

Learning to Represent Patches

Xunzhu Tang
xunzhu.tang@uni.lu
University of Luxembourg
Luxembourg

Weiguo Pian
weiguo.pian@uni.lu
University of Luxembourg
Luxembourg

Andrew Habib
andrew.a.habib@gmail.com
University of Luxembourg
Luxembourg

Haoye Tian*
haoye.tian@uni.lu
University of Luxembourg
Luxembourg

Saad Ezzini
s.ezzini@lancaster.ac.uk
Lancaster University
United Kingdom

Jacques Klein
jacques.klein@uni.lu
University of Luxembourg
Luxembourg

Zhenghan Chen
1979282882@pku.edu.cn
Peking University
China

Abdoul Kader Kaboré
abdoukader.kabore@uni.lu
University of Luxembourg
Luxembourg

Tegawendé F. Bissyandé
tegawende.bissyande@uni.lu
University of Luxembourg
Luxembourg

ABSTRACT

Patch representation plays a pivotal role in automating numerous software engineering tasks, such as classifying patch correctness or generating natural language summaries of code changes. Recent studies have leveraged deep learning to derive effective patch representation, primarily focusing on capturing changes in token sequences or Abstract Syntax Trees (ASTs). However, these current state-of-the-art representations do not explicitly calculate the intention semantic induced by the change on the AST, nor do they optimally explore the surrounding contextual information of the modified lines. To address this, we propose a new patch representation methodology, Patcherizer, which we refer to as our tool. This innovative approach explores the intention features of the context and structure, combining the context around the code change along with two novel representations. These new representations capture the sequence intention inside the code changes in the patch and the graph intention inside the structural changes of AST graphs before and after the patch. This comprehensive representation allows us to better capture the intentions underlying a patch. Patcherizer builds on graph convolutional neural networks for the structural input representation of the intention graph and on transformers for the intention sequence representation. We assess the generalizability of Patcherizer's learned embeddings on three tasks: (1) Generating patch description in NL, (2) Predicting patch correctness in program repair, and (3) Patch intention detection. Experimental results show that the learned patch representation is effective for all three tasks and achieves superior performance to SOTA approaches. For instance, on the popular task of patch description generation

(a.k.a. commit message generation), Patcherizer achieves an average improvement of 19.39%, 8.71%, and 34.03% in terms of BLEU, ROUGE-L, and METEOR metrics, respectively.

1 INTRODUCTION

A software patch represents the source code differences between two software versions. It has a dual role: on the one hand, it serves as a formal summary of the code changes that a developer intends to make on a code base; on the other hand, it is used as the main input specification for automating software evolution. Patches are thus a key artifact that is pervasive across the software development life cycle. In recent years, building on empirical findings on the repetitiveness of code changes [6], several approaches have built machine learning models based on patch datasets to automate various software engineering tasks such as patch description generation [8, 12, 31, 38, 41, 42, 44, 74], code completion [10, 39, 40, 49, 59], patch correctness assessment [65], and just-in-time defect prediction [27, 32].

Early work manually crafted a set of hand-picked features to represent a patch [32, 33]. Recently, the successful application of deep learning to learn powerful representations of text, signal, and images [14, 36, 47, 50, 52, 60, 71] led to adopting similar techniques in software engineering where researchers develop deep ML models to represent code and code changes [18, 28, 30, 43, 46, 49, 65, 75]. Initially, these approaches treated code [17, 18, 24, 24, 73] and other code-like artifacts, such as patches [15, 42, 46, 74], as a sequence of tokens and thus employ natural language processing methods to extract code in text format. But because source code is not just text, researchers noticed the importance of code structure and began to use the code structure – e.g. through Abstract Syntax Trees (ASTs) – to capture the underlying structural information in source code [3, 4, 24, 76]. However, the differences in code token sequences (usually represented by + and – lines in the textual diff format) are not sufficient to represent the full semantics of the code change as the + and – symbols do not have an inherent meaning that a DL model can learn. Therefore, recent work such as commit2vec [9], C³ [7], and CC2Vec [28] attempted to represent code changes more structurally by leveraging ASTs as well. To get the best of both worlds, more

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

xx, xx, xx, xx

© 2023 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

recent work tried to combine token information with structure information to obtain a better patch representation [15]. Finally, several such approaches of patch representation learning have been evaluated on specific tasks, e.g., BATS [62] for patch correctness assessment and FIRA [15] for patch description generation.

On the one hand, token-based approaches for patch representation [28, 46, 74] lack the rich structural information of source code and intention features inside the sequence is still unexplored. On the other hand, graph-based representation of patches [37, 42] lacks the context which is better represented by the sequence of tokens [28, 46, 74] of the patch itself and also the surrounding unchanged code and intention features inside graph changes is also still unexplored. In conclusion, approaches that try to combine context and AST information to represent patches (e.g., FIRA [15]) do not use the intention features of either sequence or graph from the patch but rather rely on representing the code before and after the change while adding some *ad-hoc* annotations to highlight the changes for the model.

This paper. We propose a novel patch representation that tackles the aforementioned problems and provides an extensive evaluation of our approach on three practical and widely used downstream software engineering tasks. Our approach, Patcherizer, learns to represent patches through a combination of (1) the context around the code change, (2) a novel SeqIntention representation of the sequential patch, and (3) a novel representation of the GraphIntention from the patch. Our approach enables us to leverage powerful DL models for the sequence intention such as Transformers and similarly powerful graph-based models such as GCN for the graph intention. Additionally, our model is pre-trained and hence task agnostic where it can be fine-tuned for many downstream tasks. We provide an extensive evaluation of our model on three popular patch representation tasks: (1) Generating patch description in NL, (2) Predicting patch correctness in program repair, and (3) Patch intention detection.

Overall, this paper makes the following contributions:

- *A novel representation learning approach for patches:* we combine the context surrounding the patch with a novel sequence intention encoder and a new graph intention encoder to represent the intention of code changes in the patch while enabling the underlying neural models to focus on the code change by representing it explicitly. To that end, we developed: ① *an adapted Transformer architecture for code sequence intention* to capture sequence intention in patches taking into account not only the changed lines (added and removed) but also the full context (i.e., the code chunk before the patch application); ② *an embedding approach for graph intention* to compute embeddings of graph intention capturing the semantics of code changes.
- *A dataset of parsable patches:* given that existing datasets only provide patches with incomplete details for readily collecting the code before and after the patch, extracting AST diffs was challenging. We therefore developed tool support to enable such collection and produced a dataset of 90k patches, which can be parsed using the Java compiler.
- *Extensive evaluation:* we evaluate our approach by assessing its performance on several downstream tasks. For each task, we show how Patcherizer outperforms carefully-selected baselines.

We further show that Patcherizer outperforms the state of the art in patch representation learning.

2 PATCHERIZER

Figure 1 presents the overview of Patcherizer. Patches are first preprocessed to split the available information about added (+) and removed (−) lines, identifying the code context (i.e., the code chunk before applying the patch) and computing the ASTs of the code before and after applying the patches (cf. Section 2.1). Then, Patcherizer deploys two encoders, which capture sequence intention semantics (cf. Section 2.2) and graph intention semantics (cf. Section 2.3). Those encoded information are aggregated (cf. Section 2.4) to produce patch embeddings that can be applied to various downstream tasks. In the rest of this section, we will detail the different components of Patcherizer before discussing the pre-training phase (cf. Section 2.5).

2.1 Patch Preprocessing

The preprocessing aims to focus on three main information within a patch for learning its representation. The code before applying the patch (which provides contextual information of the code change), plus and minus lines (which provide information about the code change operations), and the difference between AST graphs before and after patch (which provides information about graph intention in the code). Through the following steps we collect the necessary multi-modal inputs (code text, sequence intention, and graph intention) for the learning:

- (1) **Collect +/− lines in the patch.** We scan each patch line. Those starting with a + are added to a *pluslist*, while those starting with a − are added to a *minuslist*. Both lists record the line numbers in the patch.
- (2) **Reconstruct before/after code.** Besides +/- lines, a patch includes unchanged code that are part of the context. We consider that the full context is the code before applying the patch (i.e., unchanged & minuslist lines). We also construct the code after applying the patch (i.e., unchanged & pluslist lines). The reconstruction leverages the recorded line numbers for inserting each added/removed line to the proper place and ensure accuracy.
- (3) **Generate code ASTs before and after patch.** We apply the Javalang [61] tool to generate the ASTs for the reconstructed code chunks before and after applying the patch.
- (4) **Construct vocabulary.** Based on the code changes of the patches in the training data, we build a vocabulary using the Byte-Pair-Encoding (BPE) algorithm.

At the end of this preprocessing phase, for each given patch, we have a set of inputs:

$\langle cc_p, cc_m, cbp, cap, G_{cbp}, G_{cap} \rangle$, where cc_p is the sequence of added (+) lines of code, cc_m is the sequence of removed (−) lines of code, cbp is the code chunk before the patch is applied, cap is the code chunk after the patch is applied, G_{cbp} is the AST graph of cbp and G_{cap} is the graph of cap .

