

Discovering Repetitive Code Changes in Python ML Systems

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

Nikhith Sannidhi
nikhith.sannidhi@colorado.edu
University of Colorado Boulder
USA

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

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

ABSTRACT

Over the years, researchers capitalized on the repetitiveness of software changes to automate many software evolution tasks. Despite the extraordinary rise in popularity of Python-based ML systems, they do not benefit from these advances. Without knowing what are the repetitive changes that ML developers make, researchers, tool, and library designers miss opportunities for automation, and ML developers fail to learn and use best coding practices.

To fill the knowledge gap and advance the science and tooling in ML software evolution, we conducted the first and most fine-grained study on code change patterns in a diverse corpus of 1000 top-rated ML systems comprising 58 million SLOC. To conduct this study we reuse, adapt, and improve upon the state-of-the-art repetitive change mining techniques. Our novel tool, R-CPATMINER, mines over 4M commits and constructs 350K fine-grained change graphs and detects 28K change patterns. Using thematic analysis, we identified 22 pattern groups and we reveal 4 major trends of how ML developers change their code. We surveyed 650 ML developers to further shed light on these patterns and their applications, and we received a 15% response rate. We present actionable, empirically-justified implications for four audiences: (i) researchers, (ii) tool builders, (iii) ML library vendors, and (iv) developers and educators.

CCS CONCEPTS

- **Software and its engineering** → **Software maintenance tools**;
- **Computing methodologies** → **Machine learning**.

KEYWORDS

Refactoring, Repetition, Code changes, Machine learning, Python

ACM Reference Format:

Malinda Dilhara, Ameya Ketkar, Nikhith Sannidhi, and Danny Dig. 2022. Discovering Repetitive Code Changes in Python ML Systems. In *44th International Conference on Software Engineering (ICSE '22)*, May 21–29, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3510003.3510225>

*Ameya Ketkar contributed this work as a research assistant at Oregon State University.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
ICSE '22, May 21–29, 2022, Pittsburgh, PA, USA
© 2022 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-9221-1/22/05.
<https://doi.org/10.1145/3510003.3510225>

1 INTRODUCTION

Many software changes are repetitive by nature [7, 33, 53], thus forming change patterns. Like in traditional software systems, Machine Learning (ML) developers perform repetitive code changes too. For example, Listing 1 shows a common change where ML developers replaced a `for` loop that sums the list `elements` with `np.sum`, a highly optimized domain-specific abstraction provided by the library *NumPy* [61]. Since this change involves programming idioms [2, 80] at the sub-method level it is *fine-grained*. If this code change is repeated at multiple locations or in multiple commits, it is a *fine-grained code change pattern*.

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

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

Over the years, researchers in the traditional software systems have provided many applications that rely upon the repetitiveness of changes: code completion in the IDEs [13, 34, 43, 55, 56], automated program repair [6, 9, 50], API recommendation [32, 55], type migration [40], library migration [3, 18, 24, 39, 84], code refactoring [19, 28], fine-grained understanding of software evolution [3, 25, 42, 54, 58, 77, 85]. Unfortunately, these are mostly available only for Java, and do not support Python and ML systems.

Researchers [11, 22, 37, 70] observed that Python dominates the ML ecosystem in both the company-driven and the community-driven ML software systems, yet the tooling is significantly behind [22, 91]. In order to advance the science and tooling for ML code development in Python, we need to understand how developers evolve and maintain ML systems. Previous researchers have focused on high-level software evolution tasks like identifying ML bugs [35, 37, 38], updating ML libraries [22], refactoring and technical debt of ML systems [75, 82], managing version control systems for data [8], and testing [10, 31, 36]. However, there is a lack of understanding of the repetitive fine-grained code change patterns that ML developers laboriously perform. *What are fine-grained changes performed in ML systems? Which ones are ML-specific? What kinds of automation do ML developers need?*

Without answers to such questions, researchers miss opportunities to improve the state-of-the-art in automation for software evolution in ML systems, tool builders do not invest resources where automation is most needed, language and library designers

cannot make informed decisions when introducing new constructs, and ML developers fail to learn and use best practices.

In this paper, we conduct the first large-scale study and discover repetitive change patterns in Python-based ML systems. We employ both quantitative (mining repositories and thematic analysis) and qualitative methods (surveys) to answer the research questions. Blending these methods has the advantage of the results being triangulated. The quantitative methods help us discover *what* fine-grained change patterns ML developers perform. The qualitative method helps us answer *why* these changes are performed, how they are performed, and how tool builders can improve ML developer productivity.

For the quantitative analysis, we use a large data set of 1000 ML projects from GitHub, comprising 58 million source lines of code (SLOC) at the latest revisions, 1.16 million mapped code blocks, 1.5 million changed files, and 0.4 billion changed SLOCs. We extracted 28,308 fine grained code change patterns where 58% of them appear in multiple projects. We applied *thematic analysis* [12, 88] upon 2,500 most popular patterns from our dataset, and categorized them into 22 fine-grained change pattern themes that reveal 4 major trends. Moreover, we designed and conducted a survey with 650 ML developers, in which we presented 1,235 patterns for their feedback and achieved a 15% response rate. In the survey, 71% of the developers confirmed the need of automation for 22 pattern groups. Among these, we discovered four major trends: (1) *transform to Context managers* (e.g., disable or enable gradient calculation, swap ML training device), (2) *convert for loops to domain specific abstraction* (e.g., see Listing 1), (3) *update API usage* (e.g., migrate to `TensorFlow.log` from `log`, transform matrices), and (4) *use advanced language features* (e.g., transform to `list` comprehension).

The main challenge in conducting such large-scale, representative studies, is the lack of tools for mining non-Java repositories. To overcome this challenge we reuse, adapt, and extend the vast ecosystem of Java AST-level analysis tools [3, 25, 42, 54, 58, 77, 85] to Python. Most of these tools rely on techniques that are conceptually language-independent, i.e., they operate on intermediate representation of the code (e.g., AST nodes). Second, we observed that 72% of the Python AST node kinds identically overlap with those in Java (e.g., *While-Statement*, *Assignment-Statement*, etc.). Moreover, another 18% of Python AST node kinds also exist in Java with some differences (e.g., Python's `for` loop has multiple loop variables). Only 10% of the Python AST node kinds are unique to Python (e.g., *With statement*, *Generators*, etc.). Hence, one of our key ideas is to reuse the Java AST-level analysis tools to analyse 72% of the Python AST nodes and for the remaining 28% of AST nodes we either modify existing capabilities or add brand new ones.

We first developed a novel technique, `JAVAFYPY`, to transform Python AST to a format that can be processed by Java AST-level mining tools. We used `JAVAFYPY` to adapt to Python the state-of-the-art fine-grained change pattern mining tool, `CPATMINER` [58]. `CPATMINER` matches changed methods and their body statements across the commits and identifies fine-grained change patterns. Refactorings such as move, rename, and extract method, re-arrange and obfuscate the code statements, that are hard to match across the edit, leading `CPATMINER` to miss multiple occurrences of patterns. To improve the accuracy of `CPATMINER`, we integrate it with the state-of-the-art refactoring mining technique- `REFACTORINGMINER` [85],

that de-obfuscates the re-arranged code statements. Our novel tool `R-CPATMINER` performs *refactoring-aware, fine-grained* change pattern mining in the commit history of Python systems.

Our findings and tools have actionable implications for several audiences. Among others, they (i) advance our understanding of repetitive changes that the ML developers perform which helps the research community to improve the science and tools for ML software evolution, (ii) provide a rich infrastructure to automate and significantly extend the scope of existing studies on ML systems [37, 38, 75], (iii) help tool builders comprehend the ML developers' struggles and desire for automation, (iv) provide feedback to language and API designers when introducing new ML constructs, and (v) assist educators in teaching ML software evolution.

This paper makes the following contributions:

- (1) To the best of our knowledge, we conducted the first and the largest study on fine-grained 28,308 code change patterns on ML systems. We identified code changes patterns. We applied *thematic analysis* on 2,500 most popular patterns and categorized them into 22 fine-grained change pattern themes that reveal 4 major trends.
- (2) We designed and conducted a *survey* with 650 open-source ML developers to provide insights about the reasons motivating those changes, the current practices of applying those changes, and their recommendation for tool builders.
- (3) We developed novel tools to collect fine-grained change patterns applied in the evolution history of Python-based ML systems. We applied these tools on 1000 open-source projects hosted on GitHub. We make the collected information and tooling publicly available for reuse [21] so that we enable further research.
- (4) We present an *empirically-justified* set of *implications* of our findings from the perspective of four audiences: researchers, tool builders, language designers, and ML developers.

