

RAP-Gen: Retrieval-Augmented Patch Generation with CodeT5 for Automatic Program Repair

Anonymous Author(s)
Submission Id: 131

ABSTRACT

Automatic program repair (APR) is crucial to reduce manual debugging efforts for developers and improve software reliability. While conventional search-based techniques typically rely on heuristic rules or a redundancy assumption to mine fix patterns, recent years have witnessed the surge of deep learning (DL) based approaches to automate the program repair process in a data-driven manner. However, their performance is often limited by a fixed set of parameters to model the highly complex search space of APR.

To ease such burden on the parametric models, in this work, we propose a novel **R**etrieval-**A**ugmented **P**atch **G**eneration framework (RAP-Gen) by explicitly leveraging relevant fix patterns retrieved from a codebase of previous bug-fix pairs. Specifically, we build a hybrid patch retriever to account for both lexical and semantic matching based on the raw source code in a language-agnostic manner, which does not rely on any code-specific features. In addition, we adapt a code-aware language model CodeT5 as our foundation model to facilitate both patch retrieval and generation tasks in a unified manner. We adopt a stage-wise approach where the patch retriever first retrieves a relevant external bug-fix pair to augment the buggy input for the CodeT5 patch generator, which synthesizes a ranked list of repair patch candidates. Notably, RAP-Gen is a generic APR framework that can flexibly integrate different patch retrievers and generators to repair various types of bugs.

We thoroughly evaluate RAP-Gen on three benchmarks in two programming languages, including the TFix benchmark in JavaScript, and Code Refinement and Defects4J benchmarks in Java, where the bug localization information may or may not be provided. Experimental results show that RAP-Gen significantly outperforms previous state-of-the-art DL-based approaches on all benchmarks, e.g., boosting the accuracy of T5-large on TFix from 49.70% to 54.15% (corresponding to repair 478 more bugs) with a much smaller model size. Further analysis reveals that our patch retriever can search for relevant fix patterns to guide the APR systems.

1 INTRODUCTION

Program repair is one of the most important stages to maintain software quality, which however is a time-consuming and cost-dominating process in modern software development [44, 63]. Therefore, there have been huge needs for Automatic Program Repair (APR) tools to ease the difficulty and cost of program repair for developers with use cases including search of patches at program development time [40], build time [38, 58] or run time [8, 42].

A notable class of conventional techniques for APR is known as search-based (also referred to as generate-and-validate) approach [11, 18, 31, 46, 62, 64]. They often search for repairs based on the *fix*

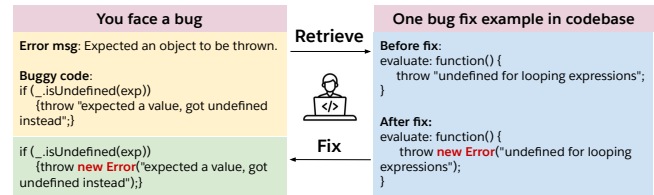


Figure 1: One example of how developers fix a bug by referring to a retrieved fix pattern in codebase.

patterns mined via manual heuristic rules [23, 46, 53] or redundancy-based techniques [11, 18, 31–33, 64]. The latter group of approaches make a *redundancy assumption* [65] that the fixed patch can often be found (or reconstructed) from elsewhere in the codebase (a donor code snippet). This hypothesis has been validated empirically by studies [2, 39] showing that a significant proportion of commits (3%–17%) are indeed composed of existing codebase.

Meanwhile, with the recent advancement in deep learning technologies, numerous deep learning (DL)-based APR approaches [3, 6, 19, 37, 57, 68] have been proposed to automate the repair process via parametric models in a purely data-driven manner. In this paradigm, the APR task is typically formulated as a neural machine translation (or sequence-to-sequence learning) problem [52] in order to translate a buggy (source) program into a correct (target) version. Despite their promising results in software intelligence tasks, their performance is often limited by the fixed set of model parameters to learn the highly complex distributional patterns for program repair, even with several hundreds of million parameters [3, 19, 68].

To ease such burden on the parametric neural models, in this work, we propose a novel retrieval-augmented patch generation framework called RAP-Gen to additionally leverage relevant fix patterns from a patch retriever. Earlier APR techniques based on the redundancy assumption have shown that mining fix patterns from existing codebase [11, 18] or even external Q&As from Stack-Overflow [31] can serve as crucial repair ingredients for APR. Our model, which is semi-parametric in nature, aims to combine both benefits of the implicit (parametric) end-to-end program repair learning and the explicit (non-parametric) fix pattern mining. One distinction from prior fix pattern mining work is that we utilize the top relevant *bug-fix pair* as a guiding fix pattern for a buggy patch instead of clustering the fix templates with hand-crafted heuristics. This retrieval-guiding strategy is also motivated by debugging behaviours of program developers, who often search for relevant bug-fix examples to distill some repair clues for bug fixing. As shown in Fig. 1, the retrieved previous repair example informs a fix pattern of wrapping the string with an “Error” object for the “throw” statement, which is applicable to fix the bug under consideration.

In addition, we propose to adapt a Transformer-based [60] encoder-decoder model CodeT5 [61] as the unified foundation model of

RAP-Gen for both patch retrieval and generation tasks. CodeT5 is a generic code-aware language model pretrained on large source code corpora in eight popular programming languages (including JavaScript and Java) curated from GitHub, achieving state-of-the-art (SoTA) performance in both code understanding and generation tasks. RAP-Gen adopts a stage-wise learning strategy to connect the patch retriever and patch generator: the patch retriever first searches for a relevant bug fix pattern and then pass it to the CodeT5 patch generator to synthesize a ranked list of fix patch candidates based on both the source buggy code and the retrieved external bug fix knowledge. While such retrieval-augmented generation paradigm has been explored in other tasks such as question answering [22] and code generation and summarization [41], we are the first to investigate its effectiveness for APR systems based on large-scale pretrained language models for code.

For the retrievers, we propose a hybrid approach that accounts for both lexical and semantic matching through sparse (BM25 [50]) and dense (DPR [22]) retrieval based on the raw source code. We employ CodeT5’s encoder as our dense DPR retriever and propose to train it with a contrastive learning objective [59] using previous bug-fix pairs as the fix patch often shares most of semantics with its buggy patch. The dense DPR retriever is expected to capture deeper code semantics while the sparse keyword-based BM25 retriever focuses more on the lexical similarity which is sensitive to the choice of naming for code identifiers. Notably, the hybrid retriever is language-agnostic as it does not require any code-specific features such as abstract syntax trees (ASTs). Experiments reveal that our patch retriever is able to retrieve lexically and semantically relevant fix patterns to guide APR systems.

We investigate the effectiveness of RAP-Gen in different APR scenarios including JavaScript linter-raised diagnostics (TFix [3]), Java bug-fix commits (Code Refinement [57]), and real Java bugs accompanied with test cases in open source projects (Defects4J [21]). Among these benchmarks, bug localization information is provided for TFix and Defects4J, while on Code Refinement, such information is absent and the model is required to implicitly localize the bugs.

For evaluation, apart from the conventional metric of the *correct repair number* by executing the fix patch against test suites, we employ a comprehensive set of metrics to suit various APR scenarios. Specifically, we further use *error removal* accuracy based on the detection result from static analyzers for JavaScript coding errors. Additionally, we employ *exact match* accuracy which has been widely adopted [3, 20, 57] to ensure the correctness with machine-generated patches based on developers’ real fix, as test suites are often unavailable in many APR benchmarks especially for the large-scale ones such as TFix and Code Refinement. From Fig. 4, We showed that exact match as a strict metric is well aligned with other looser metrics like error removal. Overall, extensive experimental results show that our RAP-Gen significantly outperforms existing DL-based methods on all these three APR benchmarks.

In summary, the paper makes the following contributions:

- We propose a novel retrieval-augmented patch generation framework (RAP-Gen) for APR. It is a generic framework that can be easily integrated with any sequence-to-sequence learning models. To the best of our knowledge, this is the first work to leverage the power of retrieval in fix pattern mining for DL-based APR systems.

- We present a hybrid patch retriever for fix pattern mining that accounts for both lexical and semantic matching through a combination of sparse and dense retrievers. It is a language-agnostic patch retriever using raw source code which does not require any code-specific features.
- We propose to adapt a generic pretrained code-aware language model CodeT5 as a foundation model for RAP-Gen to fix various bugs. Moreover, we leverage it for both patch retrieval and generation task in a unified manner.
- We extensively evaluate RAP-Gen on three APR benchmarks in JavaScript and Java. Results show RAP-Gen significantly outperforms SoTA DL-based methods on all benchmarks. Particularly, our best model yields substantial improvements (49.70 \rightarrow 54.15 on exact match accuracy and 69.30 \rightarrow 78.80 on error removal accuracy) on TFix over the previous SoTA T5-large model with a 3.5x larger model size than ours. On Code Refinement, RAP-Gen sets new SoTA exact match results of 24.80 and 15.84 over CodeT5’s 21.61 and 13.96 for the small and medium subsets. On Defects4J, RAP-Gen fixes 11 more bugs (33 \rightarrow 44) among 83 single-line bugs than CodeT5, surpassing previous SoTA’s 38 correct fixes with a much smaller patch candidate pool.