2.2 Sequence Intention Encoder

A first objective of Patcherizer is to build an encoder that is capable of capturing the semantics of the sequence intention in a patch. While prior work focuses on +/- lines, we postulate that code

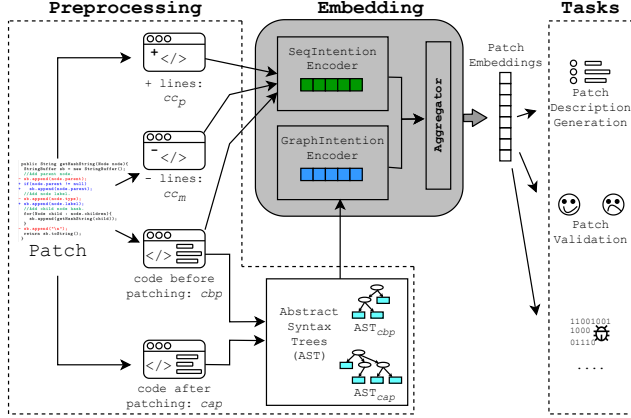


Figure 1: Overview of Patcherizer.

context is a relevant additional input for better encoding such differences. Figure 2 depicts the architecture of the Sequence intention encoder. We leverage the relevant subset of the preprocessed inputs (cf. Section 2.1) to pass to a Transformer embedding layer and further develop a specialized layer, named the *SeqIntention embedding layer*, which captures the intention features from the sequence.

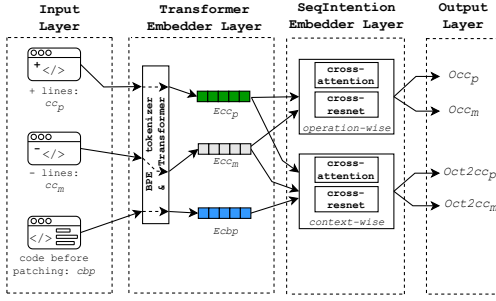


Figure 2: Architecture for the Sequence Intention Encoder.

2.2.1 Input Layer. The input, for each patch, consists of the triplet $\langle cc_m, cc_p, cbp \rangle$, where cc_m is the set of removed (-) lines, cc_p the set of added (+) lines and cbp is the code before patching, which represents the context.

2.2.2 Transformer Embedding Layer. To embed the sequence of code changes, we use a Transformer as the initial embedder. Indeed, Transformers have been designed to capture semantics in long texts and have been demonstrated to be effective for inference tasks [14, 72].

We note that $cc_m \in cbp$. Assuming that $cc_p = \{token_{p,1}, \dots, token_{p,j}\}$, $cc_m = \{token_{m,1}, \dots, token_{m,k}\}$, $cbp = \{token_{cbp,1}, \dots, token_{cbp,l}\}$, where j, k, l represent the maximum length of cc_p , cc_m , and cbp respectively, we use the initial embedding layer the Pytorch implementation of a Transformer to produce first vector representations for each input information as:

$$E_X = \text{Transformer}(\text{Init}(X; \Theta_1); \Theta_2) \quad (1)$$

where X represents an input (either cc_m , cc_p or cbp); Init is the initial embedding function; Transformer is the model based on a transformer architecture; Θ_1 and Θ_2 are the parameters of $\text{Init}(\cdot)$ and $\text{Transformer}(\cdot)$, respectively.

The Transformer embedding layer outputs $E_{cc_p} = [e_{p,1}, e_{p,2}, \dots, e_{p,j}] \in \mathbb{R}^{j \times d_e}$, $E_{cc_m} = [e_{m,1}, e_{m,2}, \dots, e_{m,k}] \in \mathbb{R}^{k \times d_e}$, $E_{cbp} = [e_{cbp,1}, e_{cbp,2}, \dots, e_{cbp,l}] \in \mathbb{R}^{l \times d_e}$, where d_e is the size of the embedding vector.

2.2.3 SeqIntention Embedding Layer. Once the Transformer embedding layer has produced the initial embeddings for the inputs cc_p , cc_m and cbp , our approach seeks to capture how they relate to each other. Prior works [14, 15, 57, 74] have proven that self-attention is effective in capturing relationships among embeddings. We thus propose to capture relationships between the added and removed sequences, with the objective of capturing the intention of the code change through the change operations. We also propose to pay attention to context information when capturing the semantics of the sequence intention.

Operation-wise. To obtain the intention of modifications in patches, we apply a cross-attention mechanism between cc_p and cc_m . To that end, we design a resnet architecture where the model performs residual learning of the importance of inputs (i.e., E_{cc_p} and its evolved v_p , which will be introduced below).

To enhance E_{cc_p} into E_{cc_m} , we apply a **cross-attention** mechanism. For the i -th token in cc_m , we compute the matrix-vector product, $E_{cc_p} e_{m,i}$, where $e_{m,i} \in \mathbb{R}^{d_e}$ is a vector parameter for i -th token i in cc_m . We then pass the resulting vector through a softmax operator, obtaining a distribution over locations in the E_{cc_p} ,

$$\alpha_i = \text{SoftMax}(E_{cc_p} e_{m,i}) \in \mathbb{R}^k, \quad (2)$$

where $\text{SoftMax}(x) = \frac{\exp(x)}{\sum_j \exp(x_j)}$. $\exp(x)$ is the element-wise exponentiation of the vector x . k is the length of cc_m . The attention α is then used to compute vectors for each token in cc_m ,

$$v_i = \sum_{n=1}^j \alpha_{i,n} h_n. \quad (3)$$

where $h_n \in E_{cc_p}$, j is the length of cc_p . In addition, v_i is the new embedding of i -th token in cc_m enhanced by semantic of E_{cc_p} .

Then, we get new cc_m embedding $v_m = [v_1, \dots, v_k] \in \mathbb{R}^{k \times d_e}$.

Similarly, following steps above, we can obtain new embedding of cc_p , $v_p \in \mathbb{R}^{j \times d_e}$, enhanced by the semantic of cc_m .

For the combination of v_p , v_m , E_{cc_p} , E_{cc_m} , inspired by [25, 58], we design a **cross-resnet** for combining v_p , v_m , v_{cc_p} , and v_{cc_m} . The pipeline of **cross-resnet** is shown in Figure 3. The process is as follows:

$$\begin{aligned} O_{cc_p} &= f(h(E_{cc_p}) + \mathcal{A}dd(E_{cc_p}, v_p)) \\ O_{cc_m} &= f(h(E_{cc_m}) + \mathcal{A}dd(E_{cc_m}, v_m)) \end{aligned} \quad (4)$$

where $h(\cdot)$ is a normalization function in [14]; $\mathcal{A}dd(\cdot)$ is the adding function, $f(\cdot)$ represents $RELU$ [23] activation function.

Finally, we obtain output O_{cc_m} and O_{cc_p} .

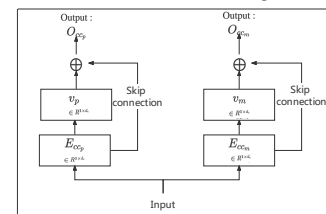


Figure 3: cross-resnet architecture.

Context-Wise. Similar to operation-wise block, we enhance the contextual information into modified lines by **cross-attention** and **cross-resnet** blocks. The computation process is as follows:

$$\begin{aligned} O_{ct2cc_p} &= f(h(E_{cc_p}) + \mathcal{A}dd(E_{cc_p}, E_{cbp})) \\ O_{ct2cc_m} &= f(h(E_{cc_m}) + \mathcal{A}dd(E_{cc_m}, E_{cbp})) \end{aligned} \quad (5)$$

where E_{cbp} is the embedding of cbp calculated by Equation 1; O_{ct2cc_p} represents the vector of context-enhanced plus embedding and O_{ct2cc_m} is the vector of context-enhanced minus embedding.

2.3 Graph Intention Encoder

Concurrently to encoding sequence information from code changes, we propose to also capture the graph intention features of the structural changes in the code when the patch is applied. To that end, we rely on a Graph Convolutional Network (GCN) architecture, which is widely used to capture dynamics in social networks, and is effective for typical graph-related tasks such as classification or knowledge injection [35, 77, 78]. Once the GCN encodes the graph nodes, the produced embeddings can be used to assess their relationship via computing their cosine similarity scores [45]. Concretely, in Patcherizer, we use a GCN-based model to capture the graph intention features. The embedder model was trained by inputting a static graph, a graph resulting from the merge of all sub-graphs from the training set. Overall, we implement this encoding phase in two steps: building the static graph, performing graph learning and encoding the graph intention (cf. Figure 4).

2.3.1 Graph Building. To start, we consider the G_{cbp} and G_{cap} trees, which represent in graph forms.

Static Graph building: Each patch in the dataset can be associated to two graphs: G_{cbp} and G_{cap} , which are obtained by parsing the cbp and cap code snippets. After collecting all graphs (which are unidirectional graphs) for the whole training set, we merge them into a “big” graph by iteratively linking the common nodes. In this big graph, each distinct code snippet AST-inferred graph is placed as a distinct sub-graph. Then, we will merge the graphs shown in Figure 5 which illustrates the merging progress of two graphs: if a node n has the same value, position, and neighbors in both ASTs, it will be merged into one (e.g., red nodes 1 and 2). However, when a common node has different neighbors between the ASTs (e.g., red nodes 3 and 4), the merge keeps one instance of the common node but includes all neighbors connected to the merged red nodes (i.e., all green and grey nodes are now connected to red nodes 3 and 4, respectively). After iterating over all graphs, we eventually build the static graph.