2 MOTIVATING EXAMPLE

Listing 2: Specifies the device (CPU) for operations executed in the context and move method `_init_model` to parent class

```

1 class _FERNeuralNet():
2 +   def _init_model(self):
3 +       with tf.device('/cpu:0'):
4 +           B, H, T, _ = q.get_shape().as_list()
5 ...
6 class TimeDelayNN(_FERNeuralNet):
7 -   def _init_model(self):
8 -       B, H, T, _ = q.get_shape().as_list()

```

The code change shown in Listing 2 specifies the hardware device using `tf.device()` (line 3) for the `TensorFlow` operation in line 4. `tf.device()` is a Context Manager [64] from the ML library, `TensorFlow`. This is a fine-grained code change and the developer has interleaved this with a *Pull up Method* refactoring that pulls `_init_model` from `TimeDelayNN` into the parent class `_FERNeuralNet`.

Is specifying hardware device for `TensorFlow` operations a *pattern*? How frequent is this pattern? Do developers need tool support to recommend and automate this *code change pattern*? We consider this *fine-grained code change instance* a repeated *pattern* if a similar code change was performed in the history of this project or any other project. Researchers have proposed advanced techniques to mine such fine-grained change patterns from the commit histories [58, 59]. However, these techniques are inapplicable to mine

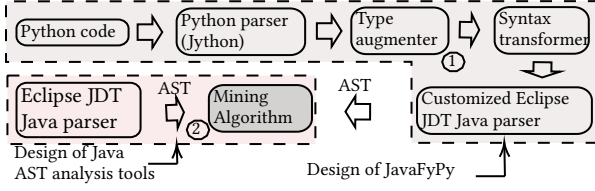


Figure 1: Design of JAVAfYpY and existing AST analysis tools

the *fine-grained code change patterns* shown in Listing 2 because (1) their techniques mine code change patterns for Java, and (2) they do not account for overlapping refactorings.

Researchers [51, 52, 78] have shown that developers often interleave many programming activities such as bug fixes, feature additions, or other refactoring operations, and often these changes overlap [54] (as shown in Listing 2). Such overlapping changes and refactorings can easily obfuscate existing fine-grained change pattern mining tools [58, 59] because they do not account for these changes when matching code across the commit. For example, CPAT-MINER [58] does not match the method body of `_init_model` in the class `_FERNeuralNet` (lines 3–4) to the body of `_init_model` in the class `TimeDelayNN` (line 7) as they are in different locations and different files. This lack of *refactoring awareness* is a serious limitation of existing pattern mining algorithms because they can miss several concrete instances of change patterns that are obfuscated by overlapping refactorings.

Re-implementing the existing Java AST mining tools for Python will require a significant amount of development effort. It is also neither feasible nor sustainable as researchers are continuously implementing new Java AST mining tools or improving existing tools. For this purpose, we propose JAVAfYpY, a technique to adapt existing Java AST mining tools to Python that leverages the similarity between the Java and Python abstract syntax trees (AST). We use JAVAfYpY to adapt the state-of-the-art fine-grained change pattern mining tool, CPATMINER [58], to Python. To make CPATMINER *refactoring aware*, we adapt the state-of-the-art Java refactoring inference tool, *RefactoringMiner* [85] (known as RMINER), to Python and integrate it with CPATMINER as R-CPATMINER. Particularly, the code-block mapping pairs (i.e., two versions of the same code-block in a method before and after the change) reported by RMINER are provided as input to CPATMINER. R-CPATMINER mines change patterns in Python software systems in a refactoring-aware manner.

3 TECHNIQUE

Most of the current code change mining tools (i.e. AST mining tools) are conceptually language-independent because they operate upon the abstract syntax trees (AST) only. However, their implementation is bound only to Java. To overcome this practical implementation limitation, we propose a very pragmatic solution - JAVAfYpY, a technique that transforms the input Python program to an AST that can be processed by the mining algorithm of existing Java AST analysis tools. JAVAfYpY will fast-track researchers and tool builders by making the AST-based mining tools implemented for Java programs applicable for Python programs. Thus, it will save several development-hours of work required for re-implementing these techniques. As shown in ① in Figure 1, JAVAfYpY takes a

Python code as an input and produces an AST, that can be used in mining algorithms of Java AST analysis tools. To achieve this, JAVAfYpY first transforms the Python code to AST and enriches the AST by augmenting type information. Then, the *Syntax transformer* maps the corresponding Java concrete syntax to the AST nodes. The Java parser (Eclipse JDT) uses it to produce the final AST. Eclipse JDT is the most popular Java parser used in AST mining research tools. Therefore, we selected Eclipse JDT as the parser that produces the final AST. This *enhanced and enriched* AST can be processed by the mining algorithms of Java AST analysis tools. Tool builders and researchers can use JAVAfYpY, and extend their tools for Python.

3.1 Python code transformation

As shown in Figure 2 JAVAfYpY first *parses* the input Python program to an AST. We define an AST as:

Definition 3.1. (AST) Let T be an AST. T has one root. Each node $N_i \in T$ has a parent (except the root node). Each node ($N_i \in T$) has a sequence of child nodes (denoted by C_{N_i}). Number of nodes in the sequence C_{N_i} is denoted by $\text{Length}_{C_{N_i}}$. Each node N_i is a specific syntax category known as AST node kind, $\text{Kind}_{N_i} = \{\text{Assignment Statement, For statement, Method Invocation ...}\}$.

We leverage the syntactic similarity between Python and Java to adapt the Java AST analysis tools to Python. We thoroughly studied the Java and Python language specifications [62, 67] and *mapped* the Python AST node kinds to those in Java based on the description in the specifications.

Definition 3.2. (Mapped AST node) Let T_j be a Java AST and T_p be a Python AST. $N_j \in T_j, N_p \in T_p$. We state that N_j and N_p are mapped AST node kind, if N_j and N_p maintain a structural similarity. Mapped node of node N_p is denoted by $M(N_p) = N_j$.

We found three kinds of mappings namely, *Identical AST node*, *Nearly identical AST node*, and *Unique AST node*.

Definition 3.3. (Identical AST node) Let C_{N_j} be a sequence of child AST nodes of a parent Java node N_j and C_{N_p} be a sequence of child AST nodes of a Python node N_p . We state that N_j and N_p are identical AST nodes if (i) $M(N_p)$ is N_j , and (ii) $\forall N_i \in C_{N_p} : M(N_p) \in C_{N_j}$.

(1) **Identical AST node** (Definition 3.3) - 72% of the Python’s AST node Kinds can be identically mapped to a Java’s AST Node. For example, we mapped Python’s `If` to Java’s `If Statement` and mapped Python’s `Assign` to Java’s `AssignmentStatement`.

Definition 3.4. (Nearly Identical AST node) We state that N_j and N_p as nearly identical AST nodes, if N_j and N_p meets conditions: (i) $M(N_p)$ is N_j , and (ii) $\exists N_i \in C_{N_p} : M(N_i) \notin C_{N_j}$.

(2) **Nearly identical AST node** (Definition 3.4) - 18% of Python’s AST Node kinds could be *partially* mapped to those of Java. For instance, both Python and Java provide the `for` construct to iterate over a collection, however unlike Java, Python allows to iterate over multiple variables (see the `for` loop in Figure 2), thus AST of Python `for` loop contains additional child AST node kinds.

Definition 3.5. (Unique AST node) Let N_p be a Python AST node. We state that N_p is unique to Python, if there is no mapped AST node in T_j . i.e., $(M(N_p) \notin T_j)$.

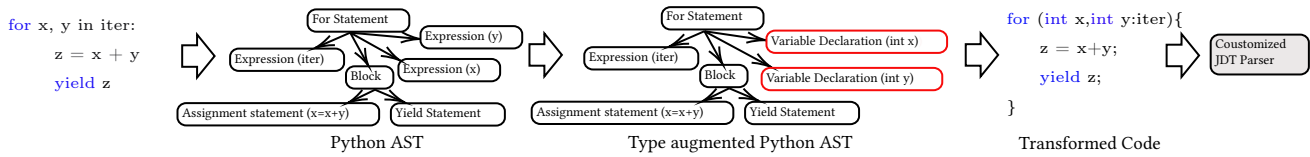


Figure 2: An example Code transformation performed by JAVAfYPy

(3) **Unique AST node** - 10% of the Python nodes had no Java counterpart. For instance, Java does not support `list` comprehensions or `yield` statement (as shown in Figure 2)

As we can observe, Java and Python syntax significantly overlaps. As shown in ② in Figure 1, AST mining tools like CPATMINER and RMINER parse the input program to Eclipse JDT AST. To adapt their tools to Python with JAVAfYPy, tool builders or researchers simply need to migrate their Java parser to our technique (JAVAfYPy). After that, we can simply reuse tools' AST-based mining algorithms to analyse 72% of the *Identical AST node* kinds, and modify the current implementation to accommodate the 18% *Nearly identical AST node* kinds and add brand-new capabilities (often involving adding new *visitors*) for handling the 10% *Unique AST node* kinds. After the changes, the tools take Python code as an input and infer the results, thus adapting Java AST mining tools to Python.