2 RELATED WORKS

2.1 Automatic Program Repair

In the past decades, automatic program repair (APR) has attracted growing attention and various APR techniques have been proposed to reduce the manual efforts in debugging. A notable class of conventional techniques for APR is known as search-based (or generate-and-validate) approach [11, 18, 31, 46, 62, 64]. Earlier search-based APR techniques are often based on program modification or mutation with heuristic algorithm [46] or genetic programming [67] to produce a large pool of candidate fixes for validating with unit tests. The search strategy has been further extended to adopt *fix patterns* mined using redundancy-based techniques [11, 18, 24, 31–33, 64]. These approaches make a *redundancy assumption* [65] that the fixed patch can often be reconstructed from elsewhere in the codebase, which has been validated empirically by studies [2, 39] showing that a significant proportion of commits (3%-17%) are indeed composed of existing codebase. More redundancy-based techniques have shown that mining fix patterns from existing codebase [11, 18] or even external Q&As from StackOverflow [31] can largely benefit APR systems.

Recently, with the recent advancement in deep learning (DL) approaches for natural language processing (NLP), many DL-based APR techniques [3, 6, 19, 37, 57, 68] have been proposed to automate the program repair process in an end-to-end data-driven manner. Motivated by the success of Neural Machine Translation (NMT), these techniques often formulate APR as a sequence-to-sequence NMT problem [52], which is to translate a buggy program into a fixed version. Various neural architectures have been explored in learning-based APR techniques. Earlier techniques [6, 57] are based on recurrent neural networks [15], which is further extended to convolution neural networks [10] in CoCoNuT [37] and Transformer-based models [60] in CURE [19], Recoder [68], RewardRepair [66], and TFix [3]. Particularly, CURE and TFix leverage the pretrained

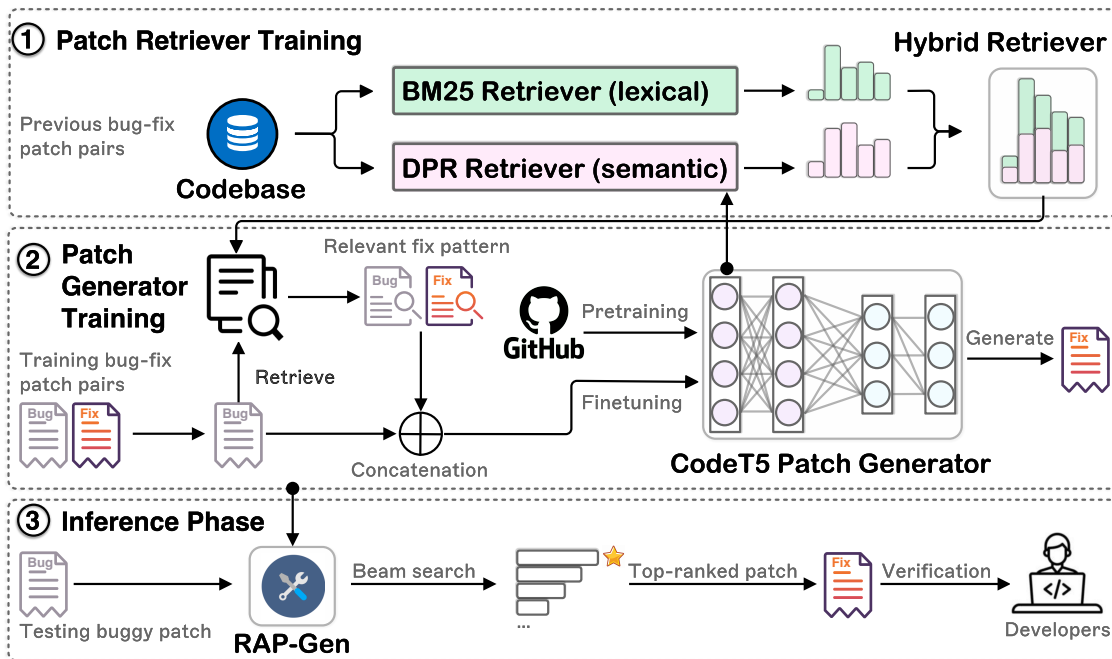


Figure 2: Retrieval-augmented patch generation (RAP-Gen) framework with CodeT5 for automatic program repair. We first retrieve the relevant bug fix patterns from the codebase through our hybrid patch retriever which takes both lexical and semantic similarity into account. We then concatenate the top-1 retrieved bug fix pattern with the query buggy patch for our patch generator to synthesize a ranked list of fix patch candidates for developers to verify.

language models in NLP such as GPT [48] and T5 [49] respectively. In this work, we propose to leverage the code-aware pretrained programming language model CodeT5 [61] for better APR.

2.2 Programming Language Models

Pretrained language models like GPT [48], BERT [7], and T5 [49] have significantly boosted performance in a broad set of NLP tasks. Inspired by their success, much recent work [1, 5, 9, 13, 36, 43, 61] attempts to adapt the NLP pretraining methods to programming language. They often rely on either an encoder-only BERT-style models (CodeBERT [9] and GraphCodeBERT [13]) or decoder-only GPT-style models (CodeGPT [36] and Codex [5]), or encoder-decoder models (PLBART [1] and CodeT5 [61]). Among these, CodeT5 is a unified language model pretrained with a code-aware pretraining objective on large-scale code corpora covering 8 different programming languages. It has been shown to achieve SoTA performance on a wide range of code understanding and generation tasks in CodeXGLUE benchmark [36], which is selected as the backbone of our RAP-Gen framework.

There are recent attempts [20, 45] to explore few-shot learning of large language models (LLMs) for APR. According to Prenner et al. [45], their method based on Codex achieves 46% EM compared to the finetuned T5’s 59% on a random sample of 200 instances from TFix, showing that there is still a gap between few-shot learning and finetuning results. Besides, few-shot learning of LLMs requires more engineering efforts for prompting tuning and post-processing [20], which is labor-intensive. Another concern is that LLMs such as

Codex [5] are not open sourced and it might be expensive to use their APIs, e.g., the Davinci version costs \$0.02 for every 1K tokens¹.

2.3 Retrieval-Augmented Generation

A general retrieval-augmented generation paradigm is comprised of three components including information retrieval, data augmentation and generation model [26]. It has been widely studied in NLP and shown to achieve SoTA performance in a wide range of NLP tasks including question answering and question generation [17, 25] and machine translation [12]. Inspired by their success, much research work adapts this paradigm (also referred as retrieve-and-edit/refine framework) to benefit software intelligence tasks, including code autocompletion [14, 35], code summarization [27, 41], and code generation [41]. To the best of our knowledge, we are the first to explore the retrieval-augmented generation method

It is not trivial to adapt such techniques to APR tasks and it requires systematic adaptation to address some unique challenges. The first challenge is how to retrieve relevant fix patterns for effectively guiding APR, where we build a hybrid retriever based on both lexical and semantic information, which is analyzed and compared with other retrievers in Table 7. The second challenge is how to build a top-performing APR model for various languages and APR scenarios. We leverage a language-agnostic pretrained model CodeT5 for both retrieval and patch generation, which is a more unified approach compared to prior work [35, 41] requiring a different retriever and generator. This is the first work to explore CodeT5 for unified retrieval-augmented generation.

¹<https://openai.com/api/pricing/>

3 APPROACH

As shown in 2, our proposed RAP-Gen aims to generate a target fix patch based on an input buggy patch together with a relevant bug fix pattern retrieved from a codebase of previous bug-fix pairs. We first introduce the task formulation of the retrieval-augmented patch generation for program repair in 3.1 and then revisit fundamental components of CodeT5 in 3.2, followed by detailing the two components of our approach: the *hybrid patch retriever* in 3.3 and the *retrieval-augmented patch generator* in 3.4.

3.1 Task Formulation

Let $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^{|\mathcal{D}|}$ be a program repair dataset consisting of $|\mathcal{D}|$ bug-fix pairs (X_i, Y_i) , where X_i and Y_i are the i -th buggy and fixed program patch, respectively. Assume that we have a codebase containing a large collection of previous bug-fix pairs $\mathcal{C} = \{(B_j, F_j)\}_{j=1}^{|\mathcal{C}|}$, where (B_j, F_j) denotes the j -th previous bug-fix pair. Given a buggy program patch X_i in \mathcal{D} , a retriever retrieves the most relevant bug-fix pair (B_j, F_j) in the codebase \mathcal{C} based on a relevance scoring function $f_\phi(X_i, B_j)$ parameterized by ϕ . Then the original input sequence X_i is augmented with the retrieved bug-fix pair to form a new input sequence $\hat{X}_i = X_i \oplus B_j \oplus F_j$, where \oplus denotes the concatenation operation. The sequence-to-sequence (seq2seq) generator then generates Y_i from \hat{X}_i in an autoregressive manner. Formally, we aim to learn the following probability with the patch seq2seq generator parameterized by θ :

$$P_\theta(Y_i|\hat{X}_i) = \prod_{k=1}^n P_\theta(Y_{i,k}|\hat{X}_i, Y_{i,1} : Y_{i,k-1})$$

where $Y_{i,1} : Y_{i,k-1}$ is the previous sequence before the k -th token and n denotes the number of tokens in the target sequence Y_i . Note that we regard the external codebase \mathcal{C} as a non-parametric memory and the retrieved bug-fix pair as a guiding fix pattern for the generator. In probabilistic terms, the retrieval can be formulated as a latent variable $Z_j = (B_j, F_j)$, which is approximated by top-1 in our case. Formally, the probability can be decomposed as:

$$P(Y_i|X_i) = \sum_{j=1}^{|\mathcal{C}|} \underbrace{P_\phi(Z_j|X_i)}_{\text{Retriever}} \underbrace{P_\theta(Y_i|X_i, Z_j)}_{\text{Generator}} \approx P_\theta(Y_i|X_i, Z_j^*)$$

where Z_j^* is the top-1 retrieved output from the retriever $P_\phi(Z_j|X_i)$. We adopt such top-1 approximation as marginalization over large k makes the training and inference complicated and inefficient [25]. We also tried to employ top- k ($k = 2, 3, 5$) with the Fusion-in-Decoding or FiD method [17] but did not observe a salient performance improvement.

