitting patches in CACHE-Large is not very hard and less powerful models can also fit the data.

### C. Just-in-Time Defect Prediction

*1) Background:* Software defects are inevitable and may substantially affect businesses and even people's lives [11]. On the other hand, the size and complexity of modern software systems grow significantly, making it hard and costly to find defects from them. To this end, just-in-time defect prediction (JIT-DP) has been proposed to identify defective code changes and provide in-time feedback when developers commit changes to the code base [4], [5], [10], [11]. Given a commit, JIT-DP targets at predicting whether it is defective and is a binary classification task just like APCA. However, different from APCA, JIT-DP targets real-world defective commits and can take as input both the code change and the commit message in a commit.

*2) Baselines:* We use three state-of-the-art JIT-DP approaches, i.e., DeepJIT [11], CC2Vec [8] and LAPredict [10] as baselines. Given a commit, DeepJIT leverages CNNs to extract feature vectors from the code change and the commit message, respectively, and concatenates such vectors as the commit vector for prediction. Hoang et al. [8] use CC2Vec to encode a code change into a feature vector and appends such vector to the commit feature extracted by DeepJIT for prediction. We also refer to this approach as CC2Vec for convenience. LAPredict, proposed by Zeng et al. [10], uses the number of added code lines as the commit feature and adopts an LR classifier to perform JIT-DP.

*3) Our Approach:* We follow Section III-D and the upper part of Figure 3 to apply CCRep in JIT-DP, similar to APCA. The only difference is that we concatenate the output of CCRep with the commit message feature produced by a CNN (i.e., the task-specific features in Figure 3) for prediction, following DeepJIT [11].

*4) Experimental Setting:* We compare CCRep with the baselines on the dataset constructed by Zeng et al. [10]. The dataset contains six large-scale projects, i.e., OpenStack, QT, Go, Gerrit, Platform and JDT, covering different programming languages. Considering that CodeBERT is not pre-trained on C++ code [57], we exclude QT for evaluation. We conduct our experiments following the within-project setting used by Zeng et al. [10] and directly borrow the evaluation results of the baselines from the LAPredict paper. Since the dataset is highly unbalanced, we follow prior work [8], [10], [11] and also use AUC as the evaluation metric. Also following prior work [11], we use the Adam optimizer to minimize the binary cross-entropy loss and employ Dropout in the classifier.

TABLE X
JIT-DP: EVALUATION RESULTS IN TERMS OF AUC

| Model | OpenStack | JDT | Go | Platform | Gerrit | Mean |
|---|---|---|---|---|---|---|
| DeepJIT [11] | 71.32 | 67.01 | 68.91 | 77.12 | 70.25 | 70.92 |
| CC2Vec [8] | 72.27 | 66.53 | 69.17 | 76.13 | 69.86 | 70.79 |
| LAPredict [10] | 74.91 | 67.57 | 68.31 | 74.61 | 74.95 | 72.07 |
| CCRep$^{token}$ | 75.37 | 68.11 | **75.63** | 81.91 | 76.89 | 75.58 |
| CCRep$^{line}$ | **76.45** | **68.96** | 75.60 | 82.08 | **77.35** | **76.09** |
| CCRep$^{hybrid}$ | 75.63 | 66.63 | 75.48 | **82.18** | 77.01 | 75.39 |

TABLE XI
JIT-DP: ABLATION RESULTS IN TERMS OF AUC

| Model | OpenStack | JDT | Go | Platform | Gerrit | Mean |
|---|---|---|---|---|---|---|
| CCRep $-$ CodeBERT | 75.78 | 68.21 | 75.04 | 81.15 | 74.89 | 75.01 |
| CCRep $-$ QueryBack | 75.70 | 66.69 | **76.30** | **82.28** | 74.82 | 75.16 |
| CCRep | **76.45** | **68.96** | 75.63 | 82.18 | **77.35** | **76.09** |

*5) Results for RQ1:* Experimental results on JIT-DP are summarized in Table X. We can see that CCRep$^{line}$ is the best-performing variant on average, but the performance differences among the three variants are small. CCRep$^{line}$ improves LAPredict, CC2Vec and DeepJIT by 5.6%, 7.4% and 7.3% on average, indicating the effectiveness of CCRep in representing code changes for JIT-DP. We also notice that different projects prefer different query-back mechanisms, which may indicate the different defect characteristics of different projects.