Figure 2 shows an example of the code transformation steps (shown in Figure 1) that JAVAfYPy performs automatically. The Python code snippet in Figure 2 contains all three AST node kinds: *Identical AST node* ($z = x + y$), *Nearly identical AST node* (`for` loop), and *Unique AST node* (`yield z`). The *Java parser* first constructs the AST of the code snippet, then the *Type Augmenter* augments the AST with type information by adding *Variable Declaration* nodes. This step is important because the Java-based AST mining tools [58, 85] rely on the syntactic richness that the Java language offers. Unlike Python, Java programmers have to explicitly declare the types of variables, fields and methods. To add this syntactic richness to the input program, JAVAfYPy augments the AST of the input program with type information (shown in Figure 2 as red nodes). We obtain the type information from PYTYPE [30], the state-of-the-practice type inference tool for Python developed by Google, which is widely adopted by the Python community. As the last step, *Syntax Transformer* transforms the AST to code and passes it to our customized *Eclipse JDT* parser which we extended to parse *Nearly identical AST node* kinds and *Unique AST node* kinds.

Can JAVAfYPy effectively transform all *Identical*, *Nearly Identical*, and *Unique AST node* kinds? We evaluated this empirically with 14 popular Python projects including *TensorFlow*, *PyTorch*, *Keras*, *NLTK*, *Scikit learn*, *Scipy*, and *NumPy* that comprise 23K Python files and 2.9M SLOC. We checked whether all AST nodes, i.e., *Identical AST nodes*, *Nearly identical AST nodes* and *Unique AST nodes* were successfully mapped and transformed to the output AST of JavafYPy. We achieved this by transforming all of the Python files in the projects, which had 12M Python AST nodes. This confirms that JAVAfYPy can effectively transform any input Python program to an Eclipse JDT AST format.

3.2 Refactoring Aware Change Pattern Mining

3.2.1 Adapting CPATMINER. CPATMINER [58] is the state-of-the-art code change pattern mining tool that uses an efficient graph-based representation of code changes to mine previously unknown fine-grained changes from git commit history. It iterates over changed methods in each commit and uses *Eclipse JDT Java parser* [27] to generate an AST of Java source code. Then, its mining algorithm builds program-dependence graphs for each AST node independently and then merges the graphs to create one big graph, called **change graph**. CPATMINER builds *change graphs* for each changed method, and it represents the before and after a source code change that can be used to mine code change patterns. Since 72% of the Python AST node kinds overlap identically with those in Java, we reused most of the capabilities for building the change graphs. We added new capabilities in CPATMINER to create program-dependence graphs for *Unique AST nodes*, and modified the existing capabilities of *Nearly identical AST nodes*. Overall, we extended CPATMINER with 2% extra code lines due to the new or modified capabilities, and reused the rest of the code. While this ratio might be different when adapting other tools, it showcases the merit of JAVAfYPy to reuse Java AST-based mining tools for Python.

3.2.2 Introducing Refactoring Awareness. As discussed in Section 2, CPATMINER [58] does not account for the overlapping refactorings applied in the commit. These refactorings move code blocks between methods or change the method signature, making it hard to match the changed code blocks. Thus missing the opportunities to build *change graphs* for the obfuscated changes. To overcome this, we made the CPATMINER refactorings aware by integrating it with RMINER [85]. We used JAVAfYPy to adapt RMINER and use it to detect 18 refactoring kinds that move code blocks. RMINER uses AST-based statement matching algorithm to match classes, methods, and statements inside method bodies, thus helping us match moved code blocks. We consult the authors of RMINER and extend its statement matching algorithm to reason about the *Unique* and *Nearly identical* AST Node kinds. For example, Listing 3 shows a variable rename refactoring in List Comprehension, a Python *Unique AST node* kind of the project "Deepmedic" detected by Python-adapted RMINER.

Listing 3: Commit 8d4be555 in DeepMedic: Variable rename in List Comprehension detected by Python-adapted RMiner

```
1 - indices = [layerNum - 1 for layerNum in layers_norm]
2 + indices = [layer_num - 1 for layer_num in layers_norm]
```

We use Python adapted RMINER to accurately match the moved code blocks. We extended CPATMINER to build change graphs for the code block pairs reported by RMINER. Hence, CPATMINER no longer misses obfuscated code-blocks that contain fine-grained changes. We developed the tool **R-CPATMINER**, to efficiently and

accurately mines source code change patterns in the version histories of Python software systems, in a refactoring-aware manner.

4 RESEARCH METHODOLOGY

We prefix all the adapted tool names with *Py* to disambiguate the tool names from their Java counterparts. We first evaluate the effectiveness of the tools we developed (or adapted). Then, we use our reliable and validated tools, to explore the repetitive code changes applied in Python ML Systems. For this purpose, we answer three research questions:

RQ1. *What are the frequent code change patterns in ML code, and what patterns need automation?* To answer this research question, we triangulate complementary empirical methods, as shown in Figure 3. (i) We mined 1000 repositories using R-CPATMINER and extracted 28,308 patterns, (ii) We applied thematic analysis on 2,500 patterns, (iii) We sent a survey to 650 ML developers to seek their opinion on automating the identified code change patterns.

RQ2. *How does the refactoring awareness improve the pattern mining over the baseline CPATMINER?* R-CPATMINER performs refactoring aware change pattern mining, thus improving baseline CPATMINER. Compared to CPATMINER, does R-CPATMINER extract (i) more change graphs? (ii) more code change patterns? and (iii) more code instances per pattern?

RQ3. *What is the runtime performance of R-CPATMINER, PyCPATMINER, and PyRMINER?* To answer this, we compare the execution time of the Python adapted tools with their Java counterparts.

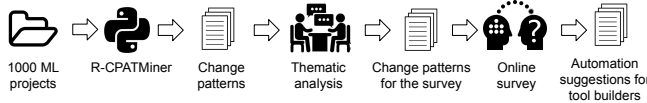


Figure 3: Schematic diagram of the research methodology to answer RQ3

4.1 Subject systems

Our corpus consists of 4,166,520 commits from 1000 large, mature, and diverse ML application systems, comprising 58M lines of source code and 150K Python files, used by other researchers [22] to understand the challenges of evolving ML systems. This corpus [22] is shown to be very diverse from the perspective of Python files, LOC, contributors, and commits. They vary widely in their domain and application, include a mix of frameworks, web utilities, databases, and robotics software systems that use ML. Further we added low-level ML libraries [10] such as Scipy, SpaCy, and high-level ML libraries [10] such as Keras, PyTorch, Caffe, NLTK, and Theano to our subject systems. This ensures our dataset is representative and large enough to answer our research questions comprehensively.

4.2 Static Analysis of Source Code History

4.2.1 Change pattern identification: Running R-CPATMINER on the ML corpus extracted 28,308 unique code change patterns, where 58% of them have code change instances in multiple projects, 63% of them have been performed by multiple authors.

Since the mined patterns are numerous, we followed the best practices from Negara et al. [53]. They ordered the patterns along three dimensions - by frequency of the pattern (F), by the size of the pattern (S), and by $F \times S$. Since the repetitive changes done by several developers and projects are stable [59] and have a higher chance of being automated, we also considered the number of projects and authors as extra two dimensions. Then we ordered the mined patterns along all five dimensions. Then, two of the authors who have more than three years of professional software development experience and extensive expertise in software evolution, manually investigated the top 500 patterns for each of the five dimensions and identified meaningful code patterns, i.e., the patterns that can be described as high-level program transformations.