3.2 Revisiting CodeT5

CodeT5 [61] is a unified pretrained Transformer-based encoder-decoder language model that achieves SoTA results in both code understanding and generation tasks. It is pretrained on 8.3 million functions in 8 different programming languages (i.e., Ruby, JavaScript, Go, Python, Java, PHP, C, C#) collected from GitHub. CodeT5 employs a set of identifier-aware pretraining objectives to incorporate the code-specific knowledge into the language model.

In this work, we adapt CodeT5 as our dense DPR retriever and patch generator to harness its powerful code understanding capability.

BPE Subword Tokenization. One benefit of using CodeT5 is that it provides a code-specific Byte-Pair Encoding (BPE) [51] tokenizer. It can avoid the prevalent Out-of-Vocabulary (OoV) problems in the code domain as programmers tend to write arbitrary identifiers [28] and it is impossible to build a fixed vocabulary to accommodate arbitrary tokens (commonly known as *open vocabulary* problem [51]). BPE is an algorithm that learns how to efficiently split tokens into subwords based on their frequency distribution. It can also help reduce the vocabulary size as it will split rare tokens into multiple subwords instead of directly adding the whole tokens into the vocabulary. Additionally, as the CodeT5 tokenizer is pretrained and optimized for eight popular programming languages, the resulting tokenization generalizes well. As pointed out by [61], it reduces the tokenized sequence by 30% - 45% on average compared to the default T5 tokenizer [49].

Encoder and Decoder Architecture. CodeT5 consists of a stack of Transformer layers [60] for its encoder and decoder. Each Transformer layer contains a multi-head self-attention for feature aggregation followed by a feed forward layer over the output of previous layer. The final layer produces the hidden states for all input tokens, which can be employed as the code presentation for classification or generation tasks. For the CodeT5 encoder, it utilizes bidirectional attention masks to learn better contextualized representation similar to BERT [7], while the CodeT5 decoder employs causal attention masks to ensure each token can only attend to the previous tokens for better sequence generation. In RAP-Gen framework, we adapt the CodeT5 as the patch generator and its encoder specifically for the dense retriever.

3.3 Hybrid Patch Retriever

The retriever module in RAP-Gen aims to retrieve relevant fix patterns to guide the APR process. It builds on a relevance scoring function $f_\phi(X_i, B_j)$ to compute the relevance between the (query) buggy patch X_i in \mathcal{D} and a previous (key) buggy patch B_j in the codebase \mathcal{C} . As shown in 2 1, we utilize a hybrid approach to combine a lexical-based BM25 [50] retriever and a semantic-based DPR [22] retriever to take both lexical and semantic information into account. Prior work like [22] show that sparse and dense retriever can complement to each other for more robust text retrieval.

Lexical-based Retriever. We employ BM25 [50], a well-known term-based retriever that uses sparse vector representation for lexical matching. BM25 converts each code patch as bag-of-words representation and computes a lexical similarity between the query patch X_i and a candidate patch B_j . The computed similarity score is represented as $f_\phi(X_i, B_j) = BM25(X_i, B_j)$. As a sparse term-based retriever, BM25 is sensitive to the choice of identifier naming in source code which does not impact the code semantics.

Semantic-based Retriever. We employ Dense Passage Retriever (DPR) [22] to retrieve relevant patches via measuring their semantic similarity. To encode the code patch, we use a Transformer-based encoder to map each patch to a fixed-size dense vector. Specifically, we initialize the DPR from a pretrained CodeT5 encoder and train it

	(a) TFix	(b) Code Refinement	(c) Defects4J (Chart-9)
Source Input	Error Information : fix guard-for-in The body of a for-in should be wrapped in an if statement to filter unwanted properties from the prototype. Patch Difference : <pre> }for (e in data.updates) { + if (data.updates.hasOwnProperty(e)) { + ext = data.updates[e]; </pre>	Patch Difference : <pre> private synchronized void METHOD_1 (java.lang.String VAR_1) { - TYPE_1 . i (VAR_2 , (STRING_1 + VAR_1)); + VAR_3 = VAR_1 ; } </pre>	Patch Difference : <pre> public class TimeSeries extends Series implements Cloneable, Serializable { endIndex = -(endIndex + 1); // this is first item AFTER end period endIndex = endIndex - 1; // so this is last item BEFORE end } - if (endIndex < 0) { + if ((endIndex < 0) (endIndex < startIndex)) { emptyRange = true; } if (emptyRange) { </pre>
Retrieval	Buggy Code : <pre> var e, ext; for (e in data.installs) { ext = data.installs[e]; } Fixed Code : var e, ext; for (e in data.installs) { if (data.installs.hasOwnProperty(e)) { ext = data.installs[e]; } } </pre>	Buggy Code : <pre> public void METHOD_1 (java.lang.String VAR_1) { TYPE_1 . i (VAR_2 , ((STRING_1 + VAR_1) + STRING_2)); VAR_3 = VAR_1 ; } Fixed Code : public void METHOD_1 (java.lang.String VAR_1 { VAR_3 = VAR_1 ; } </pre>	Buggy Code : <pre> else if (changeIndex < endIndex) { if (changeType == ListEvent.INSERT) { endIndex++; updates.addInsert(changeIndex - startIndex); } } Fixed Code : updates.addUpdate(changeIndex - startIndex); </pre>
CodeT5	<pre> }for (i = 0; i < data.updates.length; i++) { ext = data.updates[i]; } </pre>	<pre> private synchronized void METHOD_1 (java.lang.String VAR_1) { if (VAR_1 != null) { TYPE_1 . i (VAR_2 , (STRING_1 + VAR_1)); VAR_3 = VAR_1 . ; } } </pre>	<pre> if (endIndex < 0 (endIndex == 0 && startIndex < endIndex)) { } </pre>
RAP-Gen	<pre> }for (e in data.updates) { if (data.updates.hasOwnProperty(e)) { ext = data.updates[e]; } } </pre>	<pre> private synchronized void METHOD_1 (java.lang.String VAR_1) { VAR_3 = VAR_1 ; } </pre>	<pre> if (endIndex < 0 (endIndex - startIndex) < 0) { } </pre>

Figure 3: Bug fix examples on three APR benchmarks, where RAP-Gen successfully fix bugs while CodeT5 fails to do so.

for a code-to-code retrieval task. For training the DPR, we propose to use the bug-fix pairs in the codebase by considering the buggy code B_j as the query and the corresponding fixed code F_j as the key. This is based on the assumption that the buggy patch and its fixed patch often shares similar semantics (e.g., identifiers and code structures). This trick avoids the massive manual annotation efforts needed to curate a bug-to-bug search dataset.

For each query patch and candidate patch, we prepend a special token of [CLS] into its tokenized sequence and employ the final layer hidden state of the [CLS] token as the patch representation. We use a shared DPR to separately encode the query patch X_i in \mathcal{D} and a candidate patch B_j in \mathcal{C} as CLS_{X_i} and CLS_{B_j} , respectively. Then the similarity is computed by the inner product between these two patch representations as the following:

$$f_{\phi}(X_i, B_j) = \text{sim}(X_i, B_j) = [CLS_{X_i}]^T [CLS_{B_j}]$$

For training the DPR retriever, we leverage the *in-batch negatives* to optimize an InfoNCE contrastive loss [59] defined as follows:

$$\mathcal{L}_{\text{InfoNCE}} = \frac{1}{N} \sum_{i=1}^N -\log \frac{\exp(\text{sim}(B_i, F_i))}{\exp(\text{sim}(B_i, F_i)) + \sum_{j \in \mathcal{M}, j \neq i} \exp(\text{sim}(B_i, F_j))}$$

where \mathcal{M} is the current minibatch and N denotes the number of positive training examples in the minibatch. This objective aims to maximize the similarity between positive examples while minimizing the similarity between negative examples. Each positive example will have $|\mathcal{M}| - 1$ negative samples. Note that we do not adopt the hard negative mining strategy as in [22] due to the noisy nature of the training data.

In the inference stage, given a query buggy patch X_i , the DPR retrieves a relevant bug-fix pair (B_j, F_j) by computing the similarity between X_i (query) and B_j (key). We also tried to base on the similarity between X_i and F_j but it did not yield better results.