*6) Results for RQ2:* We also conduct ablation studies on JIT-DP to answer RQ2. Similar to what we do in APCA, we build and evaluate two special approaches, i.e., CCRep−CodeBERT and CCRep−QueryBack. As shown in Table XI, the average performance of our approach degrades without CodeBERT or QueryBack, highlighting the effectiveness of the two components in CCRep. We also notice that on Go and Platform, our approach is slightly worse than CCRep−QueryBack. After manual inspection, we speculate that this is because a significant number of commits in Go and Platform contain almost only deleted or added code with little context, where the query-back mechanism may bring no benefit. However, since the performance of our approach and CCRep−QueryBack is very close, we argue the query-back mechanism can be viewed as harmless on Go and Platform. In summary, both the pre-trained code model and the query-back mechanism contribute to the effectiveness of CCRep.

## V. DISCUSSION

### A. The Variants of CCRep

We propose three variants of CCRep, i.e., CCRep$^{token}$, CCRep$^{line}$ and CCRep$^{hybrid}$. We can see from Section IV that different tasks prefer different variants. On CMG, CCRep$^{token}$ performs best, which is probably because that a commit message is generated token by token and the key words in it can often be found from the changed code tokens. For APCA, the best-performing variant is CCRep$^{hybrid}$. After inspecting the dataset and the results, we think this is reasonable because a generated patch can either change a few tokens (e.g., replace

"==" with "!="), or insert/delete several lines (e.g., insert a `NullPointer` check). As for JIT-DP, CCRep$^{line}$ is the best. Based on our inspection, a possible reason is that a defective commit often adds, deletes or modifies multiple lines instead of only a few tokens. Based on these findings, before applying CCRep to other tasks, we encourage the user to analyze the characteristics of the target task and her dataset first and choose the most suitable variant. On the other hand, on each task, the three variants all outperform the state-of-the-art baselines. This indicates our approach's effectiveness and makes us believe that any of the variants can serve as a strong baseline for code-change-related tasks.

### B. Limitations

As discussed in Section IV-C6, if a code change contains little context code, the query-back mechanism may bring no benefit. Because in such situation, the changed code and the whole code change contain the same information. Another limitation of CCRep is that the lengths of the code changes that can be processed by CCRep are limited by the pre-trained code model. In detail, the length of the before-change or the after-change code should not exceed the length limit of the used pre-trained model, which is usually 512. Fortunately, committing small and coherent code changes has become a widely-recognized good practice, and many techniques have been proposed to decompose tangled code changes [58]–[62]. Such practice and techniques can help reduce long code changes and alleviate this limitation.

### C. Threats to Validity

Threats to internal validity refer to the errors and bias in our experiments. To mitigate such threats, for each task, we try our best to follow the settings used by the state-of-the-art baselines, re-use the evaluation results reported by prior work when possible and use the existing implementations of the baselines to conduct experiments. We have double checked our code and data and made them publicly available. Threats to external validity concern the generalization of CCRep. Although we have applied CCRep to three different tasks and achieved superior performance, we cannot claim that CCRep can be applied to or perform well on all code-change-related tasks. However, the three tasks have either different inputs or different outputs, care about diverse characteristics of code changes, and cover synthetic and manually-written patches, various projects and multiple programming languages. Therefore, we believe this threat is limited. In addition, based on our evaluation results, we argue that CCRep can at least serve as a strong baseline to help better solve code-change-related problems. To minimize threats to construct validity, we choose evaluation metrics following previous studies.

