

PYEVOLVE: Automating Frequent Code Changes in Python ML Systems

Malinda Dilhara
University of Colorado Boulder
USA
malinda.malwala@colorado.edu

Danny Dig
University of Colorado Boulder
USA
danny.dig@colorado.edu

Ameya Ketkar
Uber Technologies Inc.
USA
ketkara@uber.com

Abstract—Because of the naturalness of software and the rapid evolution of Machine Learning (ML) techniques, frequently repeated code change patterns (*CPATs*) occur often. They range from simple API migrations to changes involving several complex control structures such as `for` loops. While manually performing *CPATs* is tedious, the current state-of-the-art techniques for inferring transformation rules are not advanced enough to handle unseen variants of complex *CPATs*, resulting in a low recall rate. In this paper we present a novel, automated workflow that mines *CPATs*, infers the transformation rules, and then transplants them automatically to new target sites. We designed, implemented, evaluated and released this in a tool, PYEVOLVE. At its core is a novel *data-flow, control-flow* aware transformation rule inference engine. Our technique allows us to advance the state-of-the-art for transformation-by-example tools; without it, 70% of the code changes that PYEVOLVE transforms would not be possible to automate. Our thorough empirical evaluation of over 40,000 transformations shows 97% precision and 94% recall. By accepting 90% of *CPATs* generated by PYEVOLVE in famous open-source projects, developers confirmed its changes are useful.

I. INTRODUCTION

The naturalness of software [24], [25], [29], [34] leads to repeated code changes, both within and across projects. Because programmers employ the same coding idioms [26], [28] and best practices (e.g., linting [81]), they change code similarly, resulting in repeated code changes. These repeated changes are fine-grained, recurring at the method level, and have the same semantics.

Listing 1 shows an example of such a change in project *NifTK/NiftyNet*, an open-source convolutional neural network platform. The developers replaced a `for` loop that sums the list `elements` with `numpy.sum`, a best-practice for improving performance. This change is fine-grained because it involves programming idioms [27], [28] and occurs within a specific method. Since it is repeated at multiple locations in multiple commits, it is a *fine-grained code change pattern* (*CPAT*).

Listing 1: Commit c8b28432 in GitHub repository *NifTK/NiftyNet*: Replace `for` loop with `NumPy sum`

```
1 -result = 0
2 -for elem in elements:
3 -     result = elem + result
4 +result = numpy.sum(elements)
```

Writing program transformations to automate a huge variety of such changes (e.g., more than 28,000 kinds as shown in [21]) is difficult as research has shown [62]–[80]. There are

many reasons. (i) AST rewriting poses a significant barrier to entry, (ii) Matching control/data flows for real code is tedious to develop since there is too much noise and syntactic variance to account for, (iii) Moreover, as these best practices evolve, they are difficult to maintain, thus it requires a community effort. To overcome these challenges, in this paper we rely on a unique insight: if we mine the many examples of *CPATs* in the open-source community and infer transformation rules, we can feed them to automated program transformation systems, with no burden on the developer. Thus, we present an end-to-end solution, PYEVOLVE, which mines common changes from the open-source community, it automatically infers the transformation rules in a *data- and control-flow* manner, and transplants and applies them at new target sites. We show its effectiveness by applying it on Python ML projects.

According to the GitHub 2021 annual report [42] of programming language usage, Python is one of the top two most used. Moreover, it has become the lingua-franca for machine learning (ML) and data science development [8]–[12], and in many STEM fields the developers are not trained on software engineering best practices. Despite its prominence and community needs, tools for evolving Python code are significantly behind other languages.

To enhance tooling in Python, recently, researchers developed tools to mine *CPATs* in software systems. R-CPATMINER [21] is one such tool for mining *CPATs* in Python systems. The authors of the tool further conducted a large-scale empirical study on *CPATs* in 1,000 Python ML systems and found that the complexity of *CPATs* in ML systems ranged from basic API migrations to changes involving complex control structures like `for` loops. Developers perform *CPATs* for several reasons: *performance* (e.g., `for` loop \rightarrow `vectorization`), *using advanced language features* (e.g., `for` loop \rightarrow `list comprehension`), *better resource management* (e.g., using `with` statement), and *library migrations* (e.g., `numpy.mean()` \rightarrow `torch.mean()`). Manually searching and applying variations of such changes to several locations is error-prone. Moreover, developers might overlook sites that require the same edit. Indeed, over 75% of the respondents of a survey [21] sent to 650 developers indicated they needed these *CPATs* to be automated. However, the existing *CPAT* automation tools are not yet able to handle them.

Despite the existence of many program transformation sys-

tems [62]–[80], their main impediment to adoption is the need for programmers to write sophisticated program transformation rules. In recent years, we have seen an emerging trend of tools and techniques that infer transformation rules using example code edits [1], [15]–[19], [33], [48]–[52], [54], [84], [85]. These techniques infer transformation rules from before and after edits of human-adaptations, then use the inferred rules to transform target codes. This is called “*Transformation by Example*”. Despite the potential of such techniques to significantly ease code evolution, they have so far been primarily used for API migrations, such as replacing obsolete API calls with modern ones from the Android SDK [1], [5], [15], [17], [19], [54], the Linux kernel [40], [41], and Type migration of Java systems [33]. Although existing techniques work well when replacing an API call with another, they under-perform on more complex coding idioms such as the one in Listing 1. In diverse codebases, this *CPAT* will have many variations in terms of data- and control-flow.

Listing 2: The repository hachmannlab/chemm uses a `for` loop to compute the sum of an array

```

1 n_diff = 0
2 to_eval = getarray()
3 for dif in to_eval.getDif():
4     total = n_diff + dif
5     n_diff = total

```

For example, Listing 2 is semantically equivalent to Listing 1 but differs in *data- and control-flow* due to assignments and how the accumulator variable computes the sum. Therefore, existing “Transformation by Example” techniques struggle because the rule they inferred for Listing 1 cannot be used to transform the target code in Listing 2, resulting in a low recall rate. Because existing techniques are so syntax-centric, the transformation rules are prone to overfitting to the input examples. That is, the transformation rule may work well on the given examples but may be unable or erroneously transform unseen *data- or control flow* variants outside of the examples. This demonstrates a major limitation of current transformation rule inference methodologies.

To overcome these challenges, in this paper we present a novel *data- and control-flow* aware technique that infers transformation rules and adapts them to transform *even unseen variations* in the target codes. Our novel technique enables automating even unseen variations of the *CPAT* by preserving *data- and control-flow* relations. In this paper, we present a fully end-to-end pipeline that mines and automates Python *CPATs*. Our pipeline consists of four major steps: (i) mine *CPATs* from version histories, (ii) infer transformation rules, (iii) identify the new sites to apply the *CPATs*, (iv) adapt the transformation rules to the new sites. To do so, we leverage and further extend four state-of-the-art techniques: (i) to mine *CPATs* we use *R-CPATMiner* [21], [22], (ii) to infer initial transformation rules from example instances in *CPATs* we use *InferRule* [33] (iii) to identify new sites to apply *CPATs* we use fine-grained program dependence graphs *fgPDG* [34]; and (iv) to apply the *CPATs* at the new sites, even for unseen variants, we re-infer the transformation rules based

on *data- and control-flow* relations in the *fgPDG*. Lastly, we use *ComBy* [32] to declaratively rewrite programs according to the re-inferred transformation rules.

We implemented our novel technique in *PYEVOLVE*. We evaluated its effectiveness and usefulness on a corpus of *CPATs* that had previously been shown to be diverse in terms of size, frequency, authors, and projects [21]. We conducted replication case studies comprising a broad variety of 40,000 transformation trials. Using cross-validation, we tested *PYEVOLVE*’s ability to correctly transform *CPATs*. We found that *PYEVOLVE* achieves 97% precision and 94% recall when replicating developer-performed changes. In addition, we discovered that 70% of these changes are *data or control-flow* variants that cannot be automated using existing techniques, thus *PYEVOLVE* advances the state-of-the-art significantly. To evaluate *PYEVOLVE*’s usefulness, we submitted pull requests to highly-rated, best-in-class projects such as *TensorFlow*, *PyTorch*, *Scikit-Image*, and *Keras*, totaling 181 *CPAT* instances. At the time of this writing, their developers have already accepted 163 (90%) *CPATs*.

This paper makes the following contributions:

- (1) We introduce a novel *data- and control-flow* aware rule inference that effectively transforms even unseen variants that cannot be handled by existing rule inference techniques.
- (2) We designed and implemented our technique in *PYEVOLVE*. It mines *CPATs* from projects, infers transformation rules, and generates patches for Python projects. To the best of our knowledge, this is the first such pipeline developed for Python and it assists ML developers and other Python developers in keeping up with rapidly evolving best practices.
- (3) Our empirical evaluation of *PYEVOLVE* on 40,000 transformations shows that *PYEVOLVE* is *effective* (97% precision, 94% recall), *needed* (it enables automating 70% more *CPATs* than existing tools), and its patches are *useful* (developers accepted 90% of 181 *PYEVOLVE*-generated *CPATs*).
- (4) Our tool and evaluation dataset is open-source and available for others to reuse [44].

II. MOTIVATING EXAMPLES

To illustrate the challenges of using existing “Transformation by Example” techniques we use the real-world code changes, shown in Table I. The first column (Code Before→Code After) shows the code fragments before and after the code change, while the second column (Rule) presents the rules encoding the code change using the syntax of *Comby* [20], a state-of-the-art code rewrite tool. The third column shows the Guards for each rule in column-2. The last column shows examples of new target sites that we would like to transform. First, we describe a success scenario that existing techniques automate with a high recall rate, followed by three scenarios that they fail to automate (as the complexity of the code rises), but *PYEVOLVE* succeeds.