Hybrid Retriever. To take both lexical and semantic information into account, we utilize a hybrid approach following [22] to combine the BM25 and DPR. The similarity score is computed as $f_{\phi}(X_i, B_j) = \text{sim}(X_i, B_j) + \lambda \text{BM25}(X_i, B_j)$, where λ is a weight to balance the

two retrievers and was empirically set to 1 in our experiment. Based on this combined similarity score, we select the top-1 relevant bug-fix pair (B_j, F_j) as a fix pattern to guide the patch generator for bug fixing. The hybrid retriever is expected to be more robust compared to retrievers that rely only on either lexical or semantic information.

3.4 Retrieval-Augmented Patch Generator

As shown in Fig. 2, given a buggy patch X_i , we search for a top relevant fix pattern (B_j, F_j) and pass it to the patch generator to generate a fixed code patch Y_i . We adopt a simple yet effective strategy to augment X_i into $\hat{X}_i = X_i \oplus B_j \oplus F_j$ via appending the retrieved bug-fix pair into the source buggy patch. Note that the patch generator module can be any sequence generation model. Different from prior studies that directly adopt a generator optimized on natural language [3], we propose to employ CodeT5, a code-aware programming language model optimized for code.

Training. We prepare the retrieval-augmented input to CodeT5 patch generator as $\hat{X}_i = \text{[CLS]} X_i \text{[BUG]} B_j \text{[FIX]} F_j$, where [BUG] and [FIX] are special tokens to separate the retrieved bug-fix pair from the buggy patch. CodeT5’s encoder takes \hat{X}_i as input and emits the fixed patch Y_i from its decoder in an autoregressive manner (see Section 3.1). We consider two settings of the buggy patch X_i where it may or may not contain bug localization information. If it contains error information like error type, error message, and error line, the buggy patch will be augmented to “error information [SEP] X_i ” to incorporate error information to help fix the bugs. To train the patch generator, we adopt *teacher forcing* [56] to minimize the cross entropy loss \mathcal{L}_{ce} over all training instances defined as:

$$\mathcal{L}_{ce} = - \sum_{i=1}^{|\mathcal{D}|} \log(P_{\theta}(Y_i | \hat{X}_i))$$

In teacher forcing, the decoder uses ground-truth context for faster convergence. We use the training set as the search codebase following [41]. To avoid information leakage, we do not allow the retriever to access the ground-truth bug-fix pair, otherwise the training loss

Table 1: Statistics of three program repair benchmarks.

Benchmark	Version	Train	Valid	Test	All
TFix	Original	84,846	9,454	10,504	104,804
TFix	Deduplicated	84,673	9,423	10,465	104,561
Code Refinement	Small	46,680	5,835	5,835	58,350
Code Refinement	Medium	52,364	6,545	6,545	65,454
Defects4J	v1.2 (single-line)	-	-	83	83

would easily drop close to 0 as the generator can directly copy the retrieved fix as the target output. This strategy makes the training and evaluation process more compatible as the evaluation sets are not overlapped with the training set as well.

Inference with Beam Search. During inference, as shown in Fig. 2 3, we employ beam search to generate a ranked list of fixed patch candidates for an input buggy patch, where the number of predictions is determined by the beam size \mathcal{B} . Concretely, at each decoding timestep, the beam search selects the most \mathcal{B} promising fix candidates with the highest probability using a best-first search strategy. The search process is terminated when an [EOS] token notifying the end of sentence is emitted. The top ranked fix patch will be examined its correctness by comparing with ground-truth fix patches or by validating against test suites, or by manual verification by software developers.

4 EXPERIMENTAL DESIGN

We describe the studied datasets and its comparison baselines and evaluation metrics, followed by detailing implementation details, and research questions in the section.

4.1 Dataset

We evaluate RAP-Gen on three APR datasets, namely TFix [3] in JavaScript, Code Refinement [57] and Defects4J (v1.2) [21] in Java. All datasets are originally collected from open source GitHub commits but based on different criteria for bug identification, where TFix is based on diagnostics from a JavaScript static analyzer, Code Refinement is based on repair-related commit message, and Defects4J is based on running the test suites. We report their data statistics in Table 1.

4.1.1 TFix. TFix [3] is a large-scale program repair dataset comprising JavaScript code patch pairs curated from 5.5 million GitHub commits. It includes 52 error types (see Table 4) detected by a static analyzer ESLint² [55]. In addition to error types, it provides rich error annotations such as error message and localized error line so that there is no need for fault localization like prior work [19, 68]. To prepare the input sequence, as illustrated in Fig. 3(a), we follow [3] to combine all error information together with the buggy code patch into a single piece of text as the following:

**fix {error type} {error message} {error context:
Code Line N-1 + Buggy Line N + Code Line N+1}**

where error context consists of the given localized error line and its two neighboring code lines to form a buggy code patch. For the target sequence, it is obtained by replacing the error line into a

²<https://eslint.org/>

fixed line in the error context. During data processing, we observed a *duplication issue* inside each data split and between data splits. Specifically, there are 114, 2, and 4 duplicates in the train, validation, and test split respectively, and 28, 34, and 4 duplicates for inter-split duplicates between train and test, train and test, validation and test splits respectively. We filtered all these 243 duplicates to get a deduplicated version of TFix as shown in Table 1.

Baseline Models. We compare RAP-Gen with existing DL-based APR models including SequenceR [6] and CoCoNuT [37]. Besides, we compare a large pretrained model T5-large [49] which has been finetuned on TFix to achieve the SoTA performance by [3].

Evaluation Metrics. We report **Exact Match** (EM) accuracy and **BLEU-4** score to evaluate program repair performance following [61] on TFix. BLEU-4 is a looser metric to evaluate the degree of subword overlapping while EM is a more strict metric requiring the prediction to be identical to the ground-truth patch in a real commit. As a buggy program might have different ways to repair, we further employ **Error Removal** metric following [3] to take various forms of fixes into account. The prediction is counted as correct for Error Removal if the existing error is removed and no new errors (detected by the static analyzer ESLint) is introduced after the fix. For all metrics, we present their results on a scale of 0-100 (%) and a higher score represents better performance.

4.1.2 Code Refinement. Code Refinement [57] contains bug-fix pairs at the function level, which are originally collected from public GitHub Archive³ between March 2011 and October 2017. They use Google BigQuery APIs to identify all Java commits having a message containing the patterns: (“fix” or “solve”) and (“bug” or “issue” or “problem” or “error”) to ensure the quality of the collected bug-fix function pairs. They normalized the functions via obfuscating identifiers with indexed tokens such as TYPE1, VAR1, METHOD1, etc. One data example can be found in Fig. 3 (b). The dataset contains two data subsets which are determined by the number of tokens, i.e., # of code tokens ≤ 50 for the small set and $50 < \#$ of code tokens ≤ 100 for the medium set. Since the bug localization is not provided, the entire code fragment is taken as the source input sequence of our model. The target sequence is the refined version of the whole code snippet.

Baseline and Metrics. We compare our RAP-Gen with pretrained programming language models based on Transformers [60]. One group of these models is the encoder-only models such as RoBERTa (code), CodeBERT [9], and GraphCodeBERT [13]. These encoder-only models require a randomly initialized decoder to generate the fix. Besides, we compare with encoder-decoder Transformer models such as PLBART [1] and CoText [43]. NSEdit [16] is a language model with encoder and decoder initialized from CodeBERT and CodeGPT [36] respectively. It is finetuned to generate the fix via a neural-symbolic editing sequence and ranks as the current SoTA model on Code Refinement. We follow [61] to apply BLEU-4 and Exact Match to evaluate the Code Refinement datasets.

4.1.3 Defects4J. Defects4J (v1.2) [21] contains 395 real bug-fix patches in open source GitHub projects. Specifically, we focus on repairing 83 single-line bugs of Defects4J as this is the most general type of bugs evaluated by most prior work [6]. One example

³<https://www.gharchive.org/>

Table 2: Comparison results of CodeT5 on the original TFix.

Model	EM w/ spaces		EM w/o spaces	
	Avg.	W. Avg.	Avg.	W. Avg.
Naive Copy	0.00	0.00	0.00	0.00
SequenceR	17.90	-	-	-
CoCoNuT	11.70	-	-	-
T5-small	44.46	44.44	44.52	44.60
T5-base	48.54	47.63	48.72	47.70
T5-large	49.33	49.65	49.35	49.70
CodeT5-small	47.14	46.35	50.22	50.31
- error-information	45.97	45.80	49.26	49.70
CodeT5-base	50.88	49.42	54.30	53.57
- error-information	46.88	47.17	50.36	51.25

of Defects4J bugs can be found in Fig. 3(c), where “-” denotes a buggy line to be fixed and “+” represents the correct fix committed from a developer. A buggy line and its corresponding code context are combined to form the source input sequence while the target sequence is the fixed line. As Defects4J only has test set, we use the training data from [68] containing 103,585 bug-fix Java patches collected from GitHub commits. To measure the fix performance without the impact of a third-party Fault Localization tool, we focus on the setting with the correct bug location provided.