## VI. RELATED WORK

The majority of the studies related to code change representation target a specific downstream task and learn code change representations through a task-specific architecture. Some of them flatten code or code changes as token sequences

for representation learning [2], [13], [15], [37], [63]–[65]. For example, for commit message generation, applying RNN like LSTM [66] on diffs is a widely-used approach to extract code change features [2], [13], [37]. To automate comment updates with code changes, Liu et al. [63] aligned the tokens of the before-change and after-change code to form an edit sequence and fed such sequence into an LSTM-based encoder to obtain code change representations. To identify security patches, Zhou et al. [15] utilized two LSTMs to respectively learn the statement-level features of the added and deleted code in a patch and merged their features with a multi-layer convolutional neural network.

Some studies leverage the syntactic structure of code, e.g., AST, to enhance code change representation learning [14], [16], [18], [67], [68]. For example, to assess patch correctness, Lin et al. [14] proposed to extract and encode the changed AST paths and the unchanged AST paths, respectively, and merge their feature vectors as the code change representation. Yin et al. [16], Panthaplackel et al. [18] and Dong et al. [35] proposed to converted the two ASTs of a code change into a graph and use graph neural networks, e.g., GGNN [69] and GCN [36], to learn code change representations from the graph.

Recently, several studies adopted pre-trained code models to represent code changes for specific downstream tasks [3], [6], [12], [70]. For instance, to identify silent vulnerability fixes, Zhou et al. [3] leveraged CodeBERT [21] to encode each changed file in a commit and merged the feature vectors of all changed files as the code change representation. Lin et al. [12] leveraged pre-trained code models to encode commits for recovering links between issues and commits. Zhou et al. [70] also investigated the generalizability of CodeBERT on JIT-DP. However, they only considered one classification task and only used the data from two projects for evaluation.

Only a few studies aim at learning general-purpose code change representations [8], [16]. Yin et al. [16] proposed to learn distributed representations of code edits by training an auto-encoder to reconstruct code edits. Hoang et al. [8] proposed CC2Vec, which leverages a hierarchical attention network and multiple comparison functions to learn code change representations. As discussed in Section I, Yin et al.'s approach only focuses on small code edits (i.e., a single hunk with no more than 3 changed lines), while CCRep targets commit-level changes. Also, it lacks explicit interaction between the changed code and the whole code change. We have tried to compare the performance of Yin et al.'s approach with that of CCRep. However, neither Yin et al. evaluated their approach on the three tasks used in this work, nor they made their implementation publicly available. CC2Vec only considers the changed code and ignores the context, and it requires commit messages as labels, which are not always available. Besides, our evaluation results show that CCRep outperforms CC2Vec on the three tasks by substantial margins.

In summary, our work differs from prior work in several folds: First, our approach acts as a general code change encoder and can be used in diverse code-change-related tasks. Second, our approach is equipped with a pre-trained code model and the query-back mechanism, and is technically different. Third, our evaluation results show that our approach outperforms the state-of-the-art techniques on three tasks. In addition, this work also investigates the generalizability of pre-trained code models on diverse code-change-related tasks.

## VII. CONCLUSION

We propose a novel approach named CCRep to learn code change representations. It acts as a code change encoder and can be jointly trained with and used in diverse code-change-related tasks. CCRep leverages a pre-trained code model to obtain high-quality contextual embeddings and better handle datasets of different sizes, and a novel mechanism named query back to highlight the changed code and adaptively capture related context information. We evaluate CCRep on one generation task and two classification tasks. Experimental results show that CCRep outperforms the state-of-the-art approaches on each task and both the pre-trained code model and the query-back mechanism contribute to its effectiveness. In the future, we plan to apply our approach to more code-change-related tasks and improve it to encode long and structured code changes.

## VIII. ACKNOWLEDGMENTS

## REFERENCES

[1] I. I. Brudaru and A. Zeller, "What is the long-term impact of changes?" in *Proceedings of the 2008 International Workshop on Recommendation Systems for Software Engineering*, 2008, pp. 30–32.

[2] S. Jiang, A. Armaly, and C. McMillan, "Automatically generating commit messages from diffs using neural machine translation," in *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering*, 2017, pp. 135–146.