Table I, row 1 depicts an example of a *CPAT* mined from *TensorFlow*. To open a file, the code changes from `open`, which is native to Python, to `GFile` in *TensorFlow*. This adaptation is represented by the rewrite rule `:[v1]=open(:[v0]) →`

Table I: Motivating Examples

	Code Before → Code After	Rule	Guard	New Target Site
1.	<pre>f=open("f.csv") → f=tf.gfile.GFile("f.csv")</pre>	<pre>:[v1]=open(:[v0]) → :[v1]=tf.gfile.GFile(:[v0])</pre>	<pre>type=::[v1]->TextIO type=::[v1]->str</pre>	<pre>f=open("data.csv")</pre>
2.	<pre>X=numpy.dot(A,B) Y=numpy.dot(X,C) → Y=numpy.linalg .multi_dot([A,B,C])</pre>	<pre>:[v3]=:[v6]].dot (::[v1],:[v2]) :[v5]=:[v6]].dot (::[v3],:[v4]) → :[v5]=:[v6]].linalg.multi_dot (::[v1],:[v2],:[v3])</pre>	<pre>type=::[v1]->ndarray type=::[v2]->ndarray type=::[v3]->ndarray type=::[v4]->ndarray type=::[v5]->ndarray import=::[v6]->numpy</pre>	<pre>prod = numpy.dot(numpy.dot(A,B),C)</pre>
3.	<pre>olderr = np. seterr(divide ='ignore') try: actual=logit(a) finally: np.seterr(olderr) → with np.errstate(divide='ignore'): actual = logit(a)</pre>	<pre>:[l1]=:[l3].seterr(divide='ignore') try: :[l12] finally: :[l3].seterr(:[l1]) → with :[l3].errstate(divide='ignore'): :[l12]</pre>	<pre>import=::[l3] -> numpy</pre>	<pre>olderr = np. seterr(divide='ignore') computeSum([2,4,6]) try: actual = logit(a) finally: np.seterr(olderr)</pre>
4.	<pre>res = 0 for elem in elems: res = res + elem → res = numpy.sum(elems)</pre>	<pre>:[v0] = 0 for :[v1] in :[v2]: :[v0] = :[v0] + :[v1] → :[v0]=numpy.sum(:[v2])</pre>	<pre>type=::[v0]->int type=::[v1]->int</pre>	<pre>n_diff = 0 to_eval = getarray() for dif in to_eval.getDiff(): total = n_diff + dif n_diff = total</pre>

`:[v1]=tensorflow.gfile.GFile(:[v0])`. This rule transforms any target code similar to the before- change shown in column-1-Table I. The rule has a left side (indicating the “before” the change) and a right side (indicating the “after” the change) separated by an arrow. The left side contains Python statements with template variables (e.g., `:[v0]`) that bind to AST nodes from the actual source code (e.g., `v0` binds to `file.csv`). The right side of the rule also contains Python statements with template variables, where each template variable denotes the code fragments that will be used after the change is applied. Because the change in row 1 swaps one API invocation (`open`) for another (`tf.gfile.GFile`), the likelihood of finding semantically equivalent different variants (in terms of *data-flow* or *control-flow*) in new target sites is very low. Therefore, the example transformation (in columns 1&2) is sufficient to represent many variants and would be applicable to many new target sites. For these API migrations, the existing Program-by-Example techniques [15], [17]–[19], [33], [48]–[52] perform well (i.e., they achieve high recall rate).

However, many real-world *CPATs* involve multiple method invocations [21]. For example, in Table I row 2 (taken from *Scikit-learn*), developers replaced two `numpy.dot` calls that perform matrix multiplication with `multi_dot([A,B,C])`, making the code concise and efficient. We show the relevant inferred rule for this change in column 3-Table I. Now consider another site within the same *Scikit-learn* project, shown in column 4. While this new site similarly does matrix multiplication using two `numpy.dot` calls, it does so in a data flow variant by passing the result of the first `numpy.dot` directly into the second `numpy.dot`. The rule in column 3 cannot be used to transform this new target site, and because this variant was not seen during rule inference, existing rule inference techniques fail

to automate it, but PYEVOLVE changes it correctly.

Real-world *CPATs* are large in size and involve many control structures [21]. For example, as shown in row 3, in the project *SciPy*, developers transform error-ignoring code that uses `try finally` to a `with` statement. The new target site in column-4, on the other hand, has an extra method invocation, making it a control-flow variant. Now let us consider the example in row 4, which also appeared in Section I. The target site in column-4 has an additional assignment statement and variable assignment, which makes it both data- and control-flow variant. Hence, since these variants were not seen as change exemplars during rule inference, existing “Transformation by Example” techniques fail to change these new target sites, but PYEVOLVE changes them correctly.

We observed many cases where existing techniques inferred rules that are prone to over-fitting the examples and are unable to account for even minor variants. These variants are frequently employed in ML code bases that perform advanced numerical computations. Perhaps these variants make it easier to debug. For example, even though the code in Listing 2 uses a redundant variable `total` (lines 4 and 5 can be easily combined `n_diff=n_diff+ dif`), a developer can still use the variable `total` by inserting a debug point in Line 4 to test the addition of two consecutive array elements. Despite the abundance of such variants, existing “Transformation by Example” techniques fail to automate new target sites that are semantically-equivalent but syntactically-different, thus they fail to transform good candidates. PYEVOLVE provides an end-to-end solution for mining and automating code transformations in a data- and control-flow aware manner, therefore automating even previously unseen variants of changes.

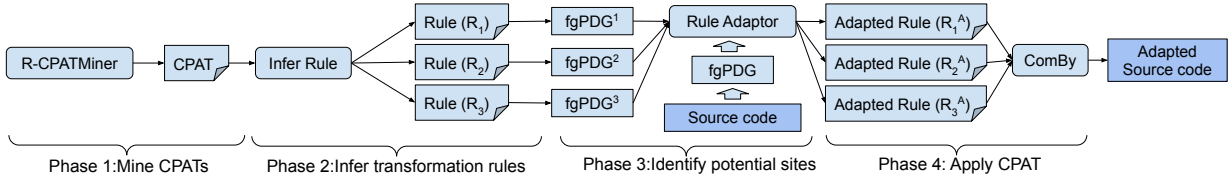


Figure 1: The high-level overview of the pipeline for mining and applying a *CPAT*

III. TECHNIQUE

PYEVOLVE is our novel technique for inferring *CPATs* from other code bases and adapting and applying them to new target sites. Figure 1 provides the high-level overview of our technique, which consists of four major phases:

Phase1. PYEVOLVE mines *CPATs* from the version history of Python projects. To do so, we adapted *R-CPATMiner* [21], the state-of-the-art *CPAT* mining tool for Python systems.

Phase2. PYEVOLVE infers transformation rules from the previously mined *CPATs*. To do so, we adapted *InferRule* [33], which is a state-of-the-art tool for inferring transformation rules in Java (that we adapted to Python).

Phase3. PYEVOLVE identifies potential target sites in the code where to apply the previously mined *CPATs*. To do so, we use a graph-based matching technique (described in Section III-B3).

Phase4. PYEVOLVE adapts the rules inferred in Phase2 to the target sites and their contexts (see Section III-B4). This allows PYEVOLVE to be effective for a wider range of unseen *CPAT* variants. These adapted rules take into account the *data*- and *control*-flow relations. Finally, it applies them to the target sites observed in Phase3.

The key idea of making PYEVOLVE resilient to both seen and unseen variants relies on a graph-based matching technique that involves two main steps: (i) We first build graphs for the target code and initial transformation rule, then mine isomorphic sub-graphs to those of the rule in the target code. This makes it possible to identify potential sites for *CPATs*, taking into account both seen and unseen *data*- and *control* flow variants (Section III-B3). Then (ii) using the sub-graphs, we adapt the transformation rules to the target site and its context. This makes it possible to create rules that take into account the target code’s *data*- and *control* flow (Section III-B4).

The rest of this section is organized as: (i) Section III-A defines fundamental terms used throughout the paper, (ii) Section III-B1 describes how we mine *CPATs*, (iii) Section III-B2 describes our graph representation for the rule and target code, and (iv) Section III-B3 describes how to mine potential sites to apply *CPAT* by mining sub-graphs, finally (v) Section III-B4 how to infer the final transformation rule.

A. Basic Concepts

We now formally define the fundamental terms.

Definition III.1 (TEMPLATE EXPRESSION: \mathbb{T}). A generic and lightweight way of searching and matching syntactic structures in the AST of a function. In Python, a template expression is made up of Python expressions and statements, as well as template variables (or *holes*) that match to a program AST.

```

TRANSFORMATIONRULE ::= GUARD RULE
RULE ::= TEMPLATEEXPRESSION → TEMPLATEEXPRESSION
GUARD ::= PRED | PRED GUARD | ∅
PRED ::= ISKIND(TEMPLATEVARIABLE, KIND) |
        ISTYPE(TEMPLATEVARIABLE, TYPE) |
        ISEQUAL(TEMPLATEVARIABLE, VALUE) | NOT(PRED)

```

Figure 2: DSL for program transformation rule

Our template expressions adhere to the COMBY syntax [32] - a state-of-the-art, lightweight multi-language syntax transformation technique for declaratively rewriting syntax. Rules written in the COMBY language are human readable and modifiable by developers, who can further customize the changes. Thus, they will appeal to developers who might not like hard-coded refactoring tools that automate code changes [45], [82] without any intervention from the developer. Details of COMBY’s syntax can be found on its website [20].

Definition III.2 (TEMPLATE VARIABLE: \mathbb{V}). This is a hole in the template expression corresponding to one or more program elements. According to COMBY [20], $[n]$ binds the source code to a template variable n . A template variable can match all characters lazily up to its suffix (like $.^*?$ in regex) within its level of balanced delimiters. COMBY supports mainly two template variables: (i) $[v]$ matches an identifier, analogous to $\backslash w^+$ in regex; we denote it as (\mathbb{V}_I) , (ii) $[v]$ matches one or more alphanumeric characters and “_”. We denote it as (\mathbb{V}_A) .

Definition III.3 (RULE: $\mathbb{R}(\mathbb{T}_{LHS} \rightarrow \mathbb{T}_{RHS})$). This defines how to transform the input AST to an output AST using *template expressions*. In our setting, *template expression* (\mathbb{T}_{RHS}) on the right side of the *Rule* contains *template variable* (\mathbb{V}) that denote the substitution with an appropriate fragment of the program AST, as matched on the left side. For example, once the rule in row 4-Table I is applied to the code that sums the elements in Listing 1, it will be rewritten as `np.sum(elements)`.

1) Transformation Rules (TR): Modern “Transformation by Example” techniques employ transformation rules that correspond to the examples defined over a predefined Domain-Specific Language (DSL). PYEVOLVE inherits the DSL of *TCInfer* [33], a tool that infers rules for Java systems, and extends it to the language shown in Figure 2. In this DSL, a program transformation rule is a pair of *Guard* and *Rule*. Essentially, the *Guard* validates *which* code fragments should be transformed while the *Rule* describes *how* those code fragments should be transformed. The *Guard* is composed of a set of conjunctive predicates (*Pred*) on the attributes (e.g., Type, Kind, Value) of the template variables. The *Guard* evaluates where a template variable satisfies its predicate *s* and returns a Boolean value accordingly. Column 2 of Table I

shows the *Rule* inferred from the code changes in column 1, and column 3 shows the relevant *Guard*. For example, the *Rule* in row 1 describes how the method call `open` should be transformed (*Rule*), but only in the places where a string is passed as an argument and `TestIO` is returned (*Guard*).