Two authors of the paper manually analyzed each change pattern, to identify the high-level programming task performed in the change patterns. Following the best practices guidelines from the literature, the authors used negotiated agreement technique to achieve agreement [14, 88]. Two authors of the paper independently coded the change patterns carefully and assigned one or more descriptive phrases (i.e., codes) to the patterns. Both authors conducted the initial meeting after coding around 25% of the data (the suggested minimum size is 10% [14]). During the meeting, the authors carefully discussed the coding process of all the patterns. Also, they negotiated any disagreements between the assigned codes and the patterns that cannot be described as high-level program transformations. After 80% inter-coder agreement was achieved (recommended inter-coder agreement level ranges from 70% to more than 90% in the literature [14]), the two authors independently coded the remaining change patterns. This process identified all the patterns for which the two authors were able to agree upon the underlying meaning of the pattern. After the coding finished, the authors held another meeting in order to finalize the codes and extract themes. Themes capture something important about the data in relation to the meaning of the pattern. It also represents some level of patterned response or meaning within the data set [12]. The two authors reviewed the initial themes against the data several times and refined their names and definitions until they both agreed that there were no further refinements possible. We identified four trends (themes) of patterns based on their structural similarity at the statement and expression level, namely (i) *transform to Context managers*, (ii) *convert FOR loops to domain specific abstraction*, (iii) *update API usage*, and (iv) *use advanced language features*.

4.3 Qualitative Study

The most reliable way to understand the motivations and challenges associated with repetitive code changes is to ask the developers who performed them. To achieve this, we surveyed 650 developers who performed the identified change patterns.

4.3.1 Contacting the developers: We contacted the developers performing repetitive code changes that we considered worthy of further investigation by sending an email to the addresses provided in their GitHub account. The body of each email message was automatically generated by the application we developed, and included the following information:

- Introduction to the research team and the purpose of the study.
- A plot of number of repetitive changes done in the project.

- A link to the frequent repetitive changes done in the respective project so that the developer can use it as an educational resource.
- The following four questions for the developer:
 - Q1. What are the reasons for performing the code changes above?
 - Q2. How often do these code changes happen in ML codes?
 - Q3. How often have you manually performed this kind of change?
 - Q4. Would you like to have this change automated by a tool?

The first question aims at discovering actual motivations behind a code change as expressed by the developers themselves. The second question focuses on the frequency of performing the code change on ML codes, and it helps to examine the need for ML-specific IDEs tools. The last three questions aim at understanding whether developers trust and use tool support for performing the code changes. This is important, as there is relatively low IDE support for performing code changes in ML codes [22, 91]. A sample email is available on the companion website [21].

In total, we sent 650 emails to developers, out of which 97 responded, bringing us to a 15% response rate. This is significantly higher than the usual response rate achieved in questionnaire-based software engineering surveys, which is around 5% [79].

5 RESULTS

5.1 Repetitive changes in ML systems (RQ1)

5.1.1 Characteristics of patterns mined by R-CPATMINER. We executed R-CPATMINER on our corpus described in Section 4.1 containing 1.5M changed source code files, comprising of over 490M lines of source code. For these changed files R-CPATMINER produced 349,406 change graphs with a total of 4.7M nodes. The tool extracted 28,308 unique code change patterns, where 63% and 58% of them are performed by multiple authors and in multiple projects, respectively. Figure 4 shows the *degree of sharing* of the patterns amongst developers and projects. We observed that 53% of the developers who performed the code change patterns share 100% of their change patterns with other developers, 79% share at least 50% of their patterns with others, and 91% share at least 10% of the patterns. Moreover, 36% of the projects share 100% of their patterns with other projects, 60% of them share at least 50% of their patterns with others, 91% of the projects share at least 10% of the patterns. This shows that R-CPATMINER extracts patterns that are pervasive amongst the developers and projects.

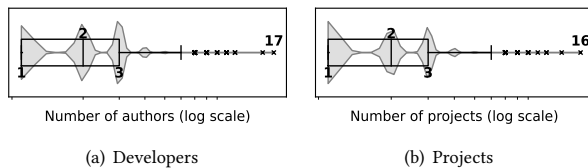


Figure 4: Degree of sharing of patterns amongst developers and projects

5.1.2 Discovering pattern trends. Understanding code change patterns that the ML developers perform is important to advance science and tooling in ML software evolution. Our thematic analysis and developer surveys reveal 22 previously unknown repetitive change patterns groups where the developers ask for automation.

Amongst these patterns, we identified four major trends based on their structural similarity (i.e., expression- and statement-level):

- (1) *transform to Context managers*- 1237 instances
- (2) *convert for loops to domain specific abstraction*- 239 instances
- (3) *update API usage*- 166 instances
- (4) *use advanced language features*- 415 instances

Next, we summarize and triangulate results obtained from source code mining, thematic analysis, and developer surveys.

Note: We use real-world code examples to describe frequent change patterns. The examples use identifiers *tf*, *np*, and *torch* as aliases of ML libraries *TensorFlow*, *NumPy*, and *PyTorch*. Due to space limitations, we provide few code examples. Our companion website [21] presents a curated repository of exemplars for each pattern, as well as *all* the instances for each pattern.

5.1.3 Trend 1 - Transform to Context managers: A Python Context manager is an abstraction for controlling the life-cycle for a code block. It declares the methods `__enter__` (initialization), and `__exit__` (finalization) which together define the desired runtime environment for the execution of a code block. The code block needs to be surrounded in a `with` statement [68] that invokes the Context manager. We observed 1,237 change instances belonging to eight pattern groups (P1–P8 in Table 1) where developers move code blocks into `with` statements and use *Context managers*.

Listing 4: Commit dfb7520c in Pytorch: Disable gradient

```
1 -input.grad.data.zero_()
2 +with torch.no_grad():
3 +     input.grad.zero_()
```

Listing 4 is an example of pattern P2 (Disable or enable gradient calculation). The survey respondent S21 said, “when we do not need gradient computation in a DL network (using `Tensor.backward()`), it is important to disable the gradient calculation globally to reduce memory consumption and increase speed”. The context manager `torch.no_grad()` from *PyTorch*, creates an execution environment for the code in line 3 and disables the gradient calculation. Like wise, the patterns (P2–P8) in Table 1 create new execution environments.

Listing 5: Commit 02ccf29b in tensorflow/datasets: Move Context managers that used to read data to with statement

```
1 -file_ = tf.gfile.GFile(label_path)
2 -dataset = csv.DictReader(file_, delimiter="")
3 +with tf.gfile.GFile(label_path) as file_:
4 +     dataset = csv.DictReader(file_, delimiter="")
```

We will now explain the most populous pattern group P1 (*Read, write, traverse data*) which moves an *already existing* Context manager within a `with` statement (see Listing 5). The survey respondent S11 said, “When we use Context managers in `with` statement, the required resources are allocated and released precisely.” Line 3 in Listing 5 uses the Context manager `tf.gfile.GFile` which handles I/O operations. The developers do not need to handle I/O operations such as file open (initialization) and file close (finalization) when they use the Context manager *inside* a `with` statement. However, if developers use the Context manager as a function call (see deleted line 1 in Listing 5), they need to handle all the initialization and finalization logic [1]. Hence, if a Context manager is used as a function call: (i) API misuses often happen, and allocated resources might not be managed efficiently [74, 83], (ii) developers need to

Table 1: Triangulating source code mining results with survey responses: 4 major trends, the pattern groups for each trend, and whether that pattern is specific to ML code (column ML). Column I shows the number of instances for each pattern. Column R shows number of survey respondents. Next columns indicate their responses to survey Q2 (How often these changes happen in ML code?) and Q3 (How often have you manually performed this change?), with frequency: Very Often (VO), Often (O), Rare (R), and Never (N). Q4 (Would you like to have this change automated by a tool?) response: Yes, No, Already Automated (AA).