Baselines and Metrics. We compare RAP-Gen with seven SoTA DL-based APR models: SequenceR [6], CODIT [4], DLFix [29], CoCoNuT [37], CURE [19], DEAR [30], and RewardRepair [66]. We did not include Recoder [68] for comparison as they did not release the correct predictions under the same setting of perfect fault localization and we failed to reproduce the results based on their implementation. For evaluation, we compute how many bugs can be correctly fixed on Defects4J based on unit testing and manual verification following prior work. We first run test suites to automatically identify plausible correct patches for each bug, followed by manual checking to completely verify its correctness. The correct predictions from our RAP-Gen are included in our artifact.

4.2 Implementation Details

We implement RAP-Gen using the deep learning framework PyTorch⁴ and train it with AdamW [34] optimizer. For the training of its neural components, we run these experiments with NVIDIA A100-40G GPUs on the Google Cloud Platform⁵. For each benchmark, we finetune a DPR retriever for 50 epochs using the contrastive loss $\mathcal{L}_{\text{infoNCE}}$ using a batch size of 64 and a learning rate of $2e-5$. We finetune RAP-Gen generator for 30 epochs using a sequence generation loss \mathcal{L}_{ce} using a batch size of 32 with a learning rate of $5e-5$. These best settings are obtained through a grid search for hyper-parameter tuning: batch size in (16, 32, 64) and learning rate in ($1e-4$, $5e-5$, $2e-5$). The training time of DPR retriever is 5-9 hours depending on the training size of the dataset, and the training time of RAP-Gen generator is within 2 days. For lexical-based retrievers, we use an open-sourced Python library⁶

⁴<https://pytorch.org/>

⁵<https://cloud.google.com/>

⁶<https://pypi.org/project/rank-bm25>

Table 3: Performance of RAP-Gen on the deduplicated TFix. Results of T5-large and CodeT5-base are different from Table 2 due to the deduplication. (EM: W. Avg. EM w/o spaces).

Model	EM	BLEU-4
T5-large (TFix)	49.58	76.96
CodeT5-base	53.46	78.92
RAP-Gen		
+Random	52.98	78.52
+BM25	53.88	79.18
+DPR	53.93	79.59
+Hybrid	54.15	79.66

of BM25, which can be efficiently trained on CPU within one hour with multi-processing. During inference, we employ beam search with a beam size of 5 for the TFix and Code Refinement, and 100 for the Defects4J.

For the foundation model of RAP-Gen, we adopt CodeT5 [61] to leverage its encoded general code understanding knowledge. We initialize the patch generator with CodeT5-base and the DPR retriever with its encoder. CodeT5-base contains 12 encoder layers and 12 decoder layers with the parameter size of 220M. We use CodeT5 to refer to this version unless specified with a size tag.

4.3 Research Questions

To investigate the effectiveness of RAP-Gen on APR tasks, we seek to answer the following research questions (RQs):

RQ1: Comparative study with DL-based APR models on TFix. How does RAP-Gen perform to repair JavaScript linter-flagged coding errors on TFix compared with other DL-based APR approaches?

RQ2: Analysis of RAP-Gen predictions on TFix. How does RAP-Gen repair TFix bugs for different error types and patch lengths? What fix operations do RAP-Gen adopt in repairing bugs?

RQ3: Comparative study with DL-based APR models on Code Refinement. How does RAP-Gen perform to repair Java commit-related bugs compared with other DL-based APR approaches?

RQ4: Analysis of our hybrid patch retriever. Can our hybrid patch retriever find relevant fix pattern to guide APR?

RQ5: Comparative study with DL-based APR models on Defects4J? How does RAP-Gen perform to repair Java bugs in open source projects compared with other DL-based APR approaches?

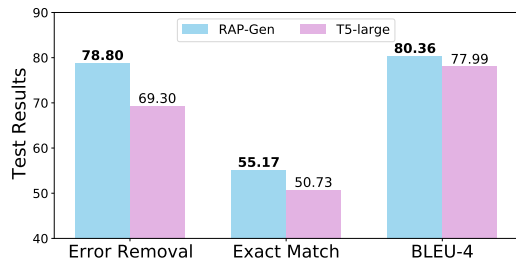
5 EXPERIMENTAL RESULT

5.1 RQ1: Comparative study with DL-based APR models on TFix

5.1.1 Improved TFix Evaluation. The original TFix benchmark employs the direct average of exact match (EM) accuracy across 52 error types as the main metric. However, as shown in the Table 4, these error types have a rather imbalanced distribution, e.g., the major error type “no-invalid-this” has 16,166 instances while the least error type “no-new-symbol” has only 10 instances. As such, it is more reasonable to employ the weighted average to take the error type distribution into account. Besides, we spot another limitation of its exact match evaluation that if the predicted fix contains one more whitespace such as a space or new line than the ground-truth

Table 4: Performance breakdown on 52 error types on TFix.

Error Type	#Samples	T5-large	RAP-Gen	Error Type	#Samples	T5-large	RAP-Gen
no-new-symbol	10	100.00	100.00	no-extra-bind	674	70.59	73.53
no-compare-neg-zero	13	0.00	0.00	no-case-declarations	723	58.90	67.12
no-ex-assign	40	25.00	25.00	no-fallthrough	743	76.00	77.33
for-direction	50	40.00	60.00	no-inner-declarations	830	38.10	46.43
no-unsafe-finally	63	42.86	14.29	no-array-constructor	980	86.73	85.71
use-isnan	71	37.50	25.00	no-constant-condition	1,251	48.78	54.47
no-class-assign	111	41.67	50.00	generator-star-spacing	1,396	67.86	72.86
no-dupe-class-members	117	8.33	8.33	no-extra-boolean-cast	1,458	54.11	58.22
no-func-assign	147	46.67	60.00	no-cond-assign	1,472	45.21	47.95
no-empty-pattern	178	27.78	44.44	no-process-exit	1,514	32.89	37.50
no-unused-labels	187	52.63	63.16	no-empty	2,055	26.70	31.55
no-duplicate-case	195	65.00	60.00	no-dupe-keys	2,181	53.42	55.25
getter-return	203	52.38	61.90	prefer-spread	2,466	45.08	45.49
no-sparse-arrays	237	25.00	45.83	no-useless-escape	2,920	35.15	40.61
no-const-assign	277	35.71	42.86	no-console	3,067	73.62	73.94
no-global-assign	318	59.38	68.75	guard-for-in	3,231	41.98	45.37
no-new-wrappers	360	27.78	38.89	no-throw-literal	4,075	72.06	74.51
no-this-before-super	413	47.62	69.05	no-debugger	4,164	94.48	94.24
no-unsafe-negation	423	72.09	76.74	prefer-rest-params	4,534	35.68	43.61
require-yield	429	72.09	76.74	no-unreachable	4,725	63.85	64.69
no-extend-native	443	31.11	26.67	no-extra-semi	5,969	82.61	83.61
no-new-object	446	71.11	66.67	no-redeclare	6,381	49.45	59.78
no-caller	446	20.00	22.22	comma-style	6,395	46.48	52.11
constructor-super	464	59.57	70.21	no-unused-vars	7,765	51.87	56.11
valid-typeof	539	51.85	51.85	no-undef	10,636	22.65	27.35
no-self-assign	610	34.43	44.26	no-invalid-this	16,166	37.48	44.13
				Sum/W. Avg.	104,561	49.58	54.15

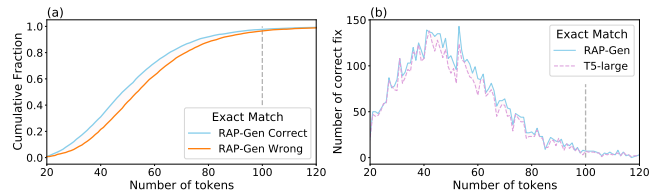
**Figure 4: Error removal comparison on TFix, where error removal is well aligned with exact match and BLEU-4 scores.**

fix, it would be regarded as a wrong exact match. However, extra whitespaces do not impact the correctness for JavaScript programs. Therefore, we propose to use the weighted average of EM w/o spaces, which normalizes the whitespaces before computing the EM to exclude the effects of the mismatch in whitespaces. As we find there is a duplication issue in the TFix dataset, we also report the results on its deduplicated version.

5.1.2 CodeT5 Results. We compare CodeT5 models with other DL-based baselines on TFix and show results in Table 2. For the original metric of average EM w/ spaces, CodeT5-base (50.88) also yields a better accuracy than T5-large (49.33), given that it has much larger model size ($\sim 3.5\times$ of CodeT5-base: 770M vs. 220M). If we focus on a more reasonable average EM w/o spaces, CodeT5-base significantly boost the performance, with around 5 absolute accuracy improvement (49.35 \rightarrow 54.30) over T5-large. Based on the weighted average EM w/o spaces, both CodeT5-small (50.31) and CodeT5-base (53.57) outperform all the baselines including T5-large (49.70). This shows CodeT5 models with code-aware pretraining on large-scale source code have a better understanding of program. For TFix evaluation, we employ EM to denote the weighted average EM w/o spaces. We perform an ablation study to remove the error information including error type and error message from the input sequence, where we observe both CodeT5-small and CodeT5-base models have a consistent performance downgrade, revealing that it is helpful to inform which types of error they need to fix for APR models.

Table 5: Analysis of error line removal operation on TFix.