[3] J. Zhou, M. Pacheco, Z. Wan, X. Xia, D. Lo, Y. Wang, and A. E. Hassan, "Finding A needle in a haystack: Automated mining of silent vulnerability fixes," in *Proceedings of the 36th IEEE/ACM International Conference on Automated Software Engineering*, 2021, pp. 705–716.

[4] S. Kim, E. J. W. Jr., and Y. Zhang, "Classifying software changes: Clean or buggy?" *IEEE Trans. Software Eng.*, vol. 34, no. 2, pp. 181–196, 2008.

[5] Y. Kamei, T. Fukushima, S. McIntosh, K. Yamashita, N. Ubayashi, and A. E. Hassan, "Studying just-in-time defect prediction using cross-project models," *Empir. Softw. Eng.*, vol. 21, no. 5, pp. 2072–2106, 2016.

[6] H. Tian, K. Liu, A. K. Kaboré, A. Koyuncu, L. Li, J. Klein, and T. F. Bissyandé, "Evaluating representation learning of code changes for predicting patch correctness in program repair," in *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, 2020, pp. 981–992.

[7] H. Ye, J. Gu, M. Martinez, T. Durieux, and M. Monperrus, "Automated classification of overfitting patches with statically extracted code features," *IEEE Trans. Software Eng.*, pp. 1–1, 2021.

[8] T. Hoang, H. J. Kang, D. Lo, and J. Lawall, "Cc2vec: distributed representations of code changes," in *Proceedings of the 42nd International Conference on Software Engineering*, 2020, pp. 518–529.

[9] Y. Tian, J. Lawall, and D. Lo, "Identifying linux bug fixing patches," in *Proceedings of the 34th International Conference on Software Engineering*, 2012, pp. 386–396.

[10] Z. Zeng, Y. Zhang, H. Zhang, and L. Zhang, "Deep just-in-time defect prediction: how far are we?" in *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2021, pp. 427–438.

[11] T. Hoang, H. K. Dam, Y. Kamei, D. Lo, and N. Ubayashi, "Deepjit: an end-to-end deep learning framework for just-in-time defect prediction," in *Proceedings of the 16th International Conference on Mining Software Repositories*, 2019, pp. 34–45.

[12] J. Lin, Y. Liu, Q. Zeng, M. Jiang, and J. Cleland-Huang, "Traceability transformed: Generating more accurate links with pre-trained BERT models," in *Proceedings of the 43rd IEEE/ACM International Conference on Software Engineering*, 2021, pp. 324–335.

[13] H. Wang, X. Xia, D. Lo, Q. He, X. Wang, and J. Grundy, "Context-aware retrieval-based deep commit message generation," *ACM Trans. Softw. Eng. Methodol.*, vol. 30, no. 4, pp. 56:1–56:30, 2021.

[14] B. Lin, S. Wang, M. Wen, and X. Mao, "Context-aware code change embedding for better patch correctness assessment," *ACM Trans. Softw. Eng. Methodol.*, pp. 1–1, 2021.

[15] Y. Zhou, J. K. Siow, C. Wang, S. Liu, and Y. Liu, "SPI: automated identification of security patches via commits," *ACM Trans. Softw. Eng. Methodol.*, vol. 31, no. 1, pp. 13:1–13:27, 2022.

[16] P. Yin, G. Neubig, M. Allamanis, M. Brockschmidt, and A. L. Gaunt, "Learning to represent edits," in *Proceedings of the 7th International Conference on Learning Representations*, 2019.

[17] T. H. M. Le, D. Hin, R. Croft, and M. A. Babar, "Deepcva: Automated commit-level vulnerability assessment with deep multi-task learning," in *Proceedings of the 36th IEEE/ACM International Conference on Automated Software Engineering*, 2021, pp. 717–729.

[18] S. Panthaplackel, J. J. Li, M. Gligoric, and R. J. Mooney, "Deep just-in-time inconsistency detection between comments and source code," in *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021, pp. 427–435.

[19] R. Dyer, H. A. Nguyen, H. Rajan, and T. N. Nguyen, "Boa: a language and infrastructure for analyzing ultra-large-scale software repositories," in *Proceedings of the 35th International Conference on Software Engineering*, 2013, pp. 422–431.