Static Analysis					Survey Responses												
Trend	P	Pattern	ML ¹	I	R	Q ₂				Q ₃				Q ₄			
						VO	O	R	N	VO	O	R	N	Yes	No	AA	
Move to with statement and use Context managers	P1	Read, write, traverse data	×	467	6	67%	17%	-	17%	83%	-	17%	-	100%	-	-	
	P2	Disable or enable gradient calculation	✓	92	5	100%	-	-	-	80%	-	20%	-	100%	-	-	
	P3	Swap ML training devices	✓	21	3	100%	-	-	-	67%	33%	-	-	67%	33%	-	
	P4	Change name and variable scopes in DL networks	✓	106	7	57%	43%	-	-	57%	29%	14%	-	57%	43%	-	
	P5	Execute dependencies of a Tensorflow graphs	✓	57	3	100%	-	-	-	100%	-	-	-	33%	67%	-	
	P6	Temporarily change configurations of libraries	×	59	2	-	50%	50%	-	100%	-	-	-	50%	50%	-	
	P7	Transform to context managers in pytest	×	335	6	33%	50%	17%	-	67%	33%	-	-	67%	33%	-	
	P8	Use context managers to open temporary directories	×	100	6	33%	50%	-	17%	67%	-	33%	-	83%	17%	-	
Dissolve for loops, into domain specific abstractions	P9	Transform to optimized operations in NumPy	✓	179	10	100%	-	-	-	90%	-	10%	-	100%	-	-	
	P10	Transform to operations in List or Dictionary	×	24	3	67%	33%	-	-	33%	-	67%	-	33%	67%	-	
	P11	Transforming to Python built in functions	×	15	2	-	100%	-	-	-	-	100%	-	100%	-	-	
	P12	Transform functions in String	×	14	2	-	50%	50%	-	100%	-	-	-	100%	-	-	
	P13	Transform to set operations	×	7	2	100%	-	-	-	50%	-	50%	-	100%	-	-	
Update API usage, (212)	P14	Migrating to APIs ML libraries	✓	26	5	100%	-	-	-	100%	-	-	-	100%	-	-	
	P15	Transform Matrix	✓	82	6	83%	17%	-	-	100%	-	-	-	33%	50%	-	
	P16	Change data visualization	✓	28	2	50%	50%	-	-	50%	-	50%	-	50%	50%	-	
	P17	Composite ML APIs	✓	22	5	100%	-	-	-	100%	-	-	-	100%	-	-	
	P18	Update Container	×	115	5	60%	-	40%	-	60%	-	40%	-	80%	20%	-	
P19	Update Type of Matrices	✓	16	4	100%	-	-	-	50%	-	50%	-	50%	50%	-		
Use advanced Language Features	P20	Simplify conditional statement	×	24	3	-	33%	67%	-	67%	-	-	33%	33%	-	67%	
	P21	Migrate from Dict, Set, List constructors to literals	×	42	4	-	25%	75%	-	25%	-	-	75%	25%	-	75%	
	P22	Transform to Python List, Dict, or Set Comprehension	×	349	6	-	33%	67%	-	33%	-	-	67%	33%	-	67%	

¹ ML specific Patterns: i.e., patterns related to ML techniques. More than 80% developers confirmed they happen in ML very often or often

update all the initialisation and finalisation code when they update the library versions (if the APIs have changed). The respondent S13 said, “I envision IDEs that automate moving Context managers to with statements.”

Table 1 tabulates the results for each major trend and pattern group, and shows survey responses for each pattern group. 90% of the survey respondents who performed Trend-1 changes confirmed that they move to with statements very often (VO) or often (O). All respondents perform the code transformation manually, and 74% of the respondents requested automation in their IDEs.

5.1.4 Trend 2 - Convert for-loops into domain specific abstraction: Listing 1 shows one such example where the developer uses `np.sum` from NumPy [61] instead of using `for` loop to compute the sum of elements in a `list`. Developers often perform this change to enhance the performance and code readability. Survey respondent S22 who performed pattern P9 said, “Sometimes, Python `for` loop is a real performance killer. I want my IDEs to suggest the optimized APIs from ML libraries that I can use instead of loops”. Moreover, as alternatives to `for` loops, developers use (i) List or Dictionary operations (P10), (ii) Built in Python functions (P11), (iii) Python `String.join()` (P12), and (iv) Set operations `union` and `intersection` (P13). Python is an interpreted language, thus compiler level optimisations do not happen in Python. The respondent S24 said, “For other programming languages, I might expect the

compiler to optimize this type of loop, so yes, I would be interested in a suggestion by the IDE.” As shown in Table 1, 95% of the respondents who performed Trend 2 confirmed they do this very often (VO) or often (O) in ML code. All the respondents manually perform the change, and 89% of the respondents requested automation support.

5.1.5 Trend 3 - Update API usage: Listing 6 shows an example API migration where the developer uses a readily-available `np.mean` instead of computing mean of the list `first_occ`. Survey respondent S35 said, “NumPy offers efficient arrays and APIs for computational operations, tools that inspect the code and suggest NumPy APIs are very much needed.”

Listing 6: Commit 8592777b in inspirehep/magpie: Migrate API to NumPy

```
1 -return sum(first_occ) / len(first_occ)
2 +return np.mean(first_occ)
```

Matrix transformations (P15) such as `transpose`, `broadcast`, `squeeze`, and `unsqueeze` are frequent in ML projects. The respondent S31 said, “bugs due to wrong matrix shapes are hard to detect and prevalent in ML systems. I like to have tools that identify these bugs and broadcast the matrices to correct shape”. Chen et al. [15] observed developers change the bit size of matrices to get good trade-off between training time and accuracy of the predictions.

Congruent to this, we observed developers update the type of matrices, e.g., migrating to an int64 matrix (P19) as shown in Listing 7.

Listing 7: Commit 8f9cabbf in stellargraph/stellargraph: Update matrix type

```
1 -np.hstack((A.row[:, None], A.col[:, None]))
2 +np.hstack((A.row[:, None], A.col[:, None]))
   ↪ .astype(np.int64)
```

Other patterns include changing data visualization with Matplotlib [48](P16), e.g., between drawing all the plots in one figure vs using an individual figure for each plot. Another pattern is using composite ML APIs (P17). Developers often traverse datasets multiple times (which is inefficient). A more efficient solution is to apply a composite operation. In Listing 8, developers compute dot product on three matrices instead of applying a `np.multi_dot`.

Listing 8: Commit 180646fa in scikit-learn: Composite APIs

```
1 -denominator = np.dot(np.dot(W.T, W), H)
2 +denominator = np.linalg.multi_dot([W.T, W, H])
```

ML libraries offer several optimized containers (e.g., NumPy arrays, Tensors) for data processing. Updating containers, e.g., from Python's `List` to `NumPy.Array`, is another frequent change in ML systems (P18). Ketkar et al. [42] discovered that in Java code, these type migrations are more common than rename refactorings. Moreover, Table 1 shows that 85% of respondents who performed Trend 3, perform it very often (VO). All the respondents manually perform these changes, and 70% of respondents sought automation in IDEs.

5.1.6 Trend 4 - Use advanced language features: Python offers powerful features: (i) functions [65] such as `bool` and `instance` that can be used to simplify a conditional statement (P20), (ii) literals such as `[]`, `{}`, `()` to efficiently create containers instead of using constructors such as `list()`, `dict()`, `tuple()` (P21) (see Listing 9). (iii) Python comprehension [66] to make code concise and inline `for` loops (P22). Researchers [10, 45] observed ML code extensively operates on data, which results in expressions that are longer and more complex than in traditional systems. Good software engineering principles [26, 47] require that developers change the code to make it concise and readable. However, 69% of the survey respondents who performed Trend 4 changes confirmed they rarely perform this in their project, and 30% of developers sought automated help.

Listing 9: Commit 15d7634d in RasaHQ/rasa: Use set literals instead of set constructor

```
1 -set(utils.module_path_from_instance(p)
   ↪ for p in agent.policy_ensemble.policies)
2 +{utils.module_path_from_instance(p)
   ↪ for p in agent.policy_ensemble.policies}
```

5.2 Improvements caused by Refactoring Awareness (RQ2)

5.2.1 Impact of refactoring awareness. To answer this question, we executed the `PyCPATMINER` and `R-CPATMINER` on the study corpus and compared the results. We compared, number of change graphs, number of patterns, and distribution of code instances per pattern reported by both tools.

As shown in Table 3, `R-CPATMINER` processed 16% more changed methods, 0.1B more AST nodes than `PyCPATMINER`. `PyCPATMINER` builds one change graph for each mapped code block pair (i.e., before and after the changed method body). Therefore, `R-CPATMINER` produces 16% more *change graphs*, thus confirming the value of de-obfuscating change graphs that were previously obfuscated by refactoring. Then the `R-CPATMINER` mines all the generated *change graphs* and extracts repeated isomorphic sub-graphs as patterns. `CPATMINER` uses minimum frequencies of repeated subgraphs σ to be three adheres to the Rule of Three [73], a standard recurrence measure in pattern analysis. Therefore, all the patterns contain at least three code instances. We compared the number of patterns generated by the `R-CPATMINER` and `PyCPATMINER` and observed that `R-CPATMINER` captures 15% more patterns than `PyCPATMINER`.