	T5-large	CodeT5	RAP-Gen
# Ground-truth EL Removal	2,381	2,381	2,381
# Predicted EL Removal	1,882	1,925	1,922
# Correct EL Removal	1,811	1,858	1,866
# False Positive	71	67	56
Precision (%)	96.23	96.52	97.09
Recall (%)	76.06	78.03	78.37
F1 (%)	84.96	86.30	86.73

**Figure 5: (a): cumulative fraction of programs by number of tokens in the source buggy patch, grouped by whether RAP-Gen can accurately fix. (b): distribution of correct fix over number of tokens for RAP-Gen and T5-large.**

5.1.3 RAP-Gen Results. We report the results of our RAP-Gen model on the deduplicated TFix benchmark in Table 3, where the results are slightly different due to data size changes after duplication. We first compare with a Random baseline via randomly retrieving bug-fix pairs from the codebase. The performance downgrade implies that randomly retrieved fix patterns cannot provide useful guiding signals for program repair. Then we compare with RAP-Gen integrated with different retrievers: either BM25 or DPR, or the hybrid approach to combine them. Results show that all retrievers significantly improve the APR performance on both exact match and BLEU-4, where the hybrid one yields the best improvement over T5-large (49.58 \rightarrow 54.15 EM). This indicates retrieval-augmented generation is a viable and effective approach for APR and both semantic information and lexical information are crucial to retrieve relevant fix patterns. We present one case in Fig. 3 (a), where we can observe RAP-gen successfully repairs the bug with the guidance of retrieved fix pattern while CodeT5 without retrieval gives a wrong fix.

5.1.4 Error Removal Evaluation. Though exact match can ensure correctness of machine-generated patches, it might be a too strict metric to consider other forms of correct fixes. Therefore, we follow [3] to employ the error removal metric, where a fix is counted as correct if the error is removed and no new error is introduced. The error detection is based on a static analyzer ESLint. We report error removal together with EM and BLEU-4 results on a large subset of 6,793 instances⁷ in Fig. 4. We observe that RAP-Gen significantly improves error removal accuracy over T5-large (69.30 \rightarrow 78.80). The larger gain compared to EM and BLEU-4 implies that RAP-Gen is more capable of producing various forms of good fixes. Additionally, EM is well aligned with the looser metric of error removal.

⁷Some source files are unavailable to reproduce this metric on the full test set.

Table 6: Performance of RAP-Gen on the Code Refinement.

Model	Small		Medium	
	EM	BLEU-4	EM	BLEU-4
Naive Copy	0.00	78.06	0.00	90.91
LSTM	10.00	76.76	2.50	72.08
Transformer	14.70	77.21	3.70	89.25
RoBERTa (code)	15.90	77.30	4.10	90.07
CodeBERT	16.40	77.42	5.16	91.07
GraphCodeBERT	17.30	80.02	9.10	91.31
PLBART	19.21	77.02	8.98	88.50
CoText	21.58	77.28	13.11	88.40
NSEdit	24.04	71.06	13.87	85.72
CodeT5	21.61	77.43	13.96	87.64
RAP-Gen				
+Random	21.25	76.97	15.13	89.83
+BM25	23.82	78.83	15.37	89.42
+DPR	24.37	78.54	15.60	89.82
+Hybrid	24.80	78.28	15.84	90.01

5.2 RQ2: Analysis of RAP-Gen Predictions on TFix

5.2.1 Performance Breakdown on Error Types. We list the performance breakdown for 52 error types on the deduplicated TFix in Table 4. RAP-Gen outperforms the previous SoTA T5-large in 40/52 error types. Especially for the major error type “no-invalid-this”, RAP-Gen improves its exact match from T5-large’s 37.48 to 44.13, corresponding to repair more 107 instances. In total, RAP-Gen correctly repairs more 478 bugs than T5-large with a much smaller model size.

5.2.2 Fix Operation Analysis. We analyze what fix patterns are performed by our models on TFix. We observe a large proportion of fix consists of deletion operations compared to the code insertion and replacement operations. We find that fix operations consist of code insertion (12.5%), replacement (8.1%), deletion (47.9%), insertion and replacement (6.9%), insertion and deletion (8.2%), replacement and deletion (7.2%), and all three manners (9.2%). Earlier studies [47, 54] also reflect that the deletion operation is one of the most common fix patterns. Besides, we find one dominating fix operation is error line (EL) removal, which is to simply remove the error line from the buggy code and accounts for around 23% in the test set. We show how models perform this operation in Table 5. We observe RAP-Gen achieves the best precision, recall, and F1 scores with a lowest false positive count of 56 compared to CodeT5’s 67 and T5-large’s 71. This indicates that RAP-Gen is able to learn more diverse bug fix patterns instead of over relying on the trivial error line removal pattern.

5.2.3 Patch Length Analysis. We analyze how the patch length influences the APR performance. Fig. 5 (a) shows the cumulative fraction of buggy patches by its patch length grouped based on their outcome. We find the patches successfully repaired by RAP-Gen tend to be shorter than those where it fails. Fig. 5 (b) shows the distribution of correct fixes by its buggy patch length, where

Table 7: Lexical (BLEU-4) and semantic (CosSim) retrieval matching results on TFix and Code Refinement benchmarks.

Retriever	TFix		Refine-Small		Refine-Medium	
	BLEU-4	CosSim	BLEU-4	CosSim	BLEU-4	CosSim
Random	0.1	35.5	14.6	35.4	14.6	30.6
BM25	23.7	70.9	41.5	68.5	39.0	66.6
DPR	21.7	75.4	54.4	84.9	44.3	81.3
Hybrid	24.4	73.4	57.4	84.2	45.0	80.9

RAP-Gen can repair more bugs than T5-large especially for patches with 40 to 60 tokens.

5.3 RQ3: Comparative study with DL-based APR models on Code Refinement

We report the comparison results on Code Refinement in Table 6. All baseline results are directly obtained from their original papers. We first observe that “Naive Copy” gives a pretty high BLEU-4 score but with a zero exact match, indicating the buggy code and its fix has a large overlap and exact match should be employed as the primary metric. Among the baselines, NSEdit is a very competitive one with a best result (24.04 EM) on the small subset and CodeT5 gives the best result (13.96 EM) on the medium set. The lower results on the medium set compared to the small set indicates that longer buggy functions are more difficult to fix, which is aligned with observations in Fig. 5 (a). For RAP-Gen, various retrievers consistently boost the performance over their CodeT5 counterparts. Our best model with the hybrid retriever establishes new SoTA results on two subsets with 24.80 EM for small set and 15.84 EM for medium set. This again confirms that retrieved fix patterns provide helpful signals to guide the program repair and the hybrid retriever is more robust by using both lexical and semantic information. Fig. 3 (b) shows one case where the retrieved fix pattern (error line removal) helps RAP-Gen to successfully fix the bug.

5.4 RQ4: Analysis of Hybrid Patch Retriever

5.4.1 Evaluation of retrievers. We analyze the performance of our retrievers in terms of lexical and semantic matching between the query and the top retrieved patches. We employ the BLEU-4 score to measure their subtoken overlap for lexical matching, while for semantic matching, we compute the cosine similarity (CosSim) between their dense vectors encoded by our fine-tuned DPR retriever. Table 7 shows the performance of our retrievers on both TFix and Code Refinement benchmarks. The first row indicates the lower-bound performance via randomly retrieving bug-fix pairs from the codebase, where we observe this Random baseline achieves much lower scores in both lexical and semantic matching.

For lexical matching, BM25 outperforms DPR on TFix but underperforms on two Code Refinement subsets. We anticipate that it is due to the data difference between TFix and Code Refinement, where the latter employs obfuscated identifiers (e.g., VAR1, VAR2, ...) that hinders the performance of the lexical-based BM25 retriever. The hybrid retriever achieves the best lexical matching on all datasets, revealing the semantic information can complement to the lexical information. For semantic matching, DPR achieves

Table 8: Comparison with DL-based APR models on Defects4J single-line bugs with perfect fault localization.

Model	Chart	Closure	Lang	Math	Mockito	Time	Total
CODIT [4]	3	2	3	4	0	0	12
SequenceR [6]	3	3	2	6	0	0	14
DLFix [29]	5	7	6	11	1	1	31
CoCoNuT [37]	5	7	6	10	4	1	33
DEAR [30]	-	-	-	-	-	-	33
CURE [19]	-	-	-	-	-	-	38
RewardRepair [66]	5	8	6	16	2	1	38
CodeT5	4	6	7	13	3	0	33
RAG-Gen	6	10	8	16	3	1	44

In the table cells, it represents the number of correct patches.
 ‘-’ indicates data unavailability.

the best results on all datasets, which is not surprising as it is optimized towards the identical objective. Notably, our hybrid retriever achieves slightly lower results than DPR but much better results than BM25, implying it can balance both lexical and semantic information and be more robust than the lexical-based retrievers, which are sensitive to the choices of identifier naming.

5.5 RQ5: Comparative study with DL-based APR models on Defects4J

Table 8 shows the comparison results with seven DL-based APR baselines on fixing single-line bugs for six open source projects on Defects4J (v1.2). Except the results of CURE [19] that are obtained from DEAR [30], all other results are from the original papers or reproduced based on their released correct predictions. For fair comparison, we adopt the same 5-hour timeout as these baselines.