However, on the one hand, some nodes in most subgraphs such as ‘prefix_operators’, ‘returnStatement’, and ‘StatementExpression’ are not related to the semantics of the patch. On the other hand, as data statistics in our study, 97.2% nodes in the initial graphs (ASTs) extracted from code are from parser tools instead of patches, which means these nodes are rarely related with the semantic of the patch. Thus, as shown in step 1 in Figure 4, we remove nodes whose children do not contain words in the patch to reduce the size of the graph because these nodes will be considered noise in our research.

In the remainder of this paper, we refer to the final graph as the *static graph* $\mathcal{G}=(\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the set of nodes and \mathcal{E} is the set of edges.

Graph Alignment to the Static Graph: GCN requires that the input graphs are all of the same size [35]. Yet, the graphs built using the graphs of cbp and cap do not have as many nodes and edges as the static graph used for training the GCN network. Consequently, we propose to use the static graph to initialize all graphs. Given the global static graph $\mathcal{G}_{global} = (\mathcal{V}_g, \mathcal{E}_g)$ and an AST graph $\mathcal{G}_{local} = (\mathcal{V}_l, \mathcal{E}_l)$, we leverage the VF graph matching algorithm [11] to find the most similar sub-graph with \mathcal{G}_{local} in \mathcal{G}_{global} :

$$subGraph = VFG(\mathcal{G}_{local}, \mathcal{G}_{global}). \quad (6)$$

where $VFG(\cdot)$ is the function representing the VF matching algorithm [1]. The matched sub-graph is a subset of both \mathcal{G}_{local} and \mathcal{G}_{global} : some nodes of $subGraph$ will be in \mathcal{G}_{local} but not in \mathcal{G}_{global} . We then align \mathcal{G}_{local} to the same size of \mathcal{G}_{global} as follows: we use the [PAD] element to pad the node of $subGraph$ to the same size of \mathcal{G}_{global} and then we obtain \mathcal{G}_{lPAD} . Therefore, \mathcal{G}_{lPAD} keeps the same size and structure of the static graph \mathcal{G}_{global} . Eventually, all graphs are aligned to the same size of \mathcal{G}_{global} and the approach can meet the requirements for GCN computation for graph learning.

2.3.2 Graph Learning. Inspired by [77], we build a deep graph convolutional network based on the undirected graph formed following the above construction steps to further encode the contextual dependencies in the graph. Specifically, for a given undirected graph $\mathcal{G}_{lPAD} = (\mathcal{V}_{lPAD}, \mathcal{E}_{lPAD})$, let \mathcal{P} be the renormalized graph laplacian matrix [35] of \mathcal{G}_{lPAD} :

$$\begin{aligned} \mathcal{P} &= \hat{\mathcal{D}}^{-1/2} \hat{\mathcal{A}} \hat{\mathcal{D}}^{-1/2} \\ &= (\mathcal{D} + \mathcal{L})^{-1/2} (\mathcal{A} + \mathcal{L}) (\mathcal{D} + \mathcal{L})^{-1/2} \end{aligned} \quad (7)$$

where \mathcal{A} denotes the adjacency matrix, \mathcal{D} denotes the diagonal degree matrix of the graph \mathcal{G}_{lPAD} , and \mathcal{L} denotes the identity matrix. The iteration of GCN through its different layers is formulated as:

$$\mathcal{H}^{(l+1)} = \sigma(((1-\alpha)\mathcal{P}\mathcal{H}^{(l)} + \alpha\mathcal{H}^{(0)})(\mathcal{L} + \beta^{(l)}\mathcal{W}^{(l)})) \quad (8)$$

where α and $\beta^{(l)}$ are two hyper parameters, σ denotes the activation function and $\mathcal{W}^{(l)}$ is a learnable weight matrix. Following GCN learning, we use the average embedding of the graph to represent the semantic of structural information in code snippet:

$$w_G = \frac{1}{L} \sum_{i=1}^L (\mathcal{H}_i). \quad (9)$$

Thus, at the end of the graph embedding, we obtain representations for G_{cbp} and G_{cap} , i.e., $w_{G_{cbp}}$ and $w_{G_{cap}} \in \mathbf{R}^{1 \times d_e}$.

2.3.3 Graph Intention Encoding. Once we have computed the embeddings of the code snippets before and after patching, (i.e., the embeddings of cbp and cap), we must get the representation of their differences to encode the intention inside the graph changes. To that end, similarly to the previous cross-resnet for sequence intention, we design a **graph-cross-resnet** operator which ensembles the semantic of $w_{G_{cbp}}$ and $w_{G_{cap}}$. Figure 6 illustrates this crossing. In this graph-cross-resnet, the model can choose and highlight a path automatically by the backpropagation mechanism. The

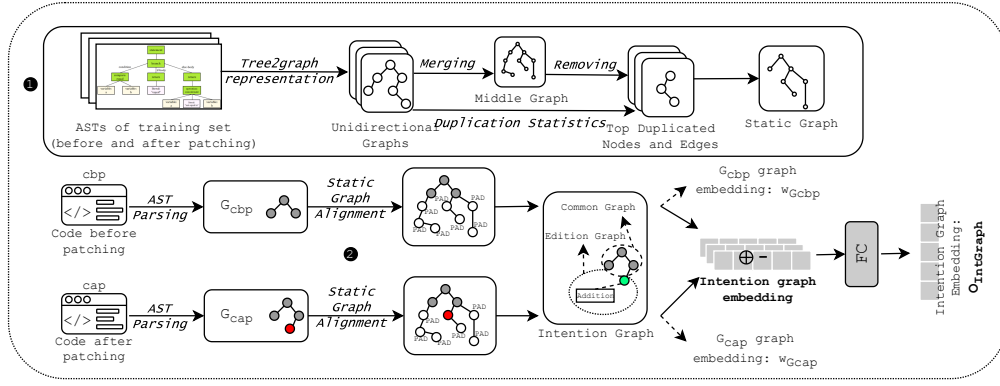


Figure 4: Architecture for the Graph Intention Encoder.

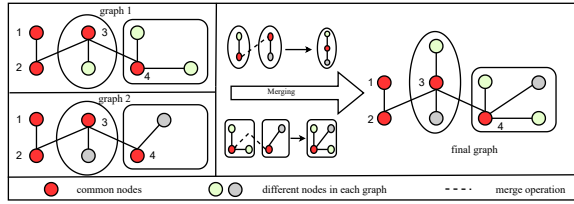


Figure 5: An example of merging graphs.

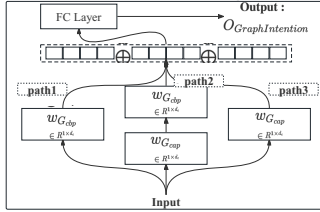


Figure 6: graph-cross-resnet architecture.

$GraphIntention$ is therefore calculated as follows:

$$\begin{aligned}
 path1 &= w_{G_{cbp}} \\
 path2 &= \mathcal{Add}(w_{G_{cbp}}, w_{G_{cap}}) \\
 path3 &= w_{G_{cap}} \\
 O_{GraphIntention} &= \mathcal{FC}(f(\mathcal{Add}(path1, path2, path3)), \Theta_3)
 \end{aligned} \tag{10}$$

where \mathcal{FC} is a fully-connected layer and Θ_3 is the parameter of \mathcal{FC} .

At the end, the graph-cross-resnet component outputs the sought graph intention embedding: $GraphIntention$.

2.4 Aggregating Multimodal Input Embeddings

With the sequence intention encoder and the graph intention encoder, we can produce for each patch several embeddings of different input modalities (code sequences and graphs) that must be aggregated into a single representation.

Concretely, we use the \mathcal{Add} aggregation function to merge the $SeqIntention$ embeddings (combination of O_{ccp} , O_{ccm} , O_{ct2ccp} and O_{ct2ccm} - cf. Equations 4 and 5) and $GraphIntention$ embedding $O_{GraphIntention}$ before outputting the final representation $E_{Patcherizer}$. Actually, we use $E_{Patcherizer}$ as the representation of patch out of the model.

2.5 Pre-training

Patcherizer is an approach that is agnostic of the downstream tasks. We propose to build a pre-trained model with the collected dataset of patches. The objective is to make the model learn contextual semantics of code, which can help improve the efficiency of representing patches.

The pre-training task is masked token prediction. Indeed, a popular bidirectional objective function is driven by the masked language model (MLM; [14]), which aims to predict masked word based on its surrounding context, i.e., compute the probability $P(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$.

Inspired by previous prior works in model pre-training [14, 17, 18, 79], we only employ predicting [MASK]s as our pre-training strategy. We randomly mask some words in the source code. The target is then to use the contextual information to predict the masked word as accurately as possible.

We take the pre-training as a translation task and employ the popular encoder-decoder architecture as our pipeline. We use our Patcherizer architecture as our encoder and the BERT Transformer model [14] as the decoder, where they share the same parameters. Apart from that, encoder and decoder share vocab. Finally, we start the decoding process with an initial $\langle s \rangle$ token to generate the code sequence word by word:

$$index = \operatorname{argmax}(p(y_t | y_{t-1}, \dots, y_1, E_{Patcherizer})) \tag{11}$$

where $y_i (i = 1, \dots, t)$ is a one-dimension distribution where the distribution size equals the length of vocab. The index is the number of the element which is the max one in y_i . Then we can return the corresponding word with the same index in vocab. Furthermore, we employ the Cross-Entropy [54] as the loss function and apply Adam algorithm [34] as the optimizer.