2) Fine-Grained Program Dependence Graph (fgPDG):

In their recent work, Nguyen et al. [34] presented *fgPDG*, a sufficiently generalized program dependency graph that can be used to mine semantically equivalent program fragments with differing data-and control relations. Researchers utilize *fgPDGs* for many applications. For example, Nguyen et al. [34], Smirnov et al. [46], and Dilhara et al. [21], used *fgPDG* on Java and Python systems to mine semantically equivalent repeated code changes, while Noda et al. [47] used *fgPDG* to mine repeated bug fixes and repair unfixed similar bugs. Our key idea for determining if the target code contains an equivalent code for \mathbb{T}_{LHS} (see Definition III.3) regardless of data-and control-flow is to construct a *fgPDG* for the \mathbb{T}_{LHS} and determine if it is a subgraph of the target code’s *fgPDG*. Then, we adapt the inferred *rule* to match the target code. For example, to match the target code in Listing 2, the learned *rule* in row 4-Table I must be adapted to Listing 3.

Listing 3: Adapted transformation rule to match with Listing 2

```

1 :[[v0]] = 0          ⇒      :[v3]
2 :[v3]              :[[v0]]=numpy.sum(:[v2])
3 for :[[v1]] in :[v2]:
4     :[[v4]] = :[[v0]] + :[[v1]]
5     :[[v0]] = :[[v4]]

```

To achieve this, we extend *fgPDG* with two new nodes: *IdentifierHole* and *AlphaHole*.

Definition III.4 (IDENTIFIERHOLE: \mathbb{I}). $\mathbb{V}_{\mathbb{I}}(:[[v]])$ represents a variable *identifier* in code (analogous to `w+` in regex). To denote $\mathbb{V}_{\mathbb{I}}$ in *fgPDG*, we add the new node, *IdentifierHole*.

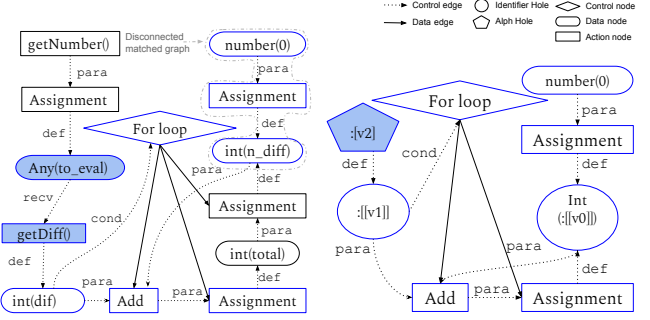
Definition III.5 (ALPHAHOLE: \mathbb{A}). $\mathbb{V}_{\mathbb{A}}(:[v])$ represents *expressions* (e.g., method calls – `np.dot()`, lists – `[1,2]`, and dictionaries – `{"one":1}`) or a *statement* (e.g., *assignments*). To denote $\mathbb{V}_{\mathbb{A}}$ in *fgPDG*, we use the new node, *AlphaHole*.

B. Data-flow Control-flow aware rule inference

PYEVOLVE generates transformation rules for applying *CPATs* to target code. First, it uses INFERRULE [33] to infer the transformation rule for the input *CPATs*, then adapts the rule based on the *data and control flow* in target codes, which we call data-and control-flow aware rule inference, formally:

Definition III.6 (DATA-AND CONTROL-FLOW AWARE RULE INFERENCE). For given input code changes $\{i_0 \rightarrow o_o, \dots, i_n \rightarrow o_n\}$, the existing “Transformation by Example” techniques infer transformation rule *TR* such that $TR_k(i_k)=o_k$ for $k \in \{0 \dots n\}$. *Control-and data-flow aware rule inference* generates an adapted TR_k^A for each TR_k that matches a target code.

1) *Input*: We use R-CPATMINER [21], a state-of-the-art tool to collect *CPATs* that other developers have performed in Python ML systems. In their empirical study [21] performed on diverse 2,500 *CPATs*, the authors further observed four



(a) $fgPDG$ of code Listing 2 ($fgPDG^T$) (b) $fgPDG$ of \mathbb{T}_{LHS} of row 4-Table I ($fgPDG^R$)

Figure 3: Sub-graph matching - nodes with blue edges are the matched sub-graph of $fgPDG^R$ in $fgPDG^T$ (G^M)

kinds of *CPATs* that frequently occur in Python ML systems, (i) *dissolve for loops to domain specific abstraction* (e.g., row-4-Table I) (ii) *update API usage* (e.g., row-2-Table I), (iii) *transform to context managers* (e.g., row-3-Table I), and (iv) *use advanced language features* (e.g., Python list comprehension). We use the same *CPATs* as input to PYEVOLVE.

2) **Generating *fgPDGs***: *fgPDG* is a directed graph consisting of three different types of nodes and two types of edges. (i) *Data Nodes* (\mathbb{N}^D) represent variables, field accesses, and constants. For example, the variable `n_diff` in Figure 3(a) is a data node. (ii) *Action Nodes* (\mathbb{N}^A) represent operations on data, e.g., array accesses, and method calls. For example, in Figure 3(a), the methods calls `getarray()` is an *Action node*, and (iii) *Control Nodes* (\mathbb{N}^C) represent control statements, such as `if` for branching, `for` for looping (see Figure 3(a)). (i) *Control edges* represent control relations between the statements/operations and the control nodes on which their executions depend, (ii) *Data edges* represent the data flow of *fgPDG* nodes. As shown in Figure 3(a), each edge has a label that specifies the type of data flow (e.g., `para`, `def`, `cond` [34]).

We first generate an *fgPDG* for the target code. For example, Figure 3(a) shows the an *fgPDG* generated for the target code in Listing 2. One of the difficulties in generating an *fgPDG* for Python codes is the lack of type information at compile time, which is critical for mining semantically equivalent code. To overcome this challenge, we employ *type inference*, a technique for inferring the type information of program elements based on data-flow information available at compile time. For this purpose, we use PYTYPE [37], developed by Google, which is widely used by the Python community.

Second, we generate an *fgPDG* for the left side of the transformation rule (\mathbb{T}_{LHS}). Figure 3(b) shows an example *fgPDG* generated for the \mathbb{T}_{LHS} in row 4-Table I. In addition to the three nodes used by Nguyen et al. [34], we added two more nodes to the *fgPDG* — *IdentifierHole* (see Definition III.4) and *AlphaHole* (see Definition III.5) — to generate the $FGPDG$ of \mathbb{T}_{LHS} . For example, the *fgPDG* nodes, `: [v2]` and `: [[v0]]` in Figure 3(b) are examples of *AlphaHole* and *IdentifierHole*.

3) **Identifying potential sites**: We now describe how we identify potential target code sites to apply *CPATs*. Conceptually, given the *fgPDG* of the target code ($fgPDG^T$) and the *fgPDG* of the \mathbb{T}_{LHS} ($fgPDG^R$), PYEVOLVE determines

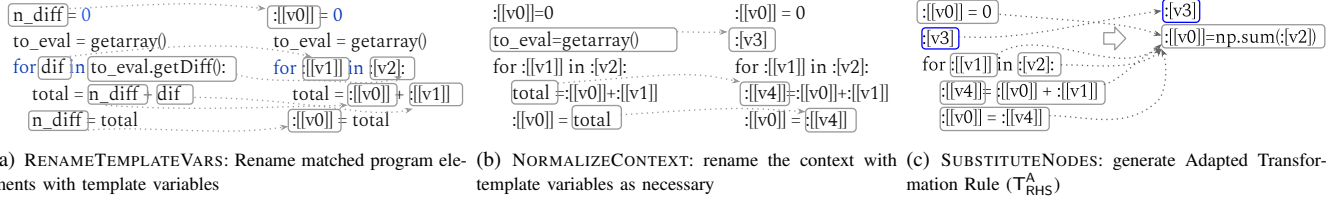


Figure 4: Steps in Algorithm 2 that are followed by Listing 2 and the rule in row 4-Table I to be adapted as a rule in Listing 3.

Algorithm 1 Match nodes from $fgPDG^T$ and $fgPDG^R$

```

1: function MATCHEDNODE(cNode, tNode)
2:   if NODETYPE(tNode) == IdentifierHole then
3:     if cNode.ASTNode == SimpleName then
4:       if GUARDSMEETS(cNode, tNode) then
5:         return True
6:   else if NODETYPE(tNode) == AlphaHole then
7:     if cNode.ASTNode == (Expression or Statement) then
8:       if GUARDSMEETS(cNode, tNode) then
9:         return True
10:  else if NODETYPE(cNode) == NODETYPE(tNode) then
11:    if NODETYPE(cNode) == ( $N^D$  Or  $N^C$  Or  $N^A$ ) then
12:      if GUARDSMEETS(cNode, tNode) then
13:        return True
14:  return False

```

whether the $fgPDG^R$ is a sub-graph of $fgPDG^T$ and if it is, it reports the matching nodes as the locations to apply the rule. This differs from typical sub-graph mining problems in two ways. (i) *AlphaHole* in $fgPDG^R$ can match against one or more nodes in $fgPDG^T$, and (ii) data nodes in $fgPDG^R$ can match transitively to data nodes in $fgPDG^T$ via *Data edges* [34]. Therefore, the matched sub-graph can be a disconnected graph.

Figure 3(a) ($fgPDG^T$) depicts the $fgPDG$ generated for the target code Listing 2, whereas Figure 3(b) ($fgPDG^R$) depicts the $fgPDG$ generated for the T_{LHS} of rule in row 4-Table I.

Algorithm 1 describes how we match two nodes from $fgPDG^R$ and $fgPDG^T$ based on node kinds and the related *Guard*. If both nodes are in same category, i.e., N^D , N^A , or N^C (see Section III-B2), the operation matches them based on Guards (line 12). For example, in Figure 3, the N^D , `number(0)` in $fgPDG^T$, matches the N^D , `number(0)`, in the $fgPDG^R$. The node kind *IdentifierHole* matches the identifiers, hence the operation determines if the target node's AST kind is *SimpleName* and whether the *Guard* matches (see 4). For example, the *IdentifierHole*, `:[v0]` in Figure 3(b) matches with the N^D , `n_diff`, which has `type->int` as the *Guard*. *AlphaHole*, matches with any expression (e.g., method invocation) or statement (e.g., for statement) if they meets the *Guard*. Hence, *AlphaHole* can match one or many nodes in $fgPDG^T$. The `:[v2]` in Figure 3(b) matches two nodes in the $fgPDG^T$ that are relevant to the method call `to_eval.getDiff()`.