From Table 8, we observe RAP-gen fixes 11 more bugs than CodeT5’s 33 correct fixes with the help of retrieval, setting a new SoTA result of 44 compared to 38 correct fixes from CURE and RewardRepair that use a larger candidate pool than us (CURE’s 10,000 and RewardRepair’s 200 vs. our 100). Compared to RAP-Gen, RewardRepair requires another training process based on program compilation and test execution information. Besides, we investigate to what extent RAP-Gen can complement to existing APR models including DLFix [29], CoCoNuT [37], and RewardRepair [66] which have released their correct fix patches. Fig. 6 shows that RAP-gen repairs 6 more unique bugs compared to other baselines. Specifically, it fixes 17, 19 and 12 more unique bugs compared to DLFix, CoCoNuT, and RewardRepair. This reveals that RAP-Gen can complement to other DL-based APR systems for Defects4J bugs. We further show a case in Fig. 3 (c) and find that RAP-Gen successfully fix the Chart-9 bug but in a different form with the developer’s fix.

6 THREATS TO VALIDITY

Construct Validity. We evaluated RAP-Gen on three APR benchmarks: TFix in JavaScript, Code Refinement and Defects4J in Java. On TFix, we spotted a duplication issue and removed 243 intra-split or inter-split duplicates out of total 104,804 data instances. This might slightly impact the comparison between our model and the TFix (T5-large) model. We mitigate this threat by reporting the results of our model on the original TFix dataset and also the results of TFix model on the deduplicated test set. On Code Refinement, unlike the pairs in TFix can be validated by a static analyzer, its

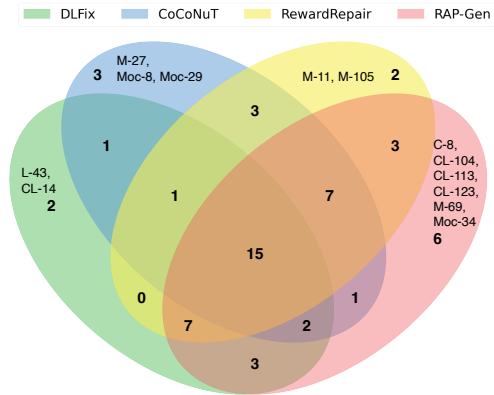


Figure 6: Overlapping analysis for six open source projects. C: Chart, CL: Closure, L: Lang, M: Math, Moc: Mockito.

bug-fix pairs are curated from GitHub commits with a bug-fix related commit message, and only a portion of them are manually verified [57]. There is a chance that some pairs are invalid (not related to the bug fix), which brings potential threats to the reliability of the evaluation on this dataset.

Internal Validity. The threats to internal validity mainly lie in the hyper-parameter search stage for RAP-Gen. As a neural model, its performance is highly affected by the choice of hyper-parameters. To alleviate such threats, we conduct a grid search to tune a better set of hyper-parameters but we still cannot claim they are the best.

External Validity. We only evaluated our RAP-Gen model on JavaScript and Java programs and do not study its generalization to other programming languages (PLs). However, our approach is language-agnostic as we do not employ any code-specific features like ASTs and can be applied in a drop-in fashion to other PLs. Besides, our evaluation on three APR datasets in two PLs should be comprehensive enough to verify the effectiveness of our approach.

7 CONCLUSION

We present a novel retrieval-augmented patch generation (RAP-Gen) framework for automatic program repair, a fundamental task in software engineering to reduce developers’ manual efforts in debugging. RAP-Gen consists of two components: a hybrid patch retriever to retrieve relevant fix patterns for a query buggy patch and a patch generator to synthesize the fixed patch based on both buggy patch and its retrieved guiding fix patterns. In addition, we propose to leverage a powerful code-aware pretrained language model CodeT5 as the backbone of RAP-Gen to facilitate both patch retrieval and generation in a unified manner. Comprehensive results on three diverse APR benchmarks in JavaScript and Java have demonstrated the effectiveness and superiority of our RAP-Gen model over existing deep learning-based APR approaches.

8 DATA AVAILABILITY

Our code and models can be found in this link (<https://figshare.com/s/a4e95baee01bba14bf4b>) and will be open-sourced upon acceptance to facilitate the development of automatic program repair.

REFERENCES

- [1] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified Pre-training for Program Understanding and Generation. In *NAACL-HLT Association for Computational Linguistics*, 2655–2668.
- [2] Earl T. Barr, Yuriy Brun, Premkumar T. Devanbu, Mark Harman, and Federica Sarro. 2014. The plastic surgery hypothesis. In *SIGSOFT FSE*. ACM, 306–317.
- [3] Berkay Berabi, Jingxuan He, Veselin Raychev, and Martin T. Vechev. 2021. TFix: Learning to Fix Coding Errors with a Text-to-Text Transformer. In *ICML (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 780–791.
- [4] Saikat Chakraborty, Yangruibo Ding, Miltiadis Allamanis, and Baishakhi Ray. 2022. CODIT: Code Editing With Tree-Based Neural Models. *IEEE Trans. Software Eng.* 48, 4 (2022), 1385–1399.
- [5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgren Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. *CoRR abs/2107.03374* (2021).
- [6] Zimin Chen, Steve Kommrusch, Michele Tufano, Louis-Noël Pouchet, Denys Poshyvanyk, and Martin Monperrus. 2021. SequenceR: Sequence-to-Sequence Learning for End-to-End Program Repair. *IEEE Trans. Software Eng.* 47, 9 (2021), 1943–1959.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT (1)*. Association for Computational Linguistics, 4171–4186.
- [8] Thomas Durieux, Youssef Hamadi, and Monperrus Martin. 2017. Production-Driven Patch Generation. *2017 IEEE/ACM 39th International Conference on Software Engineering: New Ideas and Emerging Technologies Results Track (ICSE-NIER)* (2017), 23–26.
- [9] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiao Cheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. In *EMNLP (Findings of ACL, Vol. EMNLP 2020)*. Association for Computational Linguistics, 1536–1547.
- [10] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. Convolutional Sequence to Sequence Learning. In *ICML (Proceedings of Machine Learning Research, Vol. 70)*. PMLR, 1243–1252.
- [11] Claire Le Goues, ThanhVu Nguyen, Stephanie Forrest, and Westley Weimer. 2012. GenProg: A Generic Method for Automatic Software Repair. *IEEE Trans. Software Eng.* 38, 1 (2012), 54–72.
- [12] Jiatao Gu, Yong Wang, Kyunghyun Cho, and Victor O. K. Li. 2018. Search Engine Guided Neural Machine Translation. In *AAAI*. AAAI Press, 5133–5140.
- [13] 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 *ICLR*. OpenReview.net.
- [14] Tatsunori B. Hashimoto, Kelvin Guu, Yonatan Oren, and Percy Liang. 2018. A Retrieve-and-Edit Framework for Predicting Structured Outputs. In *NeurIPS*. 10073–10083.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [16] Yaojie Hu, Xingjian Shi, Qiang Zhou, and Lee Pike. 2022. Fix Bugs with Transformer through a Neural-Symbolic Edit Grammar. *CoRR abs/2204.06643* (2022).
- [17] Gautier Izacard and Edouard Grave. 2021. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. In *EACL*. Association for Computational Linguistics, 874–880.
- [18] Jiajun Jiang, Yingfei Xiong, Hongyu Zhang, Qing Gao, and Xiangqun Chen. 2018. Shaping program repair space with existing patches and similar code. In *ISSTA*. ACM, 298–309.
- [19] Nan Jiang, Thibaud Lutellier, and Lin Tan. 2021. CURE: Code-Aware Neural Machine Translation for Automatic Program Repair. In *ICSE*. IEEE, 1161–1173.
- [20] Harshit Joshi, José Pablo Cambronero Sánchez, Sumit Gulwani, Vu Le, Ivan Radicek, and Gust Verbruggen. 2022. Repair Is Nearly Generation: Multilingual Program Repair with LLMs. *CoRR abs/2208.11640* (2022).
- [21] René Just, Dariouh Jalali, and Michael D. Ernst. 2014. Defects4J: a database of existing faults to enable controlled testing studies for Java programs. In *ISSTA*. ACM, 437–440.
- [22] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *EMNLP (1)*. Association for Computational Linguistics, 6769–6781.
- [23] Dongsun Kim, Jaechang Nam, Jaewoo Song, and Sunghun Kim. 2013. Automatic patch generation learned from human-written patches. In *ICSE*. IEEE Computer Society, 802–811.
- [24] Xuan-Bach Dinh Le, David Lo, and Claire Le Goues. 2016. History Driven Program Repair. In *SANER*. IEEE Computer Society, 213–224.
- [25] Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rottäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *NeurIPS*.
- [26] Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022. A Survey on Retrieval-Augmented Text Generation. *CoRR abs/2202.01110* (2022).
- [27] Jia Li, Yongmin Li, Ge Li, Xing Hu, Xin Xia, and Zhi Jin. 2021. EditSum: A Retrieve-and-Edit Framework for Source Code Summarization. In *ASE*. IEEE, 155–166.
- [28] Jian Li, Yue Wang, Michael R. Lyu, and Irwin King. 2018. Code Completion with Neural Attention and Pointer Networks. In *IJCAL*. ijcai.org, 4159–4165.
- [29] Yi Li, Shaohua Wang, and Tien N. Nguyen. 2020. DLFix: context-based code transformation learning for automated program repair. In *ICSE*. ACM, 602–614.
- [30] Yi Li, Shaohua Wang, and Tien N. Nguyen. 2022. DEAR: A Novel Deep Learning-based Approach for Automated Program Repair. In *ICSE*. ACM, 511–523.
- [31] Xuliang Liu and Hao Zhong. 2018. Mining stackoverflow for program repair. In *SANER*. IEEE Computer Society, 118–129.
- [32] Fan Long and Martin Rinard. 2015. Staged program repair with condition synthesis. In *ESEC/SIGSOFT FSE*. ACM, 166–178.
- [33] Fan Long and Martin C. Rinard. 2016. Automatic patch generation by learning correct code. *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (2016).
- [34] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *ICLR (Poster)*. OpenReview.net.
- [35] Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seung-won Hwang, and Alexey Svyatkovskiy. 2022. ReACC: A Retrieval-Augmented Code Completion Framework. *CoRR abs/2203.07722* (2022).
- [36] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin B. Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. 2021. CodeXGLUE: A Machine Learning Benchmark Dataset for Code Understanding and Generation. *CoRR abs/2102.04664* (2021).
- [37] Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan. 2020. CoCoNuT: combining context-aware neural translation models using ensemble for program repair. In *ISSTA*. ACM, 101–114.
- [38] Monperrus Martin, Simon Urli, Thomas Durieux, Matias Martinez, Benoît Baudry, and Lionel Seinturier. 2019. Repairator patches programs automatically. *Ubiquity 2019* (2019), 1–12.
- [39] Matias Martinez, Westley Weimer, and Monperrus Martin. 2014. Do the fix ingredients already exist? an empirical inquiry into the redundancy assumptions of program repair approaches. *Companion Proceedings of the 36th International Conference on Software Engineering* (2014).
- [40] Kivanç Muslu, Yuriy Brun, Reid Holmes, Michael D. Ernst, and David Notkin. 2012. Speculative analysis of integrated development environment recommendations. In *OOPSLA*. ACM, 669–682.
- [41] Md. Rizwan Parvez, Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval Augmented Code Generation and Summarization. In *EMNLP (Findings)*. Association for Computational Linguistics, 2719–2734.
- [42] Jeff H. Perkins, Sunghun Kim, Samuel Larsen, Saman P. Amarasinghe, Jonathan Bachrach, Michael Carbin, Carlos Pacheco, Frank Sherwood, Stelios Sidiroglou, Greg Sullivan, W. Wong, Yoav Zibin, Michael D. Ernst, and Martin C. Rinard. 2009. Automatically patching errors in deployed software. In *SOSP '09*.
- [43] Long N. Phan, Hieu Tran, Daniel Le, Hieu Nguyen, James T. Anibal, Alec Peltekian, and Yanfang Ye. 2021. CoText: Multi-task Learning with Code-Text Transformer. *CoRR abs/2105.08645* (2021).
- [44] Strategic Planning. 2002. The economic impacts of inadequate infrastructure for software testing. *National Institute of Standards and Technology* (2002), 1.
- [45] Julian Aron Preiner, Hlib Babii, and Romain Robbes. 2022. Can OpenAI’s Codex Fix Bugs?: An evaluation on QuixBugs. In *APR@ICSE*. IEEE, 69–75.
- [46] Yuhua Qi, Xiaoguang Mao, Yan Lei, Ziyang Dai, and Chengsong Wang. 2014. The strength of random search on automated program repair. In *ICSE*. ACM, 254–265.
- [47] Zichao Qi, Fan Long, Sara Achour, and Martin C. Rinard. 2015. An analysis of patch plausibility and correctness for generate-and-validate patch generation systems. In *ISSTA*. ACM, 24–36.
- [48] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