We also compared the distributions of the number of code instances per pattern in both `PyCPATMINER` and `R-CPATMINER`. To assess if there is a statistically significant difference in distributions of the number of code instances per pattern reported by `PyCPATMINER` and `R-CPATMINER`, we applied the *Wilcoxon Signed-Rank* test on the paired samples of number of code instances of each pattern. The test rejected the null hypothesis that the density of code instances of the pattern produced by `PyCPATMINER` is more than it is in `R-CPATMINER` at the significance level of 5% (p-value = 1.12×10^{-10}). We used the *Hodges-Lehman estimator* to quantify the difference between the two distributions, as it is appropriate to be used with the *Wilcoxon Signed-Rank* test. The value turned out to be 1, which is equal to the estimated median of the difference between the number of code instances per pattern from `PyCPATMINER` and `R-CPATMINER`. Therefore, `R-CPATMINER` extracts more code change instances per pattern than the `PyCPATMINER`.

5.2.2 Evaluating the precision of PyRMINEr. It is important for `PyRMINEr` to have a high precision as we use it to first match the refactored code blocks that we then pass to `R-CPATMINER` to build change graphs. We first identified 18 refactoring kinds that obfuscate fine-grained changes, i.e., the refactorings that change method signatures or shift the method bodies. First, we executed the `PyRMINEr` on our study corpus and chose a statistically significant random sample of refactoring instances for each refactoring kind. Hence, using a t-test, we conclude with 95% confidence that the precision of the refactoring detection is only $\pm 5\%$ for each refactoring kind, as shown in Table 2.

Table 2: Precision of `PyRMINEr` per refactoring kind

Refactoring Kind	Precision (#TP)	Refactoring Kind	Precision (#TP)
Rename Method	96.32% (183)	Move And Rename class	96.67% (116)
Move Method	96.3% (156)	Move class	100% (160)
Pull Up Method	86.46% (83)	Extract Class	98.82% (84)
Push Down Method	89.36% (84)	Extract Subclass	100% (45)
Extract Superclass	95.24% (40)	Parameterize Variable	87.23% (82)
Split Parameter	91.66% (22)	Move & Rename Method	88.08% (133)
Rename Class	99.2% (124)	Remove Parameter	97.56% (160)
Reorder Parameter	98.11% (104)	Rename Parameter	93.59% (146)
Add Parameter	100% (192)	Merge Parameter	91.07% (51)

Two of the authors that have more than three years of software development experience and extensive expertise in software evolution manually validated the refactorings reported by `PyRMINEr`.

Most cases were straightforward and thus were validated individually, but both authors inspected some challenging cases to reach an agreement. In total, we validated 2,062 unique refactoring instances, out of which 1,965 were true positives and 97 were false positives. This achieves an average precision of 95%, which is close to the precision of the original Java-RMINER (99.6%). This also shows the effectiveness of JAVAFYPY to adapt Java AST-analysis tools to Python. We release all the validated refactoring instances on our companion website [21]. To the best of our knowledge, this is the largest to date Python data-set of validated refactoring instances.

Recall indicates the proportion of actual refactorings identified by PYRMINER. Java researchers use previously formed unbiased oracles to compute recall of tools [42, 85]. To the best of our knowledge, there is no such unbiased refactoring oracle for Python. Considering the complexity of building an oracle and our main focus (to retrieve correct code-block mappings to make PYCPATMINER refactoring-aware), we leave computing recall as future work. However, we observed R-CPATMINER detects 16% more change-graphs (i.e., mapped code blocks) and 15% more patterns than the PYCPATMINER, which strongly indicates that the PYRMINER has a satisfactory recall.

5.3 Runtime performance of R-CPATMINER, PYCPATMINER, and PYRMINER (RQ3)

A fast execution time persuades the users that the adapted tools can be run on their systems (even on a consumer laptop) in an acceptable amount of time. It will enable the creation of bigger datasets of repetitive code changes to strengthen the validity of empirical research or train learning-based code recommendation systems, as well as enable innovative uses of code change pattern mining at commit time. Therefore, we evaluated the execution times of the adapted tools and compared them with the Java counterparts.

Table 3: Analysed data set and execution time.

	Java CPATMiner	Py- CPATMiner	R-CPATMiner
Total changed methods	824K	1M	1.16M
Total AST nodes of changed methods	92M	4.5B	4.6B
Total changed graph nodes	8M	4M	4.7M
Total patterns	17K	24K	28K
Execution time	<8hours	<12hours	<19hours

¹ The data of the Java-CPATMiner is obtained from its original paper [58]

To measure the execution time of the tools, we executed the Python adapted RMINER, CPATMINER, and R-CPATMINER on a large corpus and compared the execution time with the Java version of the tools. We executed each tool separately on the same machine with the following specifications: Intel Core i9 CPU @ 2.90GHz, 32 GB DDR4 memory, 1 TB SSD, macOS 10.14.6, and Java 13.0.1 x64.

First, we record the running time of the *type inference* tool, *PyType* [30] (version 2020.10.08). *PyType* is decoupled from JAVAFYPY. Therefore, we computed the execution time of *PyType* separately. *PyType* took on mean 360ms and 61ms on median for type inference of all changed files in a commit. We pushed the inferred type information of all the changed files in all the studied projects' commits to a repository [69] in Github. Similar to *Typeshed* [86], a type repository of Python library APIs that the library clients use for

type annotations, researchers can fast-track their analysis by simply reusing this inferred type information .

To record running time of PYRMINER, we followed the steps used by Tsantalis et al. [85] for computing the running time of their Java-RMINER. We recorded the time required for parsing the source code of the commit (and its parent), and the time required for detect refactorings. Our analysis shows that PYRMINER takes 55ms on median and 296.32ms on mean to process a Python commit. Tsantalis et al. [85] found that the Java version of RMINER takes 44ms on median and 253ms on mean to process one Java commit. Therefore, PYRMINER takes reasonable time overhead for the additional processing (i.e., AST transformation), and will not impact the primary goal of RMINER, i.e., create larger refactoring datasets to strengthen the validity of empirical studies or enable novel applications of refactoring mining [85].

Table 3 shows the size of processed data and the execution time of the tools, Java-CPATMINER, PYCPATMINER, and R-CPATMINER. The running time of the PYCPATMINER and R-CPATMINER is less than 12 hours and 19 hours respectively, whereas Java-CPATMINER takes less than 8 hours to mine patterns. However, the Java and Python corpus is diverse (see Table 3). Hence, it is hard to make a fair comparison. Nguyen et al. [58] state that the primary goal of the CPATMINER is to mine the corpus weekly to build a database of patterns. We believe that R-CPATMINER and PYCPATMINER have reasonable execution times for achieving the same goal.

6 IMPLICATIONS

We present actionable, empirically-justified implications for four audiences: (i) researchers, (ii) tool builders and IDE designers, (iii) ML library vendors, and (iv) developers and educators.

6.1 Researchers

R1. Exploit applications of change repetitiveness of Python ML software (RQ1, RQ3). In the past, researchers exploited the repetitiveness of changes in Java systems through: code completion [13, 34, 43, 55, 56], automated program repair [6, 9, 50], API recommendation [32, 55], type migration [40], library migration [3, 18, 24, 39, 84], and automated refactoring [19, 28]. Using our rich and diverse dataset of 28K change patterns instances and our JAVAFYPY, researchers can bring the same benefits to Python ML systems.

R2. Foundations to study Python ML Software Evolution (RQ1, RQ2, RQ3). Despite the widespread use of Python ML systems, its evolution and maintenance tasks are the least automated and the least studied due to the unavailability of tool support to study ML systems [22]. Previously, researchers have built infrastructure to study many aspects of software evolution. For example, *RMINER* [85] and *RefDiff* [77] mine refactorings, *TypeFactMiner* [42] mines type changes, *MigrationMiner* [3] and *APIMigrator* [25] mine API migrations, *CPATMiner* [58] and *CodingTracker* [54] mine fine-grained repetitive code changes in Java. This rich infrastructure allows researchers to study Java software systems, and there are hundreds of published papers that are built upon this infrastructure. Unfortunately, they miss the whole ecosystem of Python ML code-bases. Researchers can use our data and toolset to study whether previously-held beliefs from traditional software are still valid for ML systems or whether we need to design new tools and workflows

(e.g., version-control systems, code smells, technical debt, etc.) that are specific to ML.

R3. Enhance existing research and tools for ML systems (RQ1, RQ3): Our infrastructure can be used to enhance the existing research on ML systems. For example, Tang et al. [82] introduce a taxonomy of refactoring kinds performed in Java ML systems, Humbačová et al. [35], Islam et al. [37, 38] introduce a taxonomy of bugs in ML systems based on manual analysis of *StackOverflow* posts. These studies perform extensive manual analysis to build various taxonomies using a smaller dataset of changes. With our significantly larger dataset (of 28K change patterns), they can significantly extend or further validate their taxonomies. Moreover, researchers showed the potential of leveraging ML techniques for code completion [16, 32]. However, they train the ML models either on a noisy or small dataset which could reduce the accuracy of the recommendations. Researchers can use our tools/dataset to train their models on a large, curated dataset and improve their accuracy.