2.6 Fine-tuning for Different Tasks

Patcherizer works as an embedder for patches. It is pre-trained by unsupervised learning contextual semantics with collected datasets of patches. Thus, given a patch, Patcherizer is able to generate its task-agnostic embedding.

We can now use the generated embeddings for various downstream tasks. We now describe how patch representations produced by Patcherizer are adopted to address the three popular patch-related downstream tasks that we investigate in this work: patch

description generation [15, 28, 46, 74], patch correctness assessment [65], Just In Time (JIT) defect prediction [27, 28].

2.6.1 Patch Description Generation. This task aims to generate a natural-language description for a given patch. The patch description generation task is highly similar to the pre-training task (cf. 2.5). The training dataset of this task is bimodal. It includes patches associated to descriptions. We input the patch and fine-tune the model to generate pseudo descriptions that are as much as possible similar to the ground-truth descriptions. The fine-tuning process Patcherizer is therefore formulated as for pre-training masked word prediction of Equation 11.

2.6.2 Patch Correctness Assessment. This task is relevant in the program repair community where generated patches must be predicted as correct or not for a given bug. It is a binary classification task. Given the patch embedding produced by Patcherizer and a bug report embedding produced by BERT [14], we concatenate them and apply a fully-connected layer to classify it as correct or not:

$$\hat{y}_i = \text{sigmoid}(E_{\text{patch}_i} \oplus E_{\text{bugReport}_i}) \quad (12)$$

where E_{patch_i} is the Patcherizer’s embedding of patch i , $E_{\text{bugReport}_i}$ is the embedding of the associated bug report, and \oplus is the concatenation operator.

To optimize our model for the task of patch correctness assessment, we opt for a categorical cross-entropy loss function as follows:

$$\mathcal{L}_a = - \sum_{i=1}^n (y_i \times (\ln(\hat{y}_i)) + (1 - y_i) \times \ln(1 - \hat{y}_i)) \quad (13)$$

where \mathcal{L}_a is loss value, n is the size of the dataset, y_i is the ground-truth label, \hat{y}_i is the predicted label.

3 EXPERIMENTAL DESIGN

We provide the implementation details (cf. Sec. 3.1), discuss the research questions (cf. Sec. 3.2), and present the baselines (cf. Sec. 3.3), the datasets (cf. Sec. 3.4), and the metrics (cf. Sec. 3.5).

3.1 Implementation

In the pre-training phase used for the Sequence Intention Embedding step, we apply a beam search [67] for the best performance in predicting the masked words. The beam size was set to 3. The dimension of the hidden layer output in models is set to 512, and the default value of dropout rate is set to 0.1. For the Transformer, we apply 6 heads for the multi-header attention module and 4 layers for the attention.

For the Graph Intention Embedding step, we use javalang [61] to parse code fragments and collect ASTs. We build on graph manipulation packages (i.e., networkx [1], and dgl [70]) to represent these ASTs into graphs.

Patcherizer’s training involves the Adam optimizer [34] with learning rate 0.001. All model parameters are initialized using Xavier algorithm [22]. All experiments are performed on a server with Intel(R) Xeon(R) CPU E5-2698 v4 @ 2.20GHz, 256GB physical memory, and one NVIDIA Tesla V100 GPU with 32GB memory.

3.2 Research questions

RQ-1: *How effective is Patcherizer in learning patch representations?*

RQ-2: *What is the impact of the key design choices on the performance*

of Patcherizer?

RQ-3: *To what extent is Patcherizer effective on independent datasets?*

3.3 Baselines

We consider several SOTA models as baselines. We targeted approaches that were specifically designed for patch representation learning (e.g., CC2Vec) as well as generic techniques (e.g., NMT) that were already applied to patch-related downstream tasks. We finally consider recent SOTA for patch-representation approaches (e.g., FIRA) for specific downstream tasks.

- **NMT** technique has been leveraged by Jiang *et al.* [31] for translating code commits into commit messages.
- **NNGen** [44] is an IR-based commit message prediction technique.
- **CoDiSum** [74] is an encoder-decoder based model with multi-layer bidirectional GRU and copying mechanism [55].
- **CC2Vec** [28] learns a representation of code changes guided by commit messages. It is the incubant state of the art that we aim to outperform on all tasks.
- **CoRec** [68] is a retrieval-based context-aware encoder-decoder model for commit message generation.
- **Coregen** [46] is a **pure Transformer-based** approach for representation learning targeting commit message generation.
- **FIRA** [15] is a graph-based code change representation learning approach for commit message generation.
- **BERT** [14] is a state of the art unsupervised learning based Transformer model widely used for text processing.
- **ATOM** [42] is a commit message generation techniques, which builds on abstract syntax tree and hybrid ranking.
- **CCRep** [43] is an innovative approach that uses pre-trained models to encode code changes into feature vectors, enhancing performance in tasks like commit message generation, etc.

3.4 Datasets

Patch description generation: We build on prior benchmarks [15, 16, 28, 44] by focusing on Java samples and reconstructing snippets to make them parsable for AST collection. Eventually, our dataset includes 90,661 patches and their associated descriptions.

Patch correctness assessment: We leverage the largest dataset in the literature to date, which includes deduplicated 11,352 patches (9,092 Incorrect and 2,260 Correct) released by Tian *et al.* [65].

3.5 Metrics

ROUGE [53] is one set of metrics for comparing automatic generated text against the reference (human-produced) ones. We focus on ROUGE-L which computes the Longest Common Subsequence. **BLEU** [48] is a classical metric to evaluate the quality of machine translations. It measures how many word sequences from the reference text occur in the generated text and uses a (slightly modified) n-gram precision to generate a score.

METEOR [5] is an F-Score-Oriented metric for measuring the performance of text-generation models.

+Recall and **-Recall** [65] are specific metrics for patch correctness assessment. **+Recall** (**-Recall**) measures to what extent correct (incorrect) patches can be predicted (filtered out).

4 EXPERIMENTAL RESULTS

4.1 [RQ-1]: Performance of Patcherizer

[Experiment Goal]: We assess the effectiveness of the embeddings learned by Patcherizer on three popular and widely used software engineering tasks: (RQ-1.1) Patch description generation, (RQ-1.2) Patch correctness assessment, and (RQ-1.3) Patch intention detection. We compare Patcherizer against the relevant SOTA.

[Experiment Design (RQ-1.1)]: We employ the dataset from FIRA. As Xu et al. [15] have previously assessed FIRA and other baseline methods using this dataset, we directly reference the evaluation results of all the baselines from Table IV of the FIRA paper. The dataset contains 75K, 8K and 7.6K commit-message pairs in the training, validation and test sets, respectively. We evaluate the generated patch descriptions in the test set using the BLEU, ROUGE-L, and METEOR metrics.

Note that we distinguish between baseline generation-based methods and retrieval-based ones. In generation-based baselines, a patch description is actually synthesized, while in retrieval-based baselines, the approach selects a description text from an existing corpus (e.g., in the training set). For fairness, we build two distinct methods using Patcherizer’s embeddings. The first method is generative and follows the fine-tuning process described in Section 2.6. The second method is an IR-based approach, where, following the prior work [28], we use Patcherizer as the initial embedding tool and implement a retrieval-based approach to identifying a relevant description in the training set: the description associated with the training set patch that has the highest similarity score with the test set patch is outputted as the "retrieved" description.

[Experiment Results (RQ-1.1)]: Table 1 presents the average scores of the different metrics with the descriptions generated by Patcherizer and the relevant baselines. Patcherizer outperforms all the compared techniques on all metrics, with the exception of FIRA on the ROUGE-L metric. The superior performance of Patcherizer on generation-based and retrieval-based methods, as illustrated by the distribution of BLEU scores in Figure 7, further suggest that the produced embeddings are indeed effective.

Table 1: Performance Results of patch description generation.

Type	Approach	Rouge-L (%)	BLEU (%)	METEOR (%)
Generation	NMT [31]	7.35	8.01	7.93
	Codisum [74]	19.73	16.55	12.83
	ATOM [42]	10.17	8.35	8.73
	FIRA [15]	21.58	17.67	14.93
	CoreGen [46] (Transformer)	18.22	14.15	12.90
	CCRep [43]	23.41	19.70	15.84
	Patcherizer	25.45	23.52	21.23
Retrieval	CC2Vec [28]	12.21	12.25	11.21
	NNGen [44]	9.16	9.53	16.56
	CoRec [68]	15.47	13.03	12.04
	Patcherizer	17.32	15.21	17.25

"Generation" for generation-based strategy. Given fragments of codes, "Generation" methods generate messages from scratch.

"Retrieval" for retrieval-based approaches. Given fragments of codes, "Retrieval" approaches return the most similar message from the training dataset.