[20] Z. Liu, X. Xia, A. E. Hassan, D. Lo, Z. Xing, and X. Wang, "Neural-machine-translation-based commit message generation: how far are we?" in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018, pp. 373–384.

[21] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, "Codebert: A pre-trained model for programming and natural languages," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, ser. Findings of ACL, vol. EMNLP 2020, 2020, pp. 1536–1547.

[22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proceedings of the Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, 2017, pp. 5998–6008.

[23] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2019, pp. 4171–4186.

[24] 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," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.

[25] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[26] F. Liu, G. Li, Y. Zhao, and Z. Jin, "Multi-task learning based pre-trained language model for code completion," in *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, 2020, pp. 473–485.

[27] A. Svyatkovskiy, S. K. Deng, S. Fu, and N. Sundaresan, "Intellicode compose: code generation using transformer," in *Proceedings of the 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 1433–1443.

[28] W. U. Ahmad, S. Chakraborty, B. Ray, and K. Chang, "Unified pre-training for program understanding and generation," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021, pp. 2655–2668.

[29] Y. Wang, W. Wang, S. R. Joty, and S. C. H. Hoi, "Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 8696–8708.

[30] "Wordpress for android," https://github.com/wordpress-mobile/WordPress-Android, 2022.

[31] "Python difflib library," https://docs.python.org/3/library/difflib.html, 2022.

[32] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016.

[33] R.-M. Karampatsis, H. Babii, R. Robbes, C. Sutton, and A. Janes, "Big Code != Big Vocabulary: Open-Vocabulary Models for Source Code," in *Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering*, 2020, pp. 1073–1085.

[34] L. J. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016.

[35] J. Dong, Q. Zhu, Z. Sun, Z. Li, Y. Lou, W. Zhang, and D. Hao, "FIRA: Fine-Grained Graph-Based Code Change Representation for Automated Commit Message Generation," in *Proceedings of the IEEE/ACM 44rd International Conference on Software Engineering*, 2022, p. 12.

[36] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: https://openreview.net/forum?id=SJU4ayYgl

[37] S. Xu, Y. Yao, F. Xu, T. Gu, H. Tong, and J. Lu, "Commit message generation for source code changes," in *Proceedings of the 28th International Joint Conference on Artificiaal Intelligence*, 2019, pp. 3975–3981.

[38] "The code of CoDiSum," https://github.com/SoftWiser-group/CoDiSum, 2019.

[39] "The code of nngen," https://github.com/Tbabm/nngen, 2018.

[40] "The code of loggen," https://github.com/soarsmu/CC2Vec, 2020.

[41] "The code of corec," http://tiny.cc/o4j4oz., 2021.

[42] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the 3rd International Conference on Learning Representations*, 2015.

[43] 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 of the Association for Computational Linguistics*, 2002, pp. 311–318.

[44] S. Banerjee and A. Lavie, "Meteor: An automatic metric for mt evaluation with improved correlation with human judgments," in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.

[45] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," in *Text summarization branches out*, 2004, pp. 74–81.

[46] W. Tao, Y. Wang, E. Shi, L. Du, S. Han, H. Zhang, D. Zhang, and W. Zhang, "On the evaluation of commit message generation models: An experimental study," in *Proceedings of the 2021 IEEE International Conference on Software Maintenance and Evolution*, 2021, pp. 126–136.

[47] P. Koehn, "Statistical significance tests for machine translation evaluation," in *Proceedings of the 9th Conference on Empirical Methods in Natural Language Processing*, 2004, pp. 388–395.

[48] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," *CoRR*, vol. abs/1907.11692, 2019. [Online]. Available: http://arxiv.org/abs/1907.11692

[49] C. L. Goues, M. Pradel, and A. Roychoudhury, "Automated program repair," *Commun. ACM*, vol. 62, no. 12, pp. 56–65, 2019.

[50] Z. Qi, F. Long, S. Achour, and M. C. Rinard, "An analysis of patch plausibility and correctness for generate-and-validate patch generation systems," in *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, 2015, pp. 24–36.