R4. Build novel applications for ML developers (RQ1). Researchers can also use our tools/dataset to build novel applications. Braiek and Khomh [10] observed ML libraries are core components of ML systems and are frequently evolving, which causes developer frustration [22, 92]. A tutoring system can suggest to developers which constructs to use in order to modernize their ML code. For example, a tutoring system could recommend changing the code to use ML library APIs instead of `for` loops, or removing redundant matrix operations (or other changes from RQ1:Table 1). Moreover, a tutor system can recommend novice programmers how to use advanced language features (see Trend 4 in RQ1:Table 1).

R5. Revisit existing studies and tools for making them refactoring-aware (RQ2). There exist a plethora of research tools for mining software repositories, and hundreds of researchers used these tools to conduct empirical studies. Given that refactorings obfuscate program elements during software evolution, we showed that making a state of the art tool such as CPATMINER refactoring aware increases its overall effectiveness by as much as 15%. We call the research community to adapt similar methods and to revisit previous results obtained with tools that were not refactoring-aware.

6.2 Tool Builders and IDE Designers

T1. New inspirations for tool development (RQ1). To help tool builders invest resources where automation is most needed, in Table 1 we present 22 patterns along with the ML developers request for automation. Moving to `with` statements and using Context managers is the most prevalent change pattern among the analysed patterns (See RQ1, P1-P8). In the survey, 74% of respondents suggest tools that inspect deep learning codebases and recommend using `with` statements to (i) turn on or off gradient calculations (P2), (ii) specify hardware type (P3), (iii) change variable scopes (P4), and (iv) execute dependencies (P5). Respondents further suggested tools to (i) move Context manager invocations to `with` statements, and (ii) detect misuses of Context managers.

6.3 ML Library vendors

L1. Understand ML Library Usage (RQ1). ML library vendors continuously improve libraries and introduce new ML libraries at an unprecedented rate [11, 22]. Library developers deprecate APIs,

introduce more efficient alternative APIs (e.g., Table 1 : P18), and split ML libraries [22]. Our findings, the accompanying dataset [21], and the tools we developed can help ML library vendors to understand what APIs are most commonly used, misused, and underused, and how the developers adapt to new APIs. Thus, they can make informed, empirically-justified decisions to improve features [23, 49].

6.4 Software Developers and Educators

S1. Rich educational resource (RQ1) Developers learn and educators teach new programming constructs through examples. Robillard et al. [71] studied the challenges of learning APIs and concluded that one of the important factors is the lack of usage examples. Using our dataset of 28K code change patterns that we mined in our corpus, developers and educators can learn from real-world code transformations (e.g., transforming to `multi_dot`). We provide 22 empirically justified code change patterns that improve ML code from many aspects, including speed, code quality, and readability. ML developers can absorb these changes to their code and improve the code. We released this through an educational resource [21].

7 THREATS TO VALIDITY

Internal Validity: *Can we trust the results produced by tools?* The findings of our study depend on the accuracy of our tools to mine code change patterns in a refactoring aware manner. We rely on *type inference* for augmenting the AST with rich type information. Since *Type inference* deduces the types of elements by statically analysing the program, it may not accurately detect them (compared to the run time). This can effect the quality of mappings reported by R-CPATMINER and the statements matched by RMINER. To mitigate this threat, we use *PyType* [30], a mature tool developed by Google. Thousands of projects at Google and other places rely upon *PyType* to keep their codes well-typed [30]. Moreover, we validate JAVAFYPY’s effectiveness at transforming a variety of syntactic variations upon 12M AST nodes from 14 popular projects. Our manual validations shows that PYRMINER reports refactorings with high precision (95%). We also manually validated 2,500 most popular patterns produced by R-CPATMINER.

To avoid the *experimental bias*, we followed the best practices [14, 88] for applying thematic analysis by achieving 80% inter coder agreement when labelling the patterns.

External Validity: *Do our results generalize?* We studied 1000 projects from a wide range of application domains, making the study results generalizable to other open-source projects in similar domains. However, a study of proprietary code bases might reveal other trends. Nevertheless, we make our tools available so that others can use them to mine patterns in proprietary codebases. Moreover, R-CPATMINER reports numerous patterns, we manually analyzed a subset of them; a complete investigation is not practical. To mitigate this, we ranked the code change patterns in five dimensions and manually validated the top ones.

Verifiability: *Can others replicate our results?* To ensure replicability, we make the tools and the data publicly available [21].

8 RELATED WORK

8.0.1 Studies on evolution of ML software systems. Researchers have studied repetitive tasks of ML systems from many aspects.

Humbatova et al. [35] and Islam et al. [37, 38] introduce a taxonomy of ML-related bugs and bug fix patterns, Shen et al. [76] study bugs in Deep Learning (DL) compilers, and Yan et al. [90] study numerical bugs in DL systems. Zhang et al. [91] observe DL program failures at Microsoft and highlight the need for DL-specific tools to fix DL bugs. Lwakatare et al. [46], Zhang et al. [92] classify common challenges that ML developers face when maintaining and evolving ML systems. Wan et al. [87] found anti-patterns and misused APIs in ML libraries. Nguyen et al. [57] create a database of DL models extracted from GitHub to utilize as default models and accelerate AutoML. Amershi et al. [4] perform a field study at Microsoft and propose best practices to address common challenges that repeats among ML systems. In contrast, our focus is on understanding fine-grained repeated code changes that the ML developers perform and their motivations.

The closest related work is by Tang et al. [82]. The authors study 327 git patches from 26 *Java* ML systems and present a taxonomy of refactoring kinds that occurred in ML systems. In contrast, we quantitatively and qualitatively studied fine-grained code change patterns in a significantly larger corpus: 1000 *Python*-based ML systems comprising 4,166,520 git patches. Our findings include both refactorings (e.g., P1, P9) and other semantics-modifying code changes (e.g., P2, P6). Moreover, we cross-validated our findings by surveying 97 ML developers. We make our tools available to further enhance the science and tooling for evolving Python ML systems.

8.0.2 Studies on repetitive code changes. Researchers have conducted many studies on repetitive code changes. Nguyen et al. [58, 59] use a graph-based algorithm to mine fine-grained code changes at commit level. They conduct a large-scale study on the repetitiveness of code changes in Java software evolution and show that repetitiveness is common in small granularity (number of lines), and it drops exponentially as the granularity increases. Researchers have also studied repetitiveness from the vantage point of higher-level maintenance and evolution tasks, (i) Dig and Johnson [20], Cossette and Walker [17] study incompatibilities between API versions, (ii) Teyton et al. [84] mine the library migrations trends and observe how frequently, when, and why they are performed, and (iii) Ketkar et al. [40] conduct a large-scale study on type changes in Java systems and reveal that type changes are more frequent than renaming. All of these studies focus on repetitiveness in traditional systems and do not reveal the kinds of repetitive changes the developers perform in ML software. In contrast, we study the practices of fine-grained code changes in ML systems and found four trends of fine-grained changes. 71% of the surveyed developers requested automation support for the identified trends in their IDEs.

Previous researchers looked at Python idioms and how they were used in Python systems. Alexandru et al. [1] present a non-exhaustive list of Python idioms gleaned from a developer survey. Sakulniwat et al. [74] studied the evolution of Python *with* statement over the time. Phan-udom et al. [63] use developer forums to build a database of Python idioms and propose *Teddy*, a system that recommends idioms to developers. In comparison, we provide R-CPATMINER, which can mine repetitive code changes in Python systems in refactoring aware manner. R-CPATMINER can be used to increase related work’s idiom databases.

Researchers reimplemented certain Java AST mining tools for Python from scratch. PYREF [5] is a Python reimplement of Java RMINER [85] that identifies 9 kinds of Python refactoring. In contrast, our PYRMINER does not require reimplementing from scratch, and it finds *all* legitimate Java refactoring in Python. Moreover, we manually validated 18 kinds of refactoring, as shown in Table 2. Golubev et al. [29] developed PYTHONCHANGEMINER, a tool that uses *fgPDG* [58] to mine Python repetitive code changes involving function calls. In contrast, we present JAVAfYPY, a technique for adapting Java AST mining tools to Python without having to rewrite them from the ground up. We adapted CPATMINER [58] and RMINER [85] to Python using JavafYPy, and subsequently created R-CPATMINER, a refactoring-aware code change pattern miner for Python systems.