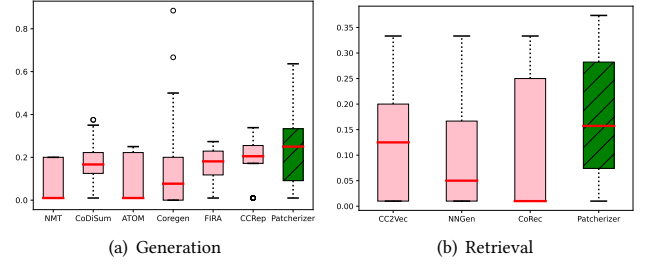
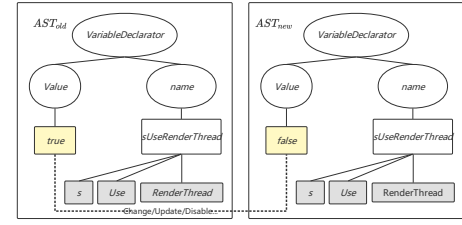


Figure 7: Comparison of the distributions of BLEU scores for different approaches in patch description generation

```

--- s0.java 2022-08-20 17:00:26.000000000 +0200
+++ t0.java 2022-08-20 17:01:04.000000000 +0200
@@ -1,6 +1,6 @@
 public abstract class HardwareRenderer {
     public static boolean sSystemRendererDisabled = false ;
-    public static boolean sUseRenderThread = true ;
+    public static boolean sUseRenderThread = false ;
     private boolean mEnabled ;
     private boolean mRequested = true ;
 }

```



Source	Patch description
Ground truth	Disable RenderThread
CC2Vec	Fix test data so that it can be compiled
FIRA	fix the in
CCRep	update suserenderthread
Patcherizer	disable renderthread

Figure 8: Illustrative example of patch description generation

In Figure 8, we provide an example result of generated description by Patcherizer, by the CC2Vec strong baseline (using retrieval-based method) and by the FIRA and CCRep state-of-the-art approach (using generation-based method) for patch description generation. Patcherizer succeeds in actually generating the exact description as the ground truth commit message, after taking into account both sequential and structural information. By observing the graph intention and sequence intention, we can see that the model found that the only change is that the node `true` has been changed/updated/disabled to `false`. Finally, the sequence intention embedding would make Patcherizer recognize that the carrier of `true` and `false` is `RenderThread` based on BPE splitting.

🔗 **Answer to RQ-1.1:** ▶ Patcherizer’s embeddings are effective for patch description generation yielding the best scores for BLEU, ROUGE-L, and METEOR metrics. ◀

[Experiment design (RQ-1.2)]: Tian et al. [63] proposed to leverage the representation learning (embeddings) of the patches to assess patch correctness. Following up on their study, we use the patch embeddings produced by CC2Vec, BERT, CCRep, and Patcherizer (cf. Section 2.6) to train three classifiers to classify APR-generated patches as correct or not and we experiment with two supervised learning algorithms: Logistic regression (LR) and XGBoost (XGB).

To perform a realistic evaluation, we split the patches dataset by bug-id into 10 groups to perform a 10-fold-cross-validation experiment similar to previous work [65]. In this splitting strategy, all patches for the same bug are either placed in the training set or the testing set to ensure that there is no data leakage between the training and testing data. We then measure the performance of the classifiers using +Recall, -Recall, AUC, and F1.

[Experiment Results (RQ-1.2)]: Table 2 shows the results of this experiment. Both classifiers, LR and XGB, when trained with Patcherizer embeddings largely outperform the classifiers that are trained with BERT or CC2Vec embeddings, which achieved SOTA results in literature [63].

Table 2: Performance of Patch Correctness Classification

Classifier	Model	AUC	F1	+Recall	-Recall
LR	CC2Vec	0.75	0.49	0.47	0.85
	BERT	0.83	0.58	0.81	0.65
	CCRep	0.86	0.67	0.74	0.83
	Patcherizer	0.96	0.82	0.87	0.91
XGB	CC2Vec	0.81	0.55	0.50	0.89
	BERT	0.84	0.61	0.64	0.85
	CCRep	0.82	0.63	0.59	0.88
	Patcherizer	0.90	0.67	0.66	0.90

✎ **Answer to RQ-1.2:** Patch embeddings generated by Patcherizer achieve SOTA results in the task of patch-correctness assessment, largely outperforming SOTA embedding models. ◀

[Experiment Goal (RQ-1.3)]: Previous work introduces that the patch has its intention and detecting the intention of the patch can help the model understand the semantics of the patch (i.e., template-based works [8, 12] and generation-based works [15, 74]). Thus, efficiency of patch intention detection can be used to measure if the patch representation model is good or not.

[Experiment Results (RQ-1.3)]:

We scan all words across three datasets in our work and figure out that patches are mainly related to four types: `fix`, `remove`, `add`, and `update`. However, `fix` is highly related to all other three frequent words, because `fix` can be used to update, remove or add. Therefore, we select `add`, `remove`, `update` as our main detected intentions. In this section, we aim to explore how Patcherizer performs against the representative models CC2Vec and CCRep on distinguishing the intention of patches.

We trained the three models (i.e., Patcherizer, CC2Vec, and CCRep) on a large dataset proposed in [15]. Then, we assess the patch intention detection ability of these models on the CC2Vec dataset [28]. We find that 572 patches contain `add`, `remove`, or `update` keywords (i.e., 201 for `add`, 341 for `remove`, 30 for `update`). Then, we use the three models to embed these 572 patches and obtain corresponding high-dimensional vectors. We employ `t-SNE` [66] to reduce the dimensionality for better visualization.

Figure 9 shows the `t-SNE` visualized results of CC2Vec, CCRep and Patcherizer. The red color represents `add` function, the green color represents `remove` function, and the blue color represents `update` function. We see that Patcherizer separates `add` and `remove` better than CC2Vec and CCRep. Furthermore, both CC2Vec and CCRep fail to separate `update` from the other two functions. The reason may be that `update` functions can be `add` or `remove` functions. Thus, the patch semantic distribution from both CC2Vec and CCRep is mixed with `add` and `update`.

✎ **Answer to RQ-1.3:** Compared with existing patch representation models, Patcherizer is more effective in detecting the intention of patches. ◀

4.2 [RQ-2]: Ablation Study

[Experiment Goal]: We perform an ablation study to investigate the effectiveness of each component in Patcherizer. The major novelty of Patcherizer is the fact that it explicitly includes and processes: ① *SeqIntention* represents intention embedding of the patch at the sequential level, and ② *GraphIntention* represents intention embedding of the patch at the structural level.

[Experiment Design]: We investigate the related contribution of *SeqIntention* and *GraphIntention* by building two variants of Patcherizer where we remove either *GraphIntention* (i.e., denoted as Patcherizer *GraphIntention*-), or *SeqIntention* (i.e., denoted as Patcherizer *SeqIntention*-). We also build a native model by removing both *GraphIntention* and *SeqIntention* components (i.e., denoted as Patcherizer *both*- for comparison. We evaluate the performance of these variants on the task of patch description generation.

[Experiment Results (RQ-2)]: Table 3 summarizes the results of our ablation test on the three variants of Patcherizer.

Table 3: Ablation study results based on the patch description generation task.

Model	ROUGE-L (%)	BLEU (%)	METEOR (%)
Patcherizer <i>GraphIntention</i> -	20.10	16.50	15.40
Patcherizer <i>SeqIntention</i> -	18.44	14.70	16.20
Patcherizer <i>both</i> -	15.00	13.00	12.00
Patcherizer	25.45	23.52	21.23

While the performance of Patcherizer is not the simple addition of the performance of each variant, we note with Patcherizer *both*- that the performance is quasi-insignificant, which means that, put together, both design choices are instrumental for the superior performance of Patcherizer.

Contribution of Graph Intention Encoding: We observe that the graph intention embedding significantly improves the model ability to generate correct patch descriptions for more patches which is evidenced by the large improvement on the ROUGE-L score (from 20.10 to 25.45), where ROUGE-L is recall oriented.

We postulate that even when token sequences (e.g., identifier names) are different among patches, the similarity of the intention graph helps the model to learn that these patches have the same intent. Nevertheless, precision in description generation (i.e., how many words are correct) is highly dependent on the model's ability to generate the exact correct tokens, which is more guaranteed by the context and sequence intention embedding.