[51] E. K. Smith, E. T. Barr, C. L. Goues, and Y. Brun, "Is the cure worse than the disease? overfitting in automated program repair," in *Proceedings of the 10th Joint Meeting on Foundations of Software Engineering*, 2015, pp. 532–543.

[52] U. Alon, M. Zilberstein, O. Levy, and E. Yahav, "code2vec: Learning distributed representations of code," *Proc. ACM Program. Lang.*, vol. 3, no. POPL, pp. 40:1–40:29, 2019.

[53] S. Wang, M. Wen, B. Lin, H. Wu, Y. Qin, D. Zou, X. Mao, and H. Jin, "Automated patch correctness assessment: How far are we?" in *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, 2020, pp. 968–980.

[54] T. Durieux, F. Madeiral, M. Martinez, and R. Abreu, "Empirical review of java program repair tools: A large-scale experiment on 2,141 bugs and 23,551 repair attempts," in *Proceedings of the 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2019, 2019, p. 302–313.

[55] R. Karampatsis and C. Sutton, "How often do single-statement bugs occur?: The manysstubs4j dataset," in *Proceedings of the 17th International Conference on Mining Software Repositories*, 2020, pp. 573–577.

[56] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.

[57] "The pre-train codebert model," https://huggingface.co/microsoft/codebert-base, 2021.

[58] M. Barnett, C. Bird, J. Brunet, and S. K. Lahiri, "Helping developers help themselves: Automatic decomposition of code review changesets," in *Proceedings of the 37th IEEE/ACM International Conference on Software Engineering*, 2015, pp. 134–144.

[59] M. Dias, A. Bacchelli, G. Gousios, D. Cassou, and S. Ducasse, "Untangling fine-grained code changes," in *Proceedings of the 22nd IEEE International Conference on Software Analysis, Evolution, and Reengineering*, 2015, pp. 341–350.

[60] Y. Tao and S. Kim, "Partitioning composite code changes to facilitate code review," in *Proceedings of the 12th IEEE/ACM Working Conference on Mining Software Repositories*, 2015, pp. 180–190.

[61] M. Wang, Z. Lin, Y. Zou, and B. Xie, "Cora: Decomposing and describing tangled code changes for reviewer," in *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering*, 2019, pp. 1050–1061.

[62] B. Shen, W. Zhang, C. Kästner, H. Zhao, Z. Wei, G. Liang, and Z. Jin, "Smartcommit: a graph-based interactive assistant for activity-oriented commits," in *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2021, pp. 379–390.

[63] Z. Liu, X. Xia, M. Yan, and S. Li, "Automating just-in-time comment updating," in *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, 2020, pp. 585–597.

[64] S. Panthaplackel, P. Nie, M. Gligoric, J. J. Li, and R. J. Mooney, "Learning to update natural language comments based on code changes," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 1853–1868.

[65] T. Hoang, J. Lawall, Y. Tian, R. J. Oentaryo, and D. Lo, "Patchnet: Hierarchical deep learning-based stable patch identification for the linux kernel," *IEEE Trans. Software Eng.*, vol. 47, no. 11, pp. 2471–2486, 2021.

[66] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.

[67] S. Liu, C. Gao, S. Chen, N. Lun Yiu, and Y. Liu, "Atom: Commit message generation based on abstract syntax tree and hybrid ranking," *IEEE Trans. Software Eng.*, pp. 1–1, 2020.

[68] S. Brody, U. Alon, and E. Yahav, "A structural model for contextual code changes," *Proc. ACM Program. Lang.*, vol. 4, no. OOPSLA, pp. 215:1–215:28, 2020.

[69] Y. Li, D. Tarlow, M. Brockschmidt, and R. S. Zemel, "Gated graph sequence neural networks," in *Proceedings of the 4th International Conference on Learning Representations*, 2016.

[70] X. Zhou, D. Han, and D. Lo, "Assessing generalizability of CodeBERT," in *2021 IEEE International Conference on Software Maintenance and Evolution*, 2021, pp. 425–436.