We recursively walk through nodes and their child nodes, and then perform the operation `MATCHEDNODE` to obtain a matched $fgPDG$ (G^M), which is a sub-graph of the $fgPDG^T$. In order to match with data flow variants, we also check whether the Data nodes of $fgPDG^R$ are transitively matched through Data edges to Data nodes in $fgPDG^T$. For example, the target code in Listing 2 *adds* the two values, `n_diff`, and

Algorithm 2 Generate Adapted Rules

```

1: function GENERATEADAPTEDRULE( $T_{LHS}$ ,  $T_{RHS}$ ,  $G^M$ , targetCode)
2:    $LHS^A \leftarrow$  RENAME_TEMPLATE_VARS(targetCode,  $G^M$ )
3:    $T_{LHS}^A \leftarrow$  NORMALIZE_CONTEXT( $LHS^A$ )
4:    $T_{RHS}^A \leftarrow$  SUBSTITUTE_NODES( $T_{LHS}^A$ ,  $G^M$ )
5:   if IS_SAFE( $T_{LHS}^A \rightarrow T_{RHS}^A$ ) then
6:     return  $T_{LHS}^A \rightarrow T_{RHS}^A$ 
7: function SUBSTITUTE_NODES
8:    $T_{RHS}^A \leftarrow$  SUBSTITUTE( $T_{LHS}^A \cap T_{LHS}$  in  $T_{LHS}^A$  with  $T_{RHS}$ )

```

`diff`, together and assigns it to a variable `total`, and then it is reassigned to `n_diff` whereas the T_{LHS} directly assigns the value to the template variable, `:[v0]`. These data flow variants can be matched if we transitively match the code through the Data edges. Therefore, the nodes `int(n_diff)`, `Assignment`, `number(0)` in Figure 3(b) are matched transitively and are disconnected from the main graph. Finally, we obtain the matched graph (G^M), a disconnected sub-graph which primarily contains information about the locations where the nodes in *CPAT* can be applied.

4) *Inferring transformation Rules* ($TR^A : T_{LHS}^A \rightarrow T_{RHS}^A$): Now we know the code elements in the *target code* that correspond to nodes in T_{LHS} , the next step is to infer the adapted TR^A (Definition III.6) that can be applied to automate the program transformation. The higher level intuition is to: (i) capture the structure of the *target code* that is matched to $fgPDG^T$, and (ii) infer an adapted transformation rule ($T_{LHS}^A \rightarrow T_{RHS}^A$) by making matched nodes in the target code as holes, if necessary, and creating or deleting the new holes in the T_{LHS} or T_{RHS} to transform the matched nodes and the context. We now explain the steps we follow to adapt the rule in row-4-Table I (inferred from Listing 1) to the rule in Listing 3, which can be applied on Listing 2, an unseen variant.

In order to generate the adapted transformation rule (TR^A), we must first rename the matched nodes with the matched *Template Variable* in $fgPDG^R$. `RENAME_TEMPLATE_VARS` in Algorithm 2 renames the matched nodes in the *target code* using G^M which contains which template variables are matched to which *target code* nodes. As shown in Figure 4(a), we rename, `:[v0]` \rightarrow `n_diff`, and `:[v1]` \rightarrow `diff`. *AlphaHole* (`:[v2]`) is used for `to_eval.getDiff()` because it matches several nodes, resulting in the renaming of an entire expression.

The next step of generating TR^A is to rename the target code's context, i.e., unmatched code nodes, using appropriate *template variables*. We use the template variable V^A , (`:[v]`) to rename statements (e.g., *for loop*), whereas the V^I (`:[v]`), is used to rename identifier names. The nodes that have

been renamed thus far in the operation in line 2 are taken into consideration when deciding which *template variable* to employ. We use a new \mathbb{V}^A , if none of the nodes in a statement have been renamed. For example, $:[v3]$ in Figure 4(b) is used to rename `to_eval = getarray()`. We use \mathbb{V}^I to rename identifiers that are still intact after part of a statement has already been replaced. Here, we use the same \mathbb{V}^I for the same identifier. For example, in Figure 4(b), we use the $:[v4]$ to rename the identifier `total`.

The next step in Algorithm 2 is to use the operation `SUBSTITUTENODES` on T_{LHS}^A to generate T_{RHS}^A which performs two actions: (i) replaces the nodes $T_{LHS}^A \cap T_{LHS}$ with T_{RHS} (line 8). The replacing statements will be ordered as specified in T_{RHS} , so that the final transformation resembles the input code change example, (ii) builds the context of the code using the new template variables ($T_{LHS}^A - T_{LHS}$). Here, we discard the transitively matched nodes (i.e., $:[v4]$) identified in the sub-graph mining algorithm described in Section III-B3. This ensures that the template expression, T_{LHS}^A , matches the data variants in the code and then transforms it to the code defined in T_{RHS} . The example in Figure 4(c) shows the two actions performed by `SUBSTITUTENODES`. It first substitutes the nodes $T_{LHS}^A \cap T_{LHS}$ in T_{LHS}^A , i.e., the `for` loop and its body from $:[v0] = np.sum(:[v2])$. Then, it uses the nodes $T_{LHS}^A - T_{LHS}$, i.e., $:[v3]$ to create the context. Finally, the algorithm outputs the adapted rule, $T_{LHS}^A \rightarrow T_{RHS}^A$.

5) **Eliminating Unsafe Transformations:** To safely transform code, refactoring researchers [30], [31], [61] use preconditions that must be satisfied before transforming a target site. In a similar spirit, the operation `ISSAFE` in Algorithm 2 evaluates *lightweight* preconditions on the identified transformable code locations (\mathbb{G}^M) before applying the rule to a target site. PYEVOLVE uses the following preconditions:

Precondition 1: Some CPATs may involve the deletion of variables used in T_{LHS} . For example, in Listing 2, the loop variable `dif` and the local variable `total` will be deleted and will no longer be available once `np.sum` is replacing the original `for` loop. However, in Python, loop and local variables can be used outside the scope of the loop or block where they are initialized, e.g., they can be used further down in the code after the `for` loop. If that was the case, transforming the loop in Listing 2 to `np.sum` and deleting those two variables would be unsafe. PYEVOLVE checks to see if any variables that are marked for deletion are later used in the code. If this is the case, PYEVOLVE does not proceed with the transformation.

Precondition 2: The *fgPDG*-based matching algorithm, as described in Section III-B3, can identify target sites in a program that contain extra statements within scopes defined in T_{LHS} . This would make the transformation unsafe. For example, suppose that Listing 2 contained an additional statement such as `print(dif)` inside the `for` loop body. The matching algorithm would identify this as a potential site to transform. However, this target site should not be replaced by `numpy.sum` because it does not preserve semantics. To prevent transforming such target sites and ensure safety, PYEVOLVE discards target sites that contain outgoing edges from action

nodes (i.e., \mathbb{N}^A) or control nodes (i.e., \mathbb{N}^C) in the matched graph (i.e., $\in \mathbb{G}^M$) to the nodes $\in \text{fgPDG}^T$ (but $\notin \mathbb{G}^M$).

Precondition 3: The T_{RHS} of CPATs might contain APIs from third-party libraries (e.g., *TensorFlow*). Transplanting such CPATs into projects that do not use those libraries would break the code. To ensure the transformation’s safety, PYEVOLVE checks whether it can fully resolve the API invocations in the T_{RHS} . For example, if the `tf.sum` is in the T_{RHS} , PYEVOLVE checks the project’s *requirements.txt* [86] to see if the *TensorFlow* is listed as one of the project’s libraries. If it is not included, it will discard the transformation.

Defining a complete list of preconditions is hard. We believe this very challenging problem could be addressed by capturing more context and analyzing dynamic traces. Therefore, our implementation is flexible to express additional preconditions as they emerge. We further elaborate on safety in Section V.

IV. EVALUATION

We empirically evaluate PYEVOLVE and we answer the following research questions:

RQ1. What is the effectiveness of PYEVOLVE in generating correct code transformations? We conduct cross-validation to determine that PYEVOLVE correctly transforms code by replicating real-world CPATs. We report PYEVOLVE’s overall effectiveness by its precision and recall.

RQ2. What is the contribution of data- and control- flow aware rule inference to overall effectiveness? To perform this analysis, we report the number of changes that would be impossible to perform without the features in PYEVOLVE.

RQ3. How do developers find PYEVOLVE’s changes useful? To answer this, we submit pull requests to open-source projects containing patches generated by PYEVOLVE and record the developers’ responses.

A. RQ1: What is the effectiveness of PYEVOLVE in generating correct code transformations?

To answer this question, we replicate with PYEVOLVE thousands of transformations that open-source developers applied manually on their projects. We conducted cross-validation to test whether PYEVOLVE correctly transforms real-world CPATs. We divided the human adaptations published by Dilhara et al. [21] into training and testing sets. The training set is used to learn the initial transformation rules, while the test set is used to apply the CPATs learned from the training set. We assessed PYEVOLVE’s effectiveness by comparing the syntactic and semantic equivalence of Python-transformed codes to those that developers manually performed. We report the overall effectiveness by computing precision and recall.

1) **Dataset:** Dilhara et al. [21] studied 2,500 CPATs that occurred in 1000 top-rated ML repositories. The authors released CPAT dataset which is shown to be diverse with respect to size, frequency, authors, and projects. With a survey of 650 developers, the authors further confirm developers’ desire to have the identified CPATs automated in their code. Hence, we use the same dataset to perform cross validation. Section III-B1 provides further details on the CPAT dataset.

2) **Experimental setup:** Each *CPAT*, denoted as cp_k , consists of three or more instances of code changes, denoted as $\{k_1^{cp_k}, k_2^{cp_k}, \dots, k_m^{cp_k}\}$. We split the change examples in each *CPAT* into training and test sets for every cross validation iteration. One iteration of our cross-validation process selects one instance (e.g., $k_1^{cp_k}$) from which to learn the initial transformation rule, and PYEVOLVE then generates adapted rules to apply the *CPAT* to the test data (e.g., $\{k_2^{cp_k}, k_3^{cp_k}, \dots, k_m^{cp_k}\}$). We identify one *transformation trial* as inferring a rule from $k_i^{cp_k}$ and applying it to $k_j^{cp_k}$ where $(i \neq j)$. Applying PYEVOLVE to our dataset yielded over 40,000 trials, providing a high degree of confidence in the effectiveness of results.