We manually checked different samples to analyze how the variants were performing. Figure 10 presents a real-world case in our dataset, including the patch, the ground truth, and the patch descriptions generated by Patcherizer, Patcherizer *GraphIntention*-, Patcherizer *SeqIntention*-, as well as three of the strongest baselines for this task (i.e., CC2Vec, FIRA, and CCRep). In this case, the embeddings of Patcherizer and Patcherizer *GraphIntention*- are effective in spotting the sub-token *trident* in class name *TridentTopologyBuilder* thanks to BPE. In addition, Patcherizer takes advantage of both the sequence intention and graph intention inside the patch. However, if we only consider the graph intention, Patcherizer *SeqIntention*-

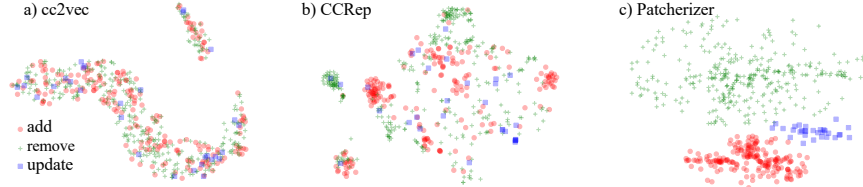


Figure 9: Visualization of Patch intention recognition by different models.

performs the worst against Patcherizer *GraphIntention*-. From the example, we find that CC2Vec, which is retrieval-based, cannot generate a proper message because there may not exist similar patches in the training set. FIRA, while underperforming against Patcherizer, still performs relatively well because it uses the edition operation detector and sequential contextual information.

```
--- a/src/[...]/topology/TridentTopologyBuilder.java
+++ b/src/[...]/topology/TridentTopologyBuilder.java
public class TridentTopologyBuilder {
    bd.allGrouping( masterCoordinator( batchGroup ),
        MasterBatchCoordinator.
        COMMIT_STREAM_ID);
    for (Map m : c.componentConfs) {
        -   scd . addConfigurations ( m );
        +   bd . addConfigurations ( m );
    }
}
```

Source	Patch description
Ground truth	set component configurations correctly for trident spouts
Patcherizer	set component configurations correctly for trident spouts
Patcherizer <i>GraphIntention</i> -	configure components for trident
Patcherizer <i>SeqIntention</i> -	set trident components.
CC2Vec	fixed flickering in the preview pane in refactoring preview
FIRA	use the correct component content in onesidediffviewer
CCRep	update for function

Figure 10: Case analysis of the ablation study

It is noteworthy that Patcherizer is able to generate the token *spouts*. This is not due to data leakage since the ground truth commit message was not part of the training set. However, our approach builds on a dictionary that considers all tokens in the dataset (just as the entire English dictionary would be considered in text generation). Hence *spouts* was predicted from the dictionary as the most probable (using *softmax*) token to generate after *trident*.

✎ **Answer to RQ-2:** ► *Evaluations of individual components of Patcherizer indicate that both GraphIntention and SeqIntention bring a significant portion of the performance.* ◀

4.3 [RQ-3]: Robustness Evaluation

[Experiment Design]: We evaluate the robustness of Patcherizer and of state-of-the-art patch representation techniques (i.e., CC2Vec [28] and CCRep [43]) on an independent dataset for the task of patch description generation. To that end, we pre-train Patcherizer on the dataset used for patch description generation task 3.4, but we use, for testing, the dataset of patches that was collected for patch correctness assessment. The test dataset is consequently independent as the patches come from a different set of projects. We indeed confirm that the samples across the two datasets are substantially different: Table 4 depicts some basic statistics, which reveal that the average patch sizes and commit message sizes are different. A comparison of patch examples shows that patches in the patch description generation task dataset are small (mostly few-line changes as in Figure 10), while the patches in the independent dataset are larger (cf. Figure 11). Conversely, the description messages are an

```
--- a/src/[...]/math/analysis/solvers/BaseSecantSolver.java
+++ b/src/[...]/math/analysis/solvers/BaseSecantSolver.java
@@ -183,14 +183,7 @@ public abstract class BaseSecantSolver
    f0 *= f1 / (f1 + fx); break;
    case REGULA_FALSI:
    -   if (x == xl) {
    -       final double delta = FastMath.max(...),
    +// Nothing.
    break;
    ...
```

Source	Patch description
Ground truth	Remove an obsolete code line
CC2Vec [28]	wrap the root cause rather than just using the message
CCRep [43]	add notes
Patcherizer	remove obsolete code

Figure 11: Case analysis: robustness evaluation

order of magnitude shorter in the independent dataset than in the patch description task dataset.

Table 4: Dataset statistics

	patch description generation	patch correctness validation
Avg. size of patches (#tokens)	43.70	217.98
Avg. size of descriptions (#tokens)	6.97	-

Following the IR-based techniques to retrieve messages for given patches, we use the generated embeddings to match highly similar test set patches with training set patches. Then, we select the descriptions associated with the matched training patches as the "generated" (actually retrieved) descriptions.

[Experiment Results (RQ-3)]: Table 5 summarizes the experimental results of Patcherizer, CC2Vec, and CCRep baselines on the independent dataset. Compared to CCRep, Patcherizer improves the score from 19.65% to 21.64%, 23.67% to 31.92%, 12.77% to 15.37% for metrics BLEU, ROUGE-L and METEOR, respectively.

Table 5: Results on independent dataset (%)

Model	ROUGE-L	BLEU	METEOR
CC2Vec [28]	17.34	9.20	5.14
CCRep [43]	23.67	19.65	12.77
Patcherizer	31.92	21.64	15.37

Figure 11 illustrates a case where Patcherizer performs better than both CC2Vec and CCRep. Patcherizer indeed manages to match and return a correct message from the retrieval target dataset. In contrast, CC2Vec fails to retrieve the precise patch description and CCRep obtains inaccurate query back. Such observations further confirm our intuition that *SeqIntention* and *GraphIntention* can help the model capture the semantics of the patch independently of the training task dataset.

📌 **Answer to RQ-3:** ► Compared to CC2Vec and CCRep on the task of patch description generation on an independent dataset, Patcherizer shows its capability to capture the semantics of patches. Patcherizer achieves 10.13%, 34.85% and 20.36% improvement on CCRep for ROUGE-L, BLEU and METEOR metrics, respectively. ◀

5 DISCUSSION

5.1 Threats to Validity

Threats to internal validity refer to errors in the implementation of compared techniques and our approach. To reduce these threats, in each task, we directly reuse the implementation of the baselines from their reproducible packages whenever available. Otherwise, we re-implement the techniques strictly following their papers. Furthermore, we also build our approach based on existing mature tools/libraries, such as javalang [61] for parsing ASTs.

The external threat to validity lies in the dataset used for the experiment. To mitigate this threat, we build a well-established dataset, which is a rewritten version based on datasets from prior works [28, 46, 65].

The construct threat involves the metrics used in evaluations. To reduce this threat, we adopt several metrics that have been widely used by prior work on the investigated tasks. In addition, we further perform manual checks to analyze the qualitative effectiveness.

5.2 Limitations

First, since Patcherizer relies on *SeqIntention* and *GraphIntention*, our approach would be less effective when patches cannot be parsed into valid AST graph. In this case, Patcherizer would only take contextual information and *SeqIntention* as sources to yield the embeddings. However, this limitation lies only when we cannot access source code repositories in which patches have been committed.

Second, for the patch description generation task, we consider two variants: generation-based and retrieval-based. Normally, we collect datasets by following fixed rules, which leads to the training set containing highly-similar patches with the test set. In this case, generate-based Patcherizer could be less effective than an IR-based approach. Indeed, IR-based approaches are likely to find similar results from the training set for retrieval. Nevertheless, as shown in Table 1, even in retrieval-based mode, Patcherizer outperforms the baselines.

Third, when a given patch contains tokens that are absent from both vocabularies of patches and messages, Patcherizer will fail to generate or recognize these tokens for all tasks.

6 RELATED WORK

6.1 Patch Representation

There are many studies on the representation of code-like texts, including source code representation [17, 18] and patch representation [28]. Previous works focus on representing given patches into latent distributed vectors. Allamanis *et al.* [2] propose a comprehensive survey on representation learning of code-like texts.

The existing works on representing code-like texts can be categorized as control-flow graph [13], and deep-learning approaches [17, 18, 28]. Before learning distributed representations, Henkel *et al.* [26]

proposes a toolchain to produce abstracted intra-procedural symbolic traces for learning word representations. They conducted their experiments on a downstream task to find and repair bugs in incorrect codes. Wang *et al.* [69] learns embeddings of code-like text by the usage of execution traces. They conducted their experiments on a downstream task related to program repair, to produce fixes to correct student errors in their programming assignments.

To leverage deep learning models, Hoang *et al.* [28] proposed CC2Vec, a sequence learning-based approach to represent patches and conduct experiments on three downstream patch tasks: patch description generation, bug fixing patch identification, and just-in-time defect prediction. Similarly, CoDiSum [74] is also a token based approach for patch representation that has been used for generating patch descriptions. CCRep [43] is an approach to learning code change representations, encoding code changes into feature vectors for a variety of tasks by utilizing pre-trained code models, contextually embedding code, and employing a mechanism called "query back" to extract, encode, and interact with changed code fragments. Our work improves on these approaches by leveraging the context around the code change and a novel graph intention embedding. CACHE [37] uses AST embeddings for the code change and its surrounding context to learn patch representation for patch correctness prediction.

The closest to our work is FIRA [15] for learning patch descriptions. It uses a special kind of graph that combines the two ASTs before and after the patch with extra special nodes to highlight the relationship (e.g., match, add, delete) between the nodes from the two ASTs. Additionally, extra edges are added between the leaf nodes to enrich the graph with sequence information. Our work is different in many aspects. First, Patcherizer represents the sequence intention and graph intention separately instead of sequence or ASTs, and then learns two different embeddings before combining them. Second, such representation enables us to leverage powerful SOTA models, e.g., Transformer for sequence learning and GCN for graph-based learning. Third, our GraphIntention representation focuses on learning an embedding of intention of graph changes between AST graphs before and after patching, and not the entire AST which enables the neural model to focus on learning the structural changes. Finally, our approach is task-agnostic and can easily be fine-tuned for any patch-related downstream tasks. We have evaluated it on three different tasks while FIRA was only assessed on patch description generation.