- 1277 [49] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang,
1278 Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the
1279 Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach.
1280 Learn. Res.* 21 (2020), 140:1–140:67.
- 1281 [50] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance
1282 Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 3, 4 (2009), 333–389.
- 1283 [51] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural Machine
1284 Translation of Rare Words with Subword Units. In *ACL (1)*. The Association for
1285 Computer Linguistics.
- 1286 [52] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning
1287 with neural networks. *Advances in neural information processing systems* 27
1288 (2014).
- 1289 [53] Shin Hwei Tan and Abhik Roychoudhury. 2015. relifix: Automated repair of
1290 software regressions. In *2015 IEEE/ACM 37th IEEE International Conference on
1291 Software Engineering*, Vol. 1. IEEE, 471–482.
- 1292 [54] Shin Hwei Tan, Hiroaki Yoshida, Mukul R. Prasad, and Abhik Roychoudhury.
1293 2016. Anti-patterns in search-based program repair. In *SIGSOFT FSE*. ACM,
1294 727–738.
- 1295 [55] Kristín Fjóra Tómasdóttir, Mauricio Finavaro Aniche, and Arie van Deursen.
1296 2017. Why and how JavaScript developers use linters. In *ASE*. IEEE Computer
1297 Society, 578–589.
- 1298 [56] Nikzad Benny Toomarian and Jacob Barhen. 1992. Learning a trajectory using
1299 adjoint functions and teacher forcing. *Neural Networks* 5, 3 (1992), 473–484.
- 1300 [57] Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin
1301 White, and Denys Poshyvanyk. 2019. An Empirical Study on Learning Bug-
1302 Fixing Patches in the Wild via Neural Machine Translation. *ACM Trans. Softw.
1303 Eng. Methodol.* 28, 4 (2019), 19:1–19:29.
- 1304 [58] Simon Urli, Zhongxing Yu, Lionel Seinturier, and Monperrus Martin. 2017. How
1305 to Design a Program Repair Bot? Insights from the Repairator Project. *2018
1306 IEEE/ACM 40th International Conference on Software Engineering: Software Engi-
1307 neering in Practice Track (ICSE-SEIP) (2017)*, 95–104.
- 1308 [59] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning
1309 with Contrastive Predictive Coding. *CoRR* abs/1807.03748 (2018).
- 1310 [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones,
1311 Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All
1312 you Need. In *NIPS*. 5998–6008.
- 1313 [61] Yue Wang, Weishi Wang, Shafiq R. Joty, and Steven C. H. Hoi. 2021. CodeT5:
1314 Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Under-
1315 standing and Generation. In *EMNLP (1)*. Association for Computational Linguis-
1316 tics, 8696–8708.
- 1317 [62] Westley Weimer, ThanhVu Nguyen, Claire Le Goues, and Stephanie Forrest.
1318 2009. Automatically finding patches using genetic programming. *2009 IEEE 31st
1319 International Conference on Software Engineering (2009)*, 364–374.
- 1320 [63] Cathrin Weiß, Rahul Premraj, Thomas Zimmermann, and Andreas Zeller. 2007.
1321 How Long Will It Take to Fix This Bug?. In *MSR*. IEEE Computer Society, 1.
- 1322 [64] Ming Wen, Junjie Chen, Rongxin Wu, Dan Hao, and Shing-Chi Cheung. 2018.
1323 Context-aware patch generation for better automated program repair. In *ICSE*.
1324 ACM, 1–11.
- 1325 [65] Martin White, Michele Tufano, Matias Martinez, Monperrus Martin, and Denys
1326 Poshyvanyk. 2019. Sorting and Transforming Program Repair Ingredients via
1327 Deep Learning Code Similarities. *2019 IEEE 26th International Conference on
1328 Software Analysis, Evolution and Reengineering (SANER) (2019)*, 479–490.
- 1329 [66] He Ye, Matias Martinez, and Martin Monperrus. 2022. Neural Program Repair
1330 with Execution-based Backpropagation. In *ICSE*. ACM, 1506–1518.
- 1331 [67] Yuan Yuan and Wolfgang Banzhaf. 2020. ARJA: Automated Repair of Java
1332 Programs via Multi-Objective Genetic Programming. *IEEE Trans. Software Eng.*
1333 46, 10 (2020), 1040–1067.
- 1334 [68] Qihao Zhu, Zeyu Sun, Yuan-an Xiao, Wenjie Zhang, Kang Yuan, Yingfei Xiong,
1335 and Lu Zhang. 2021. A syntax-guided edit decoder for neural program repair. In
1336 *ESEC/SIGSOFT FSE*. ACM, 341–353.
- 1337
- 1338
- 1339
- 1340
- 1341
- 1342
- 1343
- 1344
- 1345
- 1346
- 1347
- 1348
- 1349
- 1350
- 1351
- 1352
- 1353
- 1354
- 1355
- 1356
- 1357
- 1358
- 1359
- 1360
- 1361
- 1362
- 1363
- 1364
- 1365
- 1366
- 1367
- 1368
- 1369
- 1370
- 1371
- 1372
- 1373
- 1374
- 1375
- 1376
- 1377
- 1378
- 1379
- 1380
- 1381
- 1382
- 1383
- 1384
- 1385
- 1386
- 1387
- 1388
- 1389
- 1390
- 1391
- 1392