To evaluate the effectiveness of transformations performed by PYEVOLVE, we use the *CPATs* performed by developers as the ground-truth (the oracle). We use PYEVOLVE to replicate *CPATs* in Oracle and compare them to the changes made by the original developer. We compute precision inside one iteration as the percentage of PYEVOLVE-applied transformations (i.e., trials) that are correct (i.e., equivalent to the ground-truth). We compute recall within one iteration as the percentage of all transformations from the ground truth that PYEVOLVE was able to transform. We obtain mean values over all iterations once we have precision and recall for each iteration.

To determine if a PYEVOLVE-applied transformation is correct, we have both a manual and an automated validation. Since we validate 40,000 transformation trials, manually checking them is tedious. Therefore, we chose a statistically significant (95% confidence level) random sample of transformation trials and manually validated the semantic correctness by comparing them to the human-performed transformations.

For the automated validation, we took inspiration from the steps that Noda et al. [47] used in their distinguished study to evaluate the correctness of an automated bug patch. We denote the code AST nodes that represent the T_{LHS}^A and T_{RHS}^A as A_{LHS} and A_{RHS} , respectively. We judge that the transformed code is syntactically correct, if the transformed code: (i) contains all $N \in A_{RHS} - A_{LHS}$ (new nodes), (ii) does not contain all $N \in A_{LHS} - A_{RHS}$ (deleted nodes), and (iii) contains all $N \in A_{LHS} \cap A_{RHS}$ (unchanged nodes). Conditions (1) and (2) ensure that pattern code is successfully inserted and deleted. Condition (3) assures that no excessive changes are made.

3) **Results:** In our replication of actual *CPATs* from open-source projects, PYEVOLVE performed 40,000 transformation trials in total, and we discovered that it achieved 97% precision and 94% recall. We manually validated a statistically significant sample of 381 instances for semantic validity, and it achieved 95% precision and 91% recall, which is slightly less than the automated validation. This is because the automated validation checks whether the *CPAT* nodes are fully transplanted but does not check their semantic validity. We primarily identified three causes for why PYEVOLVE either failed to perform transformations at all or did so incorrectly.

1) **Python Either types:** We employ type inference to obtain the type information of a Python program element at compile time, which is critical in mining semantically equivalent codes. Python allows for the holding of multiple types in the same

variable through the use of *Either* types. For example, the accumulator variable `n_diff` in Listing 2 can be assigned to a `String` before the `for` loop, making `n_diff` type of `Union[int, str]`. Algorithm 1 refers to node types for Guards with types to determine whether two *fgPDG* nodes are equal. Due to the fact that it searches for `n_diff` of type `int` rather than `Union[int, str]`, PYEVOLVE will not transform such cases.

2) **Abstracting over one example:** We use *PyInfer* to infer initial transformation rules, which abstract over one example, and might generate over-specialized transformation rules. For example, `m.at1.getM()` can be generalized either to `:[v0]].getM()` or `:[v0]]:[v1]].getM()`. To decide this, other existing techniques use multiple examples, which may also result in over generalization. Our matching technique is dependent on the *fgPDG* nodes, and thus on the rule’s generalization when matching semantic variants in the target codes. Hence, PYEVOLVE missed cases where it did not match the generalization of the rule. However, developers evolve Python in a Pythonic way [56], [57], and they evolve idioms mostly in the same way. Hence, despite learning from one example, PYEVOLVE achieved a high precision and recall.

3) **Semantically nonidentical instances in CPAT:** *PyType* infers the type `Any` for program elements for which there is insufficient information to determine the correct type. This may result in unrelated examples being grouped into the same *CPAT*. PYEVOLVE under-performs when we learn rules from such cases or apply other rules to them.

PYEVOLVE achieves an overall precision of 97% and recall of 94% in the cross-validation evaluation, confirming its effectiveness in inferring rules: $T_{LHS}^A \rightarrow T_{RHS}^A$.

B. RQ2: What is the contribution of data- and control-flow aware rule inference to overall effectiveness?

To answer this question, we use PYEVOLVE to apply *CPATs* on new target sites, and then we perform a sensitivity analysis to determine the impact of our novel contribution. Researchers have developed many techniques for inferring transformation rules from human adaptations. The number of code examples for which we can apply these rules to transform codes is nonetheless constrained by the rules’ propensity to overfit the examples. The complexity of *CPATs* in Python systems ranges from simple API updates to changes involving complex control structures [21]. To automate these code modifications, the rules generated by the tools should be capable of managing data-flow or control-flow variations of the examples. To address this, we use *data-flow control-flow* aware rule inferencing in PYEVOLVE. In this RQ, we look at how many data- and control-flow variants PYEVOLVE can automate that would be impossible to automate using existing techniques.

1) **Dataset:** We chose highly rated, best-in-class 20 Python *ML* projects like *TensorFlow*, *Pytorch*, *Keras*, and *Scikit-Learn*, etc. Finding opportunities for applying *CPATs* in these high-quality, professionally-maintained projects shows that PYEVOLVE can detect subtle variants that even the most expert programmers may have missed in their programming practices.

Table II: CPATs and their number of variants

Transformation Rule		Guard	N	V	I (V/N)
<pre>:[v0] = 0 for :[v1] in :[v2]: :[v0] = :[v0] + :[v1]</pre>	\rightarrow <pre>:[v0]=np.sum()</pre>	<pre>type := :[v0] -> int type := :[v1] -> int type := :[v2] -> List[int]</pre>	32	19	59%
<pre>:[v0] = '' for :[v2] in :[v3]: :[v0] += :[v2]</pre>	\rightarrow <pre>:[v0]=:[v1].join(:[v3])</pre>	<pre>type := :[v0] -> str type := :[v3] -> List[str]</pre>	8	3	38%
<pre>:[v3] = [] :[v0] = 0 for :[v1] in :[v2]: :[v0] = :[v1] + :[v0] :[v3].append(:[v0])</pre>	\rightarrow <pre>:[v0]=np.cumsum(:[v2])</pre>	<pre>type := :[v2] -> List[int] type := :[v1] -> int type := :[v0] -> int type := :[v3] -> List[int] import := np -> numpy</pre>	6	5	83%
<pre>:[v0]=0 for :[v1], :[v2] in zip(:[v3], :[v4])\rightarrow :[v0]=np.dot(:[v3],:[v4]) :[v0] += :[v1] * :[v2]</pre>		<pre>type := :[v0] -> int type := :[v1] -> int type := :[v2] -> int</pre>	3	3	100%
<pre>sum(:[v1])/len(:[v1])</pre>	\rightarrow <pre>np.mean(:[v1])</pre>	<pre>type :[v1] -> List[int] type :[v2] -> List[int]</pre>	17	0	0%
<pre>:[v0]=np.dot(:[v1], :[v2]) np.dot(:[v0],:[v3])</pre>	\rightarrow <pre>np.linalg.multi_dot(:[v1], :[v2],:[v3])</pre>	<pre>import := np : numpy type := :[vo] -> Matrix type := :[v1] -> Matrix type := :[v2] -> Matrix type := :[v3] -> Matrix</pre>	115	96	90%
<pre>:[v0]=:[v1].apply(tf .batch_and_drop_remainder(:[v2]))</pre>	\rightarrow <pre>:[v0]=:[v1].batch(:[v2], drop_remainder=True)</pre>	<pre>type := :[vo] -> Dataset type := :[v1] -> Dataset type := :[v2] -> int import := tf : tensorflow</pre>	2	2	100%
<pre>if :[v2] in :[v3]: :[v1] = :[v3][:[v2]].strip() else: :[v1] = :[v4]</pre>	\rightarrow <pre>:[v1] = :[v3].get(:[v2], :[v4]).strip()</pre>	<pre>type := :[v1] -> str type := :[v2] -> str type := :[v3] -> List[str] type := :[v4] -> str</pre>	2	0	0%
Total			185	128	70%

N: Transformed CPAT instances

V: Number of data and control variants of the original rule

I: number of variants as a percentage of N

2) **Experimental setup:** We randomly chose eight CPATs kinds from the list provided by Dilhara et al. [21] and used PYEVOLVE to automatically apply them to new target sites in the top-rated projects. Then, to determine if the modified code is a *data* or *control-flow* variant, we manually compared it to the original example that PYEVOLVE used as input.

3) **Results:** Table II shows the evaluation results. Columns: (i) *Transformation Rule*-shows the initial transformation rule that the *InferRule* [33] generated, (ii) *Guard*-shows the *Preds* relevant to rules, (iii) *N*-shows the total number of CPATs instances transformed by PYEVOLVE, and (iv) *V*-shows the number of *data* or *control-flow* variants transformed by PYEVOLVE and (v) *I*- shows the number of variants that PYEVOLVE made possible as a percentage of all transformations.

As seen in Table II, running PYEVOLVE on 20 projects transformed a total of 185 instances. This shows that even in the top-rated projects, there are many instances to apply the best practices, and developers often overlook these instances.

Listing 4: The project Dialogue from Baidu uses a `for` loop to compute cumulative sum of `seq_lens`

```

1 cur_len = 0
2 lod = [cur_len]
3 seq_lens = [len(ids) for ids in data_ids]
4 for l in seq_lens:
5     cur_len += 1
6     lod.append(cur_len)
```

Table II shows the data-or control flow variants for each rule. Of the 185 instances, 128 (70%) are *data-flow* or *control-flow* variations of the examples used to build the initial transformation rules. These instances would not be possible to transform without PYEVOLVE. For example, Listing 4 shows

a target code that the PYEVOLVE was able to automate, which uses a `for` loop to compute the cumulative sum when `np.cumsum` should have been used instead. This is an unseen variant of the LHS of the rule given in row-3-Table II because: (i) `cur_len` and `lod` are initialized in a different order, (ii) there is an additional assignment statement, and (iii) it accumulates `cur_len` differently than the rule does. The target code was previously unseen during the rule inference, yet PYEVOLVE was able to successfully automate it. This demonstrates the significant contribution of data-and control-flow aware rule inference, and how PYEVOLVE improves over the previous state-of-the-art for automating CPATs in Python ML systems.

70% of PYEVOLVE's code transformations are data or control-flow variants that cannot be transformed using existing techniques, thus improving the state-of-the-art.

C. RQ3:How do developers find PyEvo's patches useful?

To determine how useful PYEVOLVE is for real-world developers, we automatically applied frequent CPATs from our corpus to well known open source projects using PYEVOLVE, and then submitted these patches as pull requests.