6.2 Applications of Patch Embeddings

Patch description generation: As found by prior studies [15, 16], about 14% commit messages in 23K java projects are empty. Yet patch description is very significant to developers as they help to quickly understand the intention of the patch without requiring reviewing the entire code. Techniques for patch description generation can be categorized as template-based, information-retrieval-based (IR-based), and generation-based approaches. Template-based techniques [8, 12] analyze the patch and get its correct change type, then generate messages with pre-defined patterns. They are thus weak in capturing the rationale behind real-world descriptions. IR-based approaches [28, 29, 44] leverage IR techniques to recall descriptions of the most similar patches from the train set and output

them as the "generated" descriptions for the test patches. They generally fail when there is no similar patch between the train set and the test set. Generation-based techniques [15, 42, 46, 74] try to learn the semantics of edit operations for patch description generation. Existing such approaches do not account for the bimodal nature of patches (i.e., sequence and structure), hence losing the semantics either from the sequential order information or from the semantic logic in the structural abstract syntax trees. With Patcherizer, in order to capture sufficient semantics for patches, we take advantage of both by fusing *SeqIntention* and *GraphIntention*.

Patch correctness: The state-of-the-art automated program repair techniques mainly rely on the test suite to validate generated patches. Given the weakness of test suites, validated patches are actually only plausible since they can still be incorrect [19–21, 51, 62, 63], due to overfitting. The research community is therefore investigating efficient methods of automating patch correctness assessment. While some good results can be achieved with dynamic methods [56], static methods are more scalable. Recently, Tian *et al.* [64] proposed *Panther*, which explored the feasibility of comparing overfitting and correct patches through representation learning techniques (e.g., CC2Vec [28] and Bert [14]). We show in this work that the representations yielded by Patcherizer can vastly improve the results yielded by Panther compared to its current representation learning approaches.

7 CONCLUSION

We present Patcherizer, a novel distributed patch representation learning approach, which fuses contextual, structural, and sequential information in code changes. In Patcherizer, we model sequential information by the Sequence Intention Encoder to give the model the ability to capture contextual sequence semantics and the sequential intention of patches. In addition, we model structural information by the Graph Intention Encoder to obtain the structural change semantics. Sequence Intention Encoder and Graph Intention Encoder enable Patcherizer to learn high-quality patch representations.

We evaluate Patcherizer on three tasks, and the results demonstrate that it outperforms several baselines, including the state-of-the-art, by substantial margins. An ablation study further highlights the importance of the different design choices. Finally, we compare the robustness of Patcherizer vs the CC2Vec and CCRep state-of-the-art patch representation approach on an independent dataset. The empirical result shows that Patcherizer is more effective.

Data Availability: We make our code and dataset publicly available at:

<https://anonymous.4open.science/r/Patcherizer-1E04>

REFERENCES

- [1] 2018. NetworkX. In *Encyclopedia of Social Network Analysis and Mining*, 2nd Edition, Reda Alhajj and Jon G. Rokne (Eds.). Springer. https://doi.org/10.1007/978-1-4939-7131-2_100771
- [2] Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. 2018. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 1–37.
- [3] Uri Alon, Roy Sadaka, Omer Levy, and Eran Yahav. 2020. Structural Language Models of Code. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 245–256. <http://proceedings.mlr.press/v119/alon20a.html>
- [4] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. code2vec: learning distributed representations of code. *PACMPL* 3, POPL (2019), 40:1–40:29. <https://doi.org/10.1145/3290353>
- [5] Satyanjeev Banerjee and Alon Lavie. 2005. 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*. 65–72.
- [6] Earl T Barr, Yuriy Brun, Premkumar Devanbu, Mark Harman, and Federica Sarro. 2014. The plastic surgery hypothesis. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 306–317.
- [7] Shaked Brody, Uri Alon, and Eran Yahav. 2020. A structural model for contextual code changes. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 1–28.
- [8] Raymond PL Buse and Westley R Weimer. 2010. Automatically documenting program changes. In *Proceedings of the IEEE/ACM international conference on Automated software engineering*. 33–42.
- [9] Rocio Cabrera Lozoya, Arnaud Baumann, Antonino Sabetta, and Michele Bezzi. 2021. Commit2vec: Learning distributed representations of code changes. *SN Computer Science* 2, 3 (2021), 1–16.
- [10] Matteo Ciniselli, Nathan Cooper, Luca Pascarella, Denys Poshyvanyk, Massimiliano Di Penta, and Gabriele Bavota. 2021. An empirical study on the usage of BERT models for code completion. In *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*. IEEE, 108–119.
- [11] Luigi P Cordella, Pasquale Foggia, Carlo Sansone, and Mario Vento. 1999. Performance evaluation of the VF graph matching algorithm. In *Proceedings 10th International Conference on Image Analysis and Processing*. IEEE, 1172–1177.
- [12] Luis Fernando Cortés-Coy, Mario Linares-Vásquez, Jairo Aponte, and Denys Poshyvanyk. 2014. On automatically generating commit messages via summarization of source code changes. In *2014 IEEE 14th International Working Conference on Source Code Analysis and Manipulation*. IEEE, 275–284.
- [13] Daniel DeFreez, Aditya V Thakur, and Cindy Rubio-González. 2018. Path-based function embedding and its application to error-handling specification mining. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 423–433.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [15] Jinhao Dong, Yiling Lou, Qihao Zhu, Zeyu Sun, Zhilin Li, Wenjie Zhang, and Dan Hao. 2022. FIRA: Fine-Grained Graph-Based Code Change Representation for Automated Commit Message Generation. (2022).
- [16] Robert Dyer, Hoan Anh Nguyen, Hridesh Rajan, and Tien N Nguyen. 2013. Boa: A language and infrastructure for analyzing ultra-large-scale software repositories. In *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 422–431.
- [17] Ahmed Elnaggar, Wei Ding, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Silvia Severini, Florian Matthes, and Burkhard Rost. 2021. CodeTrans: Towards Cracking the Language of Silicon's Code Through Self-Supervised Deep Learning and High Performance Computing. *arXiv preprint arXiv:2104.02443* (2021).
- [18] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).
- [19] Xiang Gao, Bo Wang, Gregory J Duck, Ruyi Ji, Yingfei Xiong, and Abhik Roychoudhury. 2021. Beyond tests: Program vulnerability repair via crash constraint extraction. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 30, 2 (2021), 1–27.
- [20] Ali Ghanbari and Andrian Marcus. 2022. Patch Correctness Assessment in Automated Program Repair Based on the Impact of Patches on Production and Test Code.
- [21] Matthias Páll Gissurarson, Leonhard Applis, Annibale Panichella, Arie van Deursen, and David Sands. 2022. PropR: property-based automatic program repair. In *Proceedings of the 44th International Conference on Software Engineering*. 1768–1780.
- [22] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 249–256.
- [23] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 315–323.
- [24] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin B. Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. 2021. GraphCodeBERT: Pre-training Code Representations with Data Flow. In *9th International Conference on Learning Representations, ICLR 2021*,