1) **Dataset:** We chose 35 best-in-class projects like *TensorFlow*, *Keras*, *PyTorch*, and *Scikit_Learn*, and apply the CPATs. Dilhara et al. [10] performed an empirical study on diverse corpus of CPATs and revealed the dominant CPAT kinds in Python systems. We chose the same set of CPATs which covers all the kinds of CPATs revealed by the authors. Section III-B1 provides an explanation of the identified CPATs kinds.

2) **Experimental setup:** PYEVOLVE transformed 181 instances of CPATs that improve the performance and quality of

the affected Python code. These patches updated 116 source code files and affected 1028 SLOC. After PYEVOLVE applied the CPATs in each project, we ran all the test cases of projects to ensure that the changes did not break the code. We then notified the open-source project maintainers via pull requests.

3) **Results:** Even the highly rated and optimized codes, like *Keras*, *PyTorch*, and *TensorFlow*, accepted our pull requests. Table II shows the number of CPAT instances submitted for each CPAT kind. We submitted 40 pull requests with 181 CPAT instances. At the time of writing the paper, 25 (65%) PRs containing 163 CPAT instances were accepted, 4 pull requests were rejected, and the rest are still under review.

Developers found that our changes either enhance performance or code quality, or both. A lead developer from the project *Prosodic*, an NLP meta-library, mentioned that “Well done, your changes are cleaner and either faster or equivalently faster.” Another developer from the ML library *Transferlearning* applauded the PR: “Your changes improve the efficiency. I did not pay attention to this efficiency before.” Developers confirmed that they were aware of the best practices, but a tool like PYEVOLVE is able to identify opportunities that even the experts missed. For example, a developer from the ML library *Ann-benchmarks* remarked that “The changes look good, I am not sure why we didn’t write it that way before.” We submitted several CPAT instances where the developers should have used efficient ML libraries that were already being imported as dependencies to the project, but instead they were employing inefficient Python constructs. For example, in 83 (37%) instances, projects had libraries as project dependencies, and in 62 (28%) cases, they already had the library imported in the changed file, yet they still missed the opportunity to use the ML library at its fullest potential.

We discovered four major reasons for pull request rejections.

- 1) **Some CPATs are dependent on matrix shape.** The performance of matrix operations is affected by their shape. For example, we submitted a pull request to project *Mne-python* to replace multiple calls to `np.dot` with `np.linalg.multi_dot`. The developers rejected this because while `multi_dot` improves performance on non-square matrices, it degrades performance on square matrices. To improve, we must update the Guards for this CPAT to account for the matrix shape. However, to the best of our knowledge, there are no tools that can infer matrix shapes in Python at compile time.
- 2) **Dependencies on Hardware.** ML library optimizations depends on hardware platforms. For example, *Pytorch* is optimized for GPU use, whereas *NumPy* is not. A pull request submitted to *Pytorch* to replace a `for` loop with *NumPy* APIs was rejected because *NumPy* is not optimized for GPUs.
- 3) **Functions are already optimized with NumBa.** *NumBa* [39] is a library that translates annotated Python codes to optimized machine code at runtime. Because the optimization occurs at runtime, the code is faster than it looks at compile time, even if the code employs inefficient constructs like Python `loops`. Our patch submitted to *Pynndescent* was rejected because their code was already optimized at runtime.

- 4) **Deprecated code no longer updated.** Developer from the project *Basenji* rejected our patch because it changed a deprecated function that they will remove in a later release.

PYEVOLVE transformed 181 CPAT instances, of which 163 have been approved as of writing.

D. Threats To Validity

- 1) **Internal Validity:** *Does our tool produce valid results?*

We thoroughly evaluated the accuracy of the transformations produced by PYEVOLVE. To understand if the inferred rules can be trusted, the authors both automatically and manually validate the transformations to identify non-conforming ones. Furthermore, we develop a comprehensive setup that semi-automatically validates the application of transformation rules for a large and diverse set of CPATs.

- 2) **External Validity:** *Do our results generalize?* Although PYEVOLVE is effective on the evaluated data-set, it may not perform well on other subjects. To address this issue, we chose a large dataset that previous researchers [21] used. It covers different scenarios and has been shown to be diverse in terms of frequency, size, authors, and projects. All our subjects are open-source, and we have yet to evaluate proprietary codes.

We performed manual steps in RQ2 (identifying the variants) and RQ3 (submitting patches). These manual steps prevented us from using all the CPATs for the evaluation of RQ2 and RQ3, so we only used a subset of CPATs. This could impact the generalizability to other CPAT kinds. To mitigate this, we used a randomly selected subset of CPATs.

- 3) **Reliability:** The collected data, source code, and executable of PYEVOLVE are publicly available [44] and ready for reuse.

V. DISCUSSION

- 1) **Safety and soundness of our approach:** As illustrated in the very recent keynote by Gulwani at FSE-2022 [87], the “Transformation by Example” systems are not intended to replace developers, nor are they designed to be sound. Our PYEVOLVE (like other “Transformation by Example” systems) sits in the middle between a regex-based find-replace tool and a refactoring tool. It has the expressivity of a find-and-replace tool, while being syntax-aware. In contrast, a refactoring tool or a compiler optimization tool have hard-coded, task-specific rules that make them safer to use, but they are expensive to develop. Determining the safety of such transformations require deep analysis of the context, hence we do not recommend blindly accepting patches from PYEVOLVE. However, in our workflow, we trade-off safety for broader applicability by relying on the developer’s insight in determining whether it is safe to perform the transformation. In their study, Ketkar et al. [33] observed that rules must sometimes be manually vetted to ensure their safety and soundness. If this is the case, our human-comprehensible rules make the process easier. In our empirical evaluation, PYEVOLVE achieved 97% precision and it did not require any intervention from us to achieve this.
- 2) **Selecting appropriate transformation rules:** We position PYEVOLVE as a companion development tool for developers.

i.e., PYEVOLVE modifies code under the developers’ supervision. It requires that developers select appropriate rules to apply to their code based on their desires. *R-CPATMiner* generates a large number of *CPATs*, some of them could apply for the same code location. For example one *CPAT* would transform a loop into a `numpy` optimized operation, while another *CPAT* would transform the same loop into a `list` comprehension. In such cases PYEVOLVE cannot determine which *CPATs* should be transplanted into the code base. Hence, PYEVOLVE generates one rule for each *CPAT* and asks the developer to make the final selection. This behavior is inherited from the “Transformation by Example” paradigm.

VI. RELATED WORK

We group the related work in two areas: (i) inferring and applying changes, and (ii) Python code automation.

Inferring rules and Applying change patterns: Researchers have developed an array of advanced program transformation systems that automatically generate programs according to given input-output examples, so called “Transformation by Example”. This technique has been used in a variety of applications, including (i) Java Type Migration [33], (ii) API migration [1], [5], [15], [19], [54], (iii) String manipulation [6], [7], (iv) Data Structure Transformation [3]. LASE [53], REFAZER [35], SPDIFF [4], TCINFER [33], and APIFIX [5] learn from multiple human adaptations, determining how to abstract the adaptations based on their commonalities and differences, whereas SYDIT [49], A4 [54] and MEDITOR [19] abstract over individual examples. APPEVOLVE [15] abstracts over a multiple examples, and then find the closest match to the target code, based on the ability of the code to pass the provided test cases. These works aim to generate a properly generalized rule by varying the number of examples and its properties. However, they do not account for previously unseen data- and control-flow variations of the examples, resulting in low recall.

Tools such as TCINFER [33], SPINFER [40] and the tool by Newman et al. [83] take into account control-flow variants to some extent. They split the rule inference task into multiple sub-tasks and solve the sub-tasks separately. These tools target transformations of short statements like *API* invocations and assignment statements. They can accommodate for an arbitrary number of unrelated/noise statements between the short statements that make up the *CPAT*. However, the *CPATs* discovered in Python systems contain not only *API* invocations but also complex control structures such as `for` loops, and cannot be handled by splitting the rules. Therefore, we suggest a novel graph-based technique that captures both *data*- and *control-flow* unseen variants, then adapts the generated rules for the new contexts and correctly handles complex *CPATs*.

Studies involving Python idioms: Researchers studied Python idioms and how they were used in Python systems. Phan-udom et al. [55] recommend 58 non-idiomatic and 55 idiomatic changes. Alexandru et al. [56] present a non-exhaustive list of Python idioms gleaned from a developer survey. Sakulniwat et al. [57] studied the evolution of Python *with* statements over time, while Wang et al. [59] studied

Python code smells. However, none of these fully automate code transformations. Using PYEVOLVE, we can infer the rules and fully automate all these idioms. Recently, Zhang et al. [23] employed AST rewriting to automatically refactor nine Python idioms. However, they must hard-code the transformations. In contrast, we present PYEVOLVE, which mines and automates idioms using “Transformation by Example”. Thus, as new idioms emerge, PYEVOLVE is future-proof and will handle new idioms. This will make it much easier for Python-ML developers—the dominant ecosystem in Python [11]—to keep up with the rapidly advancing ML techniques.

VII. CONCLUSIONS AND FUTURE WORK

Despite Python’s and ML’s meteoric rise, support for automated code evolution is still in its infancy. To advance the science and tooling for automating code changes in Python, we built PYEVOLVE, which infers transformation rules and then applies them. Unlike existing refactoring tools that are hard-coded for specific transformations, PYEVOLVE automates a wide range of best practices from code sources. Our key innovation over existing “Transformation by Example” systems is a *data*- and *control-flow* aware transformation rule inference.

Our thorough empirical evaluation of a diverse, representative corpus of 40,000 transformation trials from real-world projects shows that PYEVOLVE is effective. It has a 97% precision and 94% recall, and 70% of PYEVOLVE transformations would be impossible to automate without our technique. Developers accepted 90% of the 181 *CPATs* that PYEVOLVE produced, and their feedback shows PYEVOLVE’s usefulness.

We anticipate these future advancements for PYEVOLVE:

- 1) **Version Awareness:** *CPATs* use language and library constructs that continuously evolve based on their release versions. We will extend the transformation rules to be version aware.
- 2) **Community repository of transformation rules:** The rules that PYEVOLVE learns may change over time; some rules may become obsolete while new rules are emerging. For this reason, the researchers have called for a community-maintained central database of *CPATs* [34] and their respective rules [2], [33] that would need to be properly versioned, maintained, and evolved. The likelihood of having data-flow or control-flow variations, however, rises as the *CPAT* gets bigger. If we add more rules to the database in order to represent all variants, it would significantly increase the size of the database and, in turn, expand the search space with significant slow downs. Our novel approach uses a single rule that transforms all other data- and control-flow variants, thus it significantly reduces the number of rules while increasing the speed.
- 3) **Expanding to other domains and languages:** As new idioms emerge, PYEVOLVE is future-proof and will continue to handle new idioms. We think PYEVOLVE will aid ML developers in keeping up with the rapidly advancing ML technologies. While the corpus of programs we use in this paper is ML systems, PYEVOLVE is readily applicable to identify, recommend, and automate *CPATs* for any kind of Python software system. Moreover, the core ideas presented here are easily transferable to other programming languages.

Given how well they work even for a dynamically-typed language like Python, we think they would work even better for statically-typed languages. By making our design, implementation, and datasets available to the community [44], we hope this would inspire many others to advance the field.

A. Data Availability

We released the executable of PYEVOLVE and its evaluation data-set on an anonymous publicly accessible website [44].

REFERENCES

- [1] M. Lamothe, W. Shang and T. -H. P. Chen, "A3: Assisting Android API Migrations Using Code Examples," in *IEEE Transactions on Software Engineering*, 1 Feb. 2022, doi: 10.1109/TSE.2020.2988396.
- [2] Mattia Fazzini, Qi Xin, and Alessandro Orso. 2020. APIMigrator: an API-usage migration tool for Android apps. In *Proceedings of the IEEE/ACM 7th International Conference on Mobile Software Engineering and Systems* <https://doi.org/10.1145/3387905.3388608>
- [3] Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. 2017. Component-based synthesis of table consolidation and transformation tasks from examples. *SIGPLAN Not.* 52, 6 (June 2017), 422–436. <https://doi.org/10.1145/3140587.3062351>
- [4] Andersen, J., Lawall, J.L. Generic patch inference. *Autom Softw Eng* 17, 119–148 (2010). <https://doi.org/10.1007/s10515-010-0062-z>
- [5] Xiang Gao, Arjun Radhakrishna, Gustavo Soares, Ridwan Shariffdeen, Sumit Gulwani, and Abhik Roychoudhury. 2021. APIfix: output-oriented program synthesis for combating breaking changes in libraries. *Proc. ACM Program. Lang.* 5, OOPSLA, <https://doi.org/10.1145/3485538>
- [6] Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. *SIGPLAN Not.* 46, 1 (January 2011), 317–330. <https://doi.org/10.1145/1925844.1926423>
- [7] Rishabh Singh. 2016. BlinkFill: semi-supervised programming by example for syntactic string transformations. *Proc. VLDB Endow.* 9, 10 (June 2016), 816–827. <https://doi.org/10.14778/2977797.2977807>
- [8] Ru Zhang, Wencong Xiao, Hongyu Zhang, Yu Liu, Haoxiang Lin, and Mao Yang. 2020. An empirical study on program failures of deep learning jobs. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering (ICSE '20)*. <https://doi.org/10.1145/3377811.3380362>
- [9] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hridesh Rajan. 2019. A comprehensive study on deep learning bug characteristics. *ESEC/FSE 2019*. <https://doi.org/10.1145/3338906.3338955>
- [10] Malinda Dilhara, Ameya Ketkar, and Danny Dig. 2021. Understanding Software-2.0: A Study of Machine Learning Library Usage and Evolution. *ACM Trans. Softw. Eng. Methodol.* 30, 4, Article 55 (October 2021), 42 pages. <https://doi.org/10.1145/3453478>
- [11] H. Ben Braiek, F. Khomh and B. Adams, "The Open-Closed Principle of Modern Machine Learning Frameworks," 2018 *IEEE/ACM 15th International Conference on Mining Software Repositories (MSR)*, 2018
- [12] Raschka S, Mirjalili V. Python machine learning: Machine learning and deep learning with Python, scikit-learn, and TensorFlow 2. Packt Publishing Ltd; 2019 Dec 12.
- [13] Armando Solar-Lezama, Christopher Grant Jones, and Rastislav Bodik. 2008. Sketching concurrent data structures. *SIGPLAN Not.* 43, 6 (June 2008), 136–148. <https://doi.org/10.1145/1379022.1375599>
- [14] Armando Solar-Lezama, Rodric Rabbah, Rastislav Bodik, and Kemal Ebcioglu. 2005. Programming by sketching for bit-streaming programs. *SIGPLAN Not.* 40, 6 (June 2005), 281–294. <https://doi.org/10.1145/1064978.1065045>
- [15] Mattia Fazzini, Qi Xin, and Alessandro Orso. 2019. Automated API-usage update for Android apps. *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*. <https://doi.org/10.1145/3293882.3330571>
- [16] A. Hora, N. Anquetil, S. Ducasse and M. T. Valente, "Mining system specific rules from change patterns," 2013 20th Working Conference on Reverse Engineering (WCRE), doi: 10.1109/WCRE.2013.6671308.
- [17] Stefanus A. Haryono, Ferdian Thung, Hong Jin Kang, Lucas Serrano, Gilles Muller, Julia Lawall, David Lo, and Lingxiao Jiang. 2020. Automatic Android Deprecated-API Usage Update by Learning from Single Updated Example. *ICPC '20*. <https://doi.org/10.1145/3387904.3389285>
- [18] J. Henkel and A. Diwan, "CatchUp! Capturing and replaying refactorings to support API evolution," *Proceedings. 27th International Conference on Software Engineering*, 2005. *ICSE 2005*. doi: 10.1109/ICSE.2005.1553570.
- [19] S. Xu, Z. Dong and N. Meng, "Meditor: Inference and Application of API Migration Edits," 2019 *IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*, doi: 10.1109/ICPC.2019.00052.
- [20] Comby. 2022. Comby Syntax Reference. <https://comby.dev/docs/syntax-reference> Accessed: 11 July 2022.
- [21] M. Dilhara, A. Ketkar, N. Sannidhi and D. Dig, "Discovering Repetitive Code Changes in Python ML Systems," 2022 *IEEE/ACM 44th International Conference on Software Engineering (ICSE)*, 2022, doi: 10.1145/3510003.3510225.
- [22] Malinda Dilhara. 2021. Discovering repetitive code changes in ML systems. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2021)*. <https://doi.org/10.1145/3468264.3473493>
- [23] Zhang, Z., Xing, Z., Xia, X., Xu, X. and Zhu, L., "Making Python Code Idiomatic by Automatic Refactoring Non-Idiomatic Python Code with Pythonic Idioms". *Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2022)*. arXiv preprint [arXiv:2207.05613](https://arxiv.org/abs/2207.05613)
- [24] Abram Hindle, Earl T. Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. 2016. On the naturalness of software. *Commun. ACM* 59, 5 (May 2016), 122–131. <https://doi.org/10.1145/2902362>
- [25] Earl T. Barr, Yuriy Brun, Premkumar Devanbu, Mark Harman, and Federica Sarro. 2014. The plastic surgery hypothesis, *FSE 2014*. <https://doi.org/10.1145/2635868.2635898>
- [26] James O Coplien. 1991. *Advanced C++ programming styles and idioms*. Addison-Wesley Longman Publishing Co., Inc. 30
- [27] Slatkin B. *Effective python: 90 specific ways to write better python*. Addison-Wesley Professional; 2019 Oct 25.
- [28] Miltiadis Allamanis and Charles Sutton. 2014. Mining idioms from source code. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014)*. <https://doi.org/10.1145/2635868.2635901>
- [29] Stas Negara, Mihai Codoban, Danny Dig, and Ralph E. Johnson. 2014. Mining fine-grained code changes to detect unknown change patterns. In *Proceedings of the 36th International Conference on Software Engineering (ICSE 2014)*. <https://doi.org/10.1145/2568225.2568317>
- [30] Ittai Balaban, Frank Tip, and Robert Fuhrer. 2005. Refactoring support for class library migration. *SIGPLAN Not.* 40, 10 (October 2005), 265–279. <https://doi.org/10.1145/1103845.1094832>
- [31] D. Dig, J. Marrero and M. D. Ernst, "Refactoring sequential Java code for concurrency via concurrent libraries," 2009 *IEEE 31st International Conference on Software Engineering*, doi: 10.1109/ICSE.2009.5070539.
- [32] Rijnard van Tonder and Claire Le Goues. 2019. Lightweight multi-language syntax transformation with parser parser combinators, *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019)*. <https://doi.org/10.1145/3314221.3314589>
- [33] Ameya Ketkar, Oleg Smirnov, Nikolaos Tsantalis, Danny Dig, and Timofey Bryksin. 2022. Inferring and applying type changes. In *Proceedings of the 44th International Conference on Software Engineering (ICSE '22)*. <https://doi.org/10.1145/3510003.3510115>
- [34] H. A. Nguyen, T. N. Nguyen, D. Dig, S. Nguyen, H. Tran and M. Hilton, "Graph-Based Mining of In-the-Wild, Fine-Grained, Semantic Code Change Patterns," 2019 *IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, doi: 10.1109/ICSE.2019.00089.
- [35] Rolim, Reudismam, Gustavo Soares, Loris D'Antoni, Oleksandr Polozov, Sumit Gulwani, Rohit Gheyi, Ryo Suzuki, and Björn Hartmann, "Learning Syntactic Program Transformations from Examples," 2017 *IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, doi: 10.1109/ICSE.2017.44.
- [36] M. Nita and D. Notkin, "Using twinning to adapt programs to alternative APIs," 2010 *ACM/IEEE 32nd International Conference on Software Engineering*, doi: 10.1145/1806799.1806832.
- [37] Google.PyType. <https://github.com/google/pytype> Accessed: 2022-07-16.
- [38] J. Jiang, L. Ren, Y. Xiong and L. Zhang, "Inferring Program Transformations From Singular Examples via Big Code," 2019 34th *IEEE/ACM International Conference on Automated Software Engineering (ASE)*, doi: 10.1109/ASE.2019.00033.
- [39] NumBa. <https://numba.pydata.org> Accessed: 2022-07-16.