- Virtual Event, Austria, May 3-7, 2021. OpenReview.net. <https://openreview.net/forum?id=jLoC4ez43PZ>
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
 - [26] Jordan Henkel, Shuvendu K Lahiri, Ben Liblit, and Thomas Reps. 2018. Code vectors: Understanding programs through embedded abstracted symbolic traces. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 163–174.
 - [27] Thong Hoang, Hoa Khanh Dam, Yasutaka Kamei, David Lo, and Naoyasu Ubayashi. 2019. DeepJIT: an end-to-end deep learning framework for just-in-time defect prediction. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. IEEE, 34–45.
 - [28] Thong Hoang, Hong Jin Kang, David Lo, and Julia Lawall. 2020. Cc2vec: Distributed representations of code changes. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 518–529.
 - [29] Yuan Huang, Nan Jia, Hao-Jie Zhou, Xiang-Ping Chen, Zi-Bin Zheng, and Ming-Dong Tang. 2020. Learning human-written commit messages to document code changes. *Journal of Computer Science and Technology* 35, 6 (2020), 1258–1277.
 - [30] Nan Jiang, Thibaud Lutellier, and Lin Tan. 2021. CURE: Code-aware neural machine translation for automatic program repair. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 1161–1173.
 - [31] Siyuan Jiang, Ameer Armaly, and Collin McMillan. 2017. Automatically generating commit messages from diffs using neural machine translation. In *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering*. IEEE, 135–146.
 - [32] Yasutaka Kamei, Takafumi Fukushima, Shane McIntosh, Kazuhiro Yamashita, Naoyasu Ubayashi, and Ahmed E Hassan. 2016. Studying just-in-time defect prediction using cross-project models. *Empirical Software Engineering* 21, 5 (2016), 2072–2106.
 - [33] Yasutaka Kamei, Emad Shihab, Bram Adams, Ahmed E Hassan, Audris Mockus, Anand Sinha, and Naoyasu Ubayashi. 2012. A large-scale empirical study of just-in-time quality assurance. *IEEE Transactions on Software Engineering* 39, 6 (2012), 757–773.
 - [34] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
 - [35] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
 - [36] Min Li, Zhenjiang Miao, Xiao-Ping Zhang, and Wanru Xu. 2021. An Attention-Seq2Seq Model Based on CRNN Encoding for Automatic Labanotation Generation from Motion Capture Data. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 4185–4189.
 - [37] Bo Lin, Shangwen Wang, Ming Wen, and Xiaoguang Mao. 2022. Context-aware code change embedding for better patch correctness assessment. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 31, 3 (2022), 1–29.
 - [38] Mario Linares-Vásquez, Luis Fernando Cortés-Coy, Jairo Aponte, and Denys Poshyvanyk. 2015. Changelog: A tool for automatically generating commit messages. In *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, Vol. 2. IEEE, 709–712.
 - [39] Fang Liu, Ge Li, Bolin Wei, Xin Xia, Zhiyi Fu, and Zhi Jin. 2020. A self-attentional neural architecture for code completion with multi-task learning. In *Proceedings of the 28th International Conference on Program Comprehension*. 37–47.
 - [40] Fang Liu, Ge Li, Yunfei Zhao, and Zhi Jin. 2020. Multi-task learning based pre-trained language model for code completion. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*. 473–485.
 - [41] Qin Liu, Zihui Liu, Hongming Zhu, Hongfei Fan, Bowen Du, and Yu Qian. 2019. Generating commit messages from diffs using pointer-generator network. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. IEEE, 299–309.
 - [42] Shangqing Liu, Cuiyun Gao, Sen Chen, Nie Lun Yiu, and Yang Liu. 2020. ATOM: Commit message generation based on abstract syntax tree and hybrid ranking. *IEEE Transactions on Software Engineering* (2020).
 - [43] Zhongxin Liu, Zhijie Tang, Xin Xia, and Xiaohu Yang. 2023. CCRep: Learning Code Change Representations via Pre-Trained Code Model and Query Back. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 17–29. <https://doi.org/10.1109/ICSE48619.2023.00014>
 - [44] Zhongxin Liu, Xin Xia, Ahmed E Hassan, David Lo, Zhenchang Xing, and Xinyu Wang. 2018. Neural-machine-translation-based commit message generation: how far are we?. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 373–384.
 - [45] Chunjie Luo, Jianfeng Zhan, Xiaohu Xue, Lei Wang, Rui Ren, and Qiang Yang. 2018. Cosine normalization: Using cosine similarity instead of dot product in neural networks. In *International Conference on Artificial Neural Networks*. Springer, 382–391.
 - [46] Lun Yiu Nie, Cuiyun Gao, Zhicong Zhong, Wai Lam, Yang Liu, and Zenglin Xu. 2021. CoreGen: Contextualized Code Representation Learning for Commit Message Generation. *Neurocomputing* 459 (2021), 97–107.
 - [47] Michael Niemeyer and Andreas Geiger. 2021. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11453–11464.
 - [48] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
 - [49] Weiguo Pian, Hanyu Peng, Xunzhu Tang, Tiezhu Sun, Haoye Tian, Andrew Habib, Jacques Klein, and Tegawendé F Bissyandé. 2022. MetaTPTrans: A Meta Learning Approach for Multilingual Code Representation Learning. *arXiv preprint arXiv:2206.06460* (2022).
 - [50] Weiguo Pian, Hanyu Peng, Xunzhu Tang, Tiezhu Sun, Haoye Tian, Andrew Habib, Jacques Klein, and Tegawendé F Bissyandé. 2023. MetaTPTrans: A meta learning approach for multilingual code representation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 5239–5247.
 - [51] Zichao Qi, Fan Long, Sara Achour, and Martin Rinard. 2015. 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*. 24–36.
 - [52] Libo Qin, Tailu Liu, Wanxiang Che, Bingbing Kang, Sendong Zhao, and Ting Liu. 2021. A co-interactive transformer for joint slot filling and intent detection. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 8193–8197.
 - [53] Lin CY ROUGE. 2004. A package for automatic evaluation of summaries. In *Proceedings of Workshop on Text Summarization of ACL, Spain*.
 - [54] Reuven Rubinfeld. 1999. The cross-entropy method for combinatorial and continuous optimization. *Methodology and computing in applied probability* 1, 2 (1999), 127–190.
 - [55] Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368* (2017).
 - [56] Ridwan Shariffdeen, Yannic Noller, Lars Grunske, and Abhik Roychoudhury. 2021. Concolic program repair. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*. 390–405.
 - [57] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-Attention with Relative Position Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)*, Marilyn A. Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, 464–468. <https://doi.org/10.18653/v1/n18-2074>
 - [58] Ensheng Shi, Yanlin Wang, Lun Du, Hongyu Zhang, Shi Han, Dongmei Zhang, and Hongbin Sun. 2021. CAST: Enhancing Code Summarization with Hierarchical Splitting and Reconstruction of Abstract Syntax Trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 4053–4062. <https://doi.org/10.18653/v1/2021.emnlp-main.332>
 - [59] Alexey Svyatkovskiy, Ying Zhao, Shengyu Fu, and Neel Sundaresan. 2019. Pythia: Ai-assisted code completion system. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2727–2735.
 - [60] Xunzhu Tang, Rujie Zhu, Tiezhu Sun, and Shi Wang. 2021. Moto: Enhancing Embedding with Multiple Joint Factors for Chinese Text Classification. In *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2882–2888.
 - [61] Chris Thunes. 2013. javalang: pure Python Java parser and tools.
 - [62] Haoye Tian, Yinghua Li, Weiguo Pian, Abdoul Kader Kaboré, Kui Liu, Andrew Habib, Jacques Klein, and Tegawendé F Bissyandé. 2022. Predicting Patch Correctness Based on the Similarity of Failing Test Cases. *ACM Transactions on Software Engineering and Methodology* (2022).
 - [63] Haoye Tian, Kui Liu, Abdoul Kader Kaboré, Anil Koyuncu, Li Li, Jacques Klein, and Tegawendé F Bissyandé. 2020. Evaluating representation learning of code changes for predicting patch correctness in program repair. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 981–992.
 - [64] Haoye Tian, Kui Liu, Yinghua Li, Abdoul Kader Kaboré, Anil Koyuncu, Andrew Habib, Li Li, Junhao Wen, Jacques Klein, and Tegawendé F Bissyandé. 2022. The Best of Both Worlds: Combining Learned Embeddings with Engineered Features for Accurate Prediction of Correct Patches. *arXiv preprint arXiv:2203.08912* (2022).
 - [65] Haoye Tian, Xunzhu Tang, Andrew Habib, Shangwen Wang, Kui Liu, Xin Xia, Jacques Klein, and Tegawendé F Bissyandé. 2022. Is this Change the Answer to that Problem? Correlating Descriptions of Bug and Code Changes for Evaluating Patch Correctness. *arXiv preprint arXiv:2208.04125* (2022).
 - [66] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
 - [67] Ashwin K Vijayakumar, Michael Cogswell, Ramprasad R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424* (2016).

- [68] Haoye Wang, Xin Xia, David Lo, Qiang He, Xinyu Wang, and John Grundy. 2021. Context-aware retrieval-based deep commit message Generation. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 30, 4 (2021), 1–30.
- [69] Ke Wang, Rishabh Singh, and Zhendong Su. 2017. Dynamic neural program embedding for program repair. *arXiv preprint arXiv:1711.07163* (2017).
- [70] Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma, Lingfan Yu, Yu Gai, et al. 2019. Deep graph library: A graph-centric, highly-performant package for graph neural networks. *arXiv preprint arXiv:1909.01315* (2019).
- [71] Shi Wang, Daniel Tang, Luchen Zhang, Huilin Li, and Ding Han. 2022. Hienet: Bidirectional hierarchy framework for automated icd coding. In *International Conference on Database Systems for Advanced Applications*. Springer, 523–539.
- [72] Shi Wang, Daniel Tang, Luchen Zhang, Huilin Li, and Ding Han. 2022. HieNet: Bidirectional Hierarchy Framework for Automated ICD Coding. In *Database Systems for Advanced Applications - 27th International Conference, DASFAA 2022, Virtual Event, April 11-14, 2022, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 13246)*, Arnab Bhattacharya, Janice Lee, Mong Li, Divyakant Agrawal, P. Krishna Reddy, Mukesh K. Mohania, Anirban Mondal, Vikram Goyal, and Rage Uday Kiran (Eds.). Springer, 523–539. https://doi.org/10.1007/978-3-031-00126-0_38
- [73] Yue Wang, Weishi Wang, Shafiq R. Joty, and Steven C. H. Hoi. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 8696–8708. <https://doi.org/10.18653/v1/2021.emnlp-main.685>
- [74] Shengbin Xu, Yuan Yao, Feng Xu, Tianxiao Gu, Hanghang Tong, and Jian Lu. 2019. Commit message generation for source code changes. In *IJCAI*.
- [75] Pengcheng Yin, Graham Neubig, Miltiadis Allamanis, Marc Brockschmidt, and Alexander L. Gaunt. 2019. Learning to Represent Edits. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=BJl6AjC5F7>
- [76] Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. 2019. A novel neural source code representation based on abstract syntax tree. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 783–794.
- [77] Si Zhang, Hanghang Tong, Jiejun Xu, and Ross Maciejewski. 2019. Graph convolutional networks: a comprehensive review. *Computational Social Networks* 6, 1 (2019), 1–23.
- [78] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced language representation with informative entities. *arXiv preprint arXiv:1905.07129* (2019).
- [79] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 1441–1451. <https://doi.org/10.18653/v1/P19-1139>