- [40] Serrano L, Nguyen VA, Thung F, Jiang L, Lo D, Lawall J, Muller G. SPINFER: Inferring Semantic Patches for the Linux Kernel. In 2020 USENIX Annual Technical Conference (USENIX ATC 20) 2020.
- [41] Yoann Padioleau, René Rydhof Hansen, Julia L. Lawall, and Gilles Muller. 2006. Semantic patches for documenting and automating collateral evolutions in Linux device drivers. (PLOS '06). <https://doi.org/10.1145/1215995.1216005>
- [42] Github. 2022. Octoverse. <https://octoverse.github.com> Accessed: 2022-07-16.
- [43] IntelliJ. 2021. IntelliJ: Structural Search and Replace. <https://www.jetbrains.com/help/idea/structural-search-and-replace.html> Accessed: 18 Jul 2022.
- [44] Artifacts, <https://pythoninfer.github.io>: 1 September 2022.
- [45] Danilo Silva, Nikolaos Tsantalis, and Marco Tulio Valente. 2016. Why we refactor? confessions of GitHub contributors. In Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2016). <https://doi.org/10.1145/2950290.2950305>
- [46] O. Smirnov, A. Lobanov, Y. Golubev, E. Tikhomirova and T. Bryksin, "Revizor: A Data-Driven Approach to Automate Frequent Code Changes Based on Graph Matching," 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE), doi: 10.1109/ASE51524.2021.9678635.
- [47] K. Noda, H. Yokoyama and S. Kikuchi, "Sirius: Static Program Repair with Dependence Graph-Based Systematic Edit Patterns," ICSME-2021, doi: 10.1109/ICSME52107.2021.00045.
- [48] Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: learning to fix bugs automatically. Proc. ACM Program. Lang. 3, OOPSLA, Article 159 (October 2019), 27 pages. <https://doi.org/10.1145/3360585>
- [49] Na Meng, Miryung Kim, and Kathryn S. McKinley. 2011. Systematic editing: generating program transformations from an example. SIGPLAN Not. 46, 6 (June 2011), 329–342. <https://doi.org/10.1145/1993316.1993537>
- [50] Anders Miltner, Sumit Gulwani, Vu Le, Alan Leung, Arjun Radhakrishna, Gustavo Soares, Ashish Tiwari, and Abhishek Udupa. 2019. On the fly synthesis of edit suggestions. Proc. ACM Program. Lang. 3, OOPSLA, Article 143 (October 2019), 29 pages. <https://doi.org/10.1145/3360569>
- [51] Reudismam Sousa, Gustavo Soares, Rohit Gheyi, Titus Barik, and Loris D'Antoni. 2021. Learning Quick Fixes from Code Repositories. In Brazilian Symposium on Software Engineering (SBES '21). <https://doi.org/10.1145/3474624.3474650>
- [52] Barthélemy Dagenais and Martin P. Robillard. 2011. Recommending Adaptive Changes for Framework Evolution. ACM Trans. Softw. Eng. Methodol. 20, 4, Article 19 (September 2011), 35 pages. <https://doi.org/10.1145/2000799.2000805>
- [53] N. Meng, M. Kim and K. S. McKinley, "Lase: Locating and applying systematic edits by learning from examples," ICSE-2013, doi: 10.1109/ICSE.2013.6606596.
- [54] Lamothe, Maxime, Weiyi Shang and Tse-Hsun Peter Chen. "A4: Automatically Assisting Android API Migrations Using Code Examples." ArXiv abs/1812.04894 (2018): n. pag.
- [55] Phan-Udom, P., Wattanakul, N., Sakulniwat, T., Ragkhitwetsagul, C., Sunetnanta, T., Choetkiertikul, M., Kula, R. G., "Teddy: Automatic Recommendation of Pythonic Idiom Usage For Pull-Based Software Projects," ICSME-2020, doi: 10.1109/ICSME46990.2020.00098.
- [56] Carol V. Alexandru, José J. Merchante, Sebastiano Panichella, Sebastian Proksch, Harald C. Gall, and Gregorio Robles. 2018. On the usage of pythonic idioms. In Proceedings of the 2018 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward! 2018). <https://doi.org/10.1145/3276954.3276960>
- [57] Sakulniwat, T., Kula, R. G., Ragkhitwetsagul, C., Choetkiertikul, M., Sunetnanta, T., Wang, D., Matsumoto, K., "Visualizing the Usage of Pythonic Idioms Over Time: A Case Study of the with open Idiom," IWSEF, doi: 10.1109/IWSEF49350.2019.00016.
- [58] Y. Tang, R. Khatchadourian, M. Bagherzadeh, R. Singh, A. Stewart and A. Raja, "An Empirical Study of Refactorings and Technical Debt in Machine Learning Systems," ICSE-2021, doi: 10.1109/ICSE43902.2021.00033.
- [59] T. Wang, Y. Golubev, O. Smirnov, J. Li, T. Bryksin and I. Ahmed, "PyNose: A Test Smell Detector For Python," ASE-2021, doi: 10.1109/ASE51524.2021.9678615.
- [60] Vélez, T. C., Khatchadourian, R., Bagherzadeh, M., Raja, A. "Challenges in Migrating Imperative Deep Learning Programs to Graph Execution: An Empirical Study." arXiv preprint arXiv:2201.09953 (2022).
- [61] Rick Ossendrijver, Stephan Schroovers, and Clemens Grelck. 2022. Towards automated library migrations with error prone and refaster. In Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing (SAC '22). <https://doi.org/10.1145/3477314.3507153>
- [62] Balland, E., Brauner, P., Kopetz, R., Moreau, P. E., Reilles, A. "Tom: Piggybacking rewriting on java." International Conference on Rewriting Techniques and Applications. Springer, Berlin, Heidelberg, 2007.
- [63] I. D. Baxter, C. Pidgeon and M. Mehlich, "DMS/spl reg/: program transformations for practical scalable software evolution," ICSE-2004, doi: 10.1109/ICSE.2004.1317484.
- [64] Marat Boshernitsan and Susan L. Graham. 2004. IXj: interactive source-to-source transformations for java. In Companion to the 19th annual ACM SIGPLAN conference on Object-oriented programming systems, languages, and applications (OOPSLA '04). 212–213. <https://doi.org/10.1145/1028664.1028755>
- [65] Bravenboer, M., Kalleberg, K. T., Vermaas, R., Visser, E. "Stratego/XT 0.17. A language and toolset for program transformation." Science of computer programming 72.1-2 (2008): 52-70.
- [66] Cordy, James R. "The TXL source transformation language." Science of Computer Programming 61.3 (2006): 190-210.
- [67] Mark Hills, Paul Klint, and Jurgen J. Vinju. 2012. Scripting a refactoring with Rascal and Eclipse. In Proceedings of the Fifth Workshop on Refactoring Tools (WRT '12). <https://doi.org/10.1145/2328876.2328882>
- [68] Li, Huiqing, and Simon Thompson. "A domain-specific language for scripting refactorings in erlang." International Conference on Fundamental Approaches to Software Engineering. Berlin, Heidelberg, 2012.
- [69] Mens, Tom, and Tom Tourwé. "A Declarative Evolution Framework for Object-Oriented Design Patterns." icsm. Vol. 1. 2001.
- [70] Steimann, Friedrich, Christian Kollee, and Jens von Pilgrim. "A refactoring constraint language and its application to Eiffel." European Conference on Object-oriented Programming. Berlin, Heidelberg, 2011.
- [71] van den Brand, M. G., van Deursen, A., Heering, J., de Jong, H. A., de Jonge, M., Kuipers, T., Visser, J. "The ASF+ SDF meta-environment: A component-based language development environment." Electronic Notes in Theoretical Computer Science 44.2 (2001): 3-8.
- [72] Mathieu Verbaere, Ran Ettinger, and Oege de Moor. 2006. JunGL: a scripting language for refactoring. ICSE '06. <https://doi.org/10.1145/1134285.1134311>
- [73] John Brant and Don Roberts. 2009. The SmaCC transformation engine: how to convert your entire code base into a different programming language. OOPSLA '09. <https://doi.org/10.1145/1639950.1640026>
- [74] L. Frenzel, "The Language Toolkit: An API for Automated Refactorings in Eclipse-based IDEs," <https://eclipse.org/articles/Article-LTK/ltk.html>.
- [75] Max Schaefer and Oege de Moor. 2010. Specifying and implementing refactorings. SIGPLAN Not. 45, 10 (October 2010), 286–301. <https://doi.org/10.1145/1932682.1869485>
- [76] Roberts, Donald Bradley. Practical analysis for refactoring. University of Illinois at Urbana-Champaign, 1999.
- [77] A.Garrido, "Program Refactoring in the Presence of Preprocessor Directives," Ph.D. dissertation, University of Illinois at Urbana-Champaign, 2005.
- [78] MLA Miller, Robert C., and Brad A. Myers. "Interactive Simultaneous Editing of Multiple Text Regions." USENIX Annual Technical Conference, General Track. 2001.
- [79] M. Toomim, A. Begel and S. L. Graham, "Managing Duplicated Code with Linked Editing," 2004 IEEE Symposium on Visual Languages - Human Centric Computing, doi: 10.1109/VLHCC.2004.35.
- [80] M. van den Brand, M. Bruntink, G. Economopoulos, H. de Jong, P. Klint, T. Kooiker, T. van der Storm, and J. Vinju, "Using The Meta-Environment for Maintenance and Renovation," in CSMR, 2007.
- [81] Pylint. 2022. Pylint tool. <https://pylint.org/> [Online. Accessed 01-Sept-2022]
- [82] Jevgenija Pantuchina, Bin Lin, Fiorella Zampetti, Massimiliano Di Penta, Michele Lanza, and Gabriele Bavota. 2021. Why Do Developers Reject Refactorings in Open-Source Projects? <https://doi.org/10.1145/3487062>
- [83] Newman, C. D., Bartman, B., Collard, M. L., Maletic, J. I. "Simplifying the construction of source code transformations via automatic syntactic restructurings." Journal of Software: Evolution and Process 29.4 (2017): e1831.

- [84] Haryono, S. A., Thung, F., Lo, D., Jiang, L., Lawall, J., Kang, H. J., Muller, G. AndroEvolve: automated Android API update with data flow analysis and variable denormalization. *Empir Software Eng* <https://doi.org/10.1007/s10664-021-10096-0>
- [85] Stefanus A. Haryono, Ferdian Thung, Hong Jin Kang, Lucas Serrano, Gilles Muller, Julia Lawall, David Lo, and Lingxiao Jiang. 2020. Automatic Android Deprecated-API Usage Update by Learning from Single Updated Example. *ICPC '20*, <https://doi.org/10.1145/3387904.3389285>
- [86] PIP. 2019. Requirements. Retrieved December 12, 2022 from <https://pip.pypa.io/en/stable/reference/requirements-file-format/>.
- [87] FSE-2022 KeyNote Summit Gulwani November 14, 2022 <https://youtu.be/9g9Ss-SvhME>.