Researchers proposed several techniques that infer *specific kinds* of code change patterns. For example, REVISAR [72], GETAFIX [6], and DEEPDELTA [50] infer repeated bug fixes and compilation errors from commit histories. LIBSYNC [60], MEDITOR [89], and A3 [44] infer the adaptations required to perform library migration. Ketkar et al. [41] developed TCINFER a technique to accurately infer rewrite rules for type changes from a project’s version history. They also proposed INTELLITC [81], a configurable IDE refactoring plugin that automates these inferred type changes. In contrast, we discover *previously unknown* patterns in ML systems that involve adapting ML libraries (E.g. P9 - `for` loop to `NumPy`).

9 CONCLUSIONS

This paper presents the first and the largest study of fine-grained code change patterns in Python-based ML software systems. To provide unique insights, we use complementary empirical methods: (i) mining 1000 software repositories containing over 58 million LOC, (ii) using thematic analysis to identify the groups and trends, and (iii) surveying 97 ML developers. To conduct this study and advance the science and tooling in Python ML software evolution, we designed a novel technique, JAVAfYPY, to reuse, adapt and improve upon the Java state-of-the-art AST mining tools. We introduce a novel tool R-CPATMINER that performs refactoring-aware change pattern mining in the version history of Python projects. We present 22 code change pattern groups in four trends, where 10 of them are specific to ML. In the developer survey, 71% of the respondents requested these patterns automated in their IDEs. The results and the tools presented in this study have actionable implications for researchers, tool builders, library designers, ML developers, and educators. We hope that this paper and our readily available dataset and tools [21] catalyzes the community to advance the science and tooling for the evolution of Python-based ML systems.

10 ACKNOWLEDGEMENTS

We would like to thank Ellick Chan, Rahul Khanna, Bob Banfield, Julia Romero, Carla Pomian, Dorin Pomian, students from CSCI 7000-08-Spring 2022 at CU Boulder, CUPLV group at CU Boulder, and the anonymous reviewers for their insightful and constructive feedback for improving the paper. This research was partially funded through the NSF grants CCF-1553741, CNS-1941898, and the Industry-University Cooperative Research Center on Pervasive Personalized Intelligence.

REFERENCES

- [1] Carol V. Alexandru, José J. Merchante, Sebastiano Panichella, Sebastian Proksch, Harald C. Gall, and Gregorio Robles. 2018. On the Usage of Pythonic Idioms. In *ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software* (Boston, MA, USA) (Onward! 2018). ACM, New York, NY, USA, 1–11. <https://doi.org/10.1145/3276954.3276960>
- [2] Miltiadis Allamanis and Charles Sutton. 2014. Mining Idioms from Source Code. In *FSE 2014* (Hong Kong, China). Association for Computing Machinery, New York, NY, USA, 472–483. <https://doi.org/10.1145/2635868.2635901>
- [3] Hussein Alrubaye, Mohamed Wiem Mkaouer, and Ali Ouni. 2019. Migration-Miner: An Automated Detection Tool of Third-Party Java Library Migration at the Method Level. In *ICSME 2019*. 414–417. <https://doi.org/10.1109/ICSME.2019.00072>
- [4] Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. 2019. Software Engineering for Machine Learning: A Case Study. In *ICSE* (Montreal, Quebec, Canada) (*ICSE-SEIP '19*). IEEE Press, Piscataway, NJ, USA, 291–300. <https://doi.org/10.1109/ICSE-SEIP.2019.00042>
- [5] Hassan Atwi, Bin Lin, Nikolaos Tsantalis, Yutaro Kashiwa, Yasutaka Kamei, Naoyasu Ubayashi, Gabriele Bavota, and Michele Lanza. 2021. PYREF: Refactoring Detection in Python Projects. In *2021 IEEE 21st International Working Conference on Source Code Analysis and Manipulation (SCAM)*. 136–141. <https://doi.org/10.1109/SCAM52516.2021.00025>
- [6] Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: Learning to Fix Bugs Automatically. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 159 (Oct. 2019), 27 pages. <https://doi.org/10.1145/3360585>
- [7] Earl T. Barr, Yuriy Brun, Premkumar Devanbu, Mark Harman, and Federica Sarro. 2014. The Plastic Surgery Hypothesis. In *FSE 2014* (Hong Kong, China) (*FSE 2014*). Association for Computing Machinery, New York, NY, USA, 306–317. <https://doi.org/10.1145/2635868.2635898>
- [8] Amine Barrak, Ellis E. Eghan, and Bram Adams. 2021. On the Co-evolution of ML Pipelines and Source Code - Empirical Study of DVC Projects. In *SANER 2021*. 422–433. <https://doi.org/10.1109/SANER50967.2021.00046>
- [9] Rohan Bavishi, Hiroaki Yoshida, and Mukul R. Prasad. 2019. Phoenix: Automated Data-Driven Synthesis of Repairs for Static Analysis Violations. In *ESEC/FSE 2019* (Tallinn, Estonia) (*ESEC/FSE 2019*). ACM, New York, NY, USA, 613–624. <https://doi.org/10.1145/3338906.3338952>
- [10] Houssein Ben Braiek and Foutse Khomh. 2020. On testing machine learning programs. *Journal of Systems and Software* 164 (2020), 110542. <https://doi.org/10.1016/j.jss.2020.110542>
- [11] Houssein Ben Braiek, Foutse Khomh, and Bram Adams. 2018. The Open-Closed Principle of Modern Machine Learning Frameworks. In *MSR '18* (Gothenburg, Sweden) (*MSR '18*). Association for Computing Machinery, New York, NY, USA, 353–363. <https://doi.org/10.1145/3196398.3196445>
- [12] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology* 3, 2 (2006), 77–101.
- [13] Marcel Bruch, Martin Monperrus, and Mira Mezini. 2009. Learning from Examples to Improve Code Completion Systems. In *ESEC/FSE '09* (Amsterdam, The Netherlands) (*ESEC/FSE '09*). Association for Computing Machinery, New York, NY, USA, 213–222. <https://doi.org/10.1145/1595696.1595728>
- [14] John L Campbell, Charles Quincy, Jordan Osserman, and Ove K Pedersen. 2013. Coding in-depth semistructured interviews: Problems of unitization and inter-coder reliability and agreement. *Sociological Methods & Research* 42, 3 (2013), 294–320. <https://doi.org/10.1177/0049124113500475>
- [15] Zhenpeng Chen, Yanbin Cao, Yuanqiang Liu, Haoyu Wang, Tao Xie, and Xuanzhe Liu. 2020. A Comprehensive Study on Challenges in Deploying Deep Learning Based Software. In *FSE* (Virtual Event, USA) (*ESEC/FSE 2020*). Association for Computing Machinery, New York, NY, USA, 750–762. <https://doi.org/10.1145/3368089.3409759>
- [16] Zimin Chen, Steve James Kommrusch, Michele Tufano, Louis-Noël Pouchet, Denys Poshyvanyk, and Martin Monperrus. 2019. SEQUENCER: Sequence-to-Sequence Learning for End-to-End Program Repair. *TSE 2019* (2019), 1–1. <https://doi.org/10.1109/TSE.2019.2940179>
- [17] Bradley E. Cossette and Robert J. Walker. 2012. Seeking the Ground Truth: A Retroactive Study on the Evolution and Migration of Software Libraries (*FSE '12*). ACM, New York, NY, USA, Article 55, 11 pages. <https://doi.org/10.1145/2393596.2393661>
- [18] Barthélemy Dagenais and Martin P. Robillard. 2011. Recommending Adaptive Changes for Framework Evolution. *ACM Trans. Softw. Eng. Methodol.* 20, 4, Article 19 (Sept. 2011), 35 pages. <https://doi.org/10.1145/2000799.2000805>
- [19] Danny Dig, Can Comertoglu, Darko Marinov, and Ralph Johnson. 2006. Automated Detection of Refactorings in Evolving Components. In *ECOOP'06* (Nantes, France) (*ECOOP'06*). Springer-Verlag, Berlin, Heidelberg, 404–428. https://doi.org/10.1007/11785477_24
- [20] Danny Dig and Ralph Johnson. 2006. How Do APIs Evolve? A Story of Refactoring: Research Articles. *J. Softw. Maint. Evol.* 18, 2 (March 2006), 83–107.
- [21] Malinda Dilhara. 2022. *Discovering Repetitive Code Changes in Python-based ML Systems*. <https://mlcodepatterns.github.io> Accessed: 2022-02-07.
- [22] 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 (July 2021), 42 pages. <https://doi.org/10.1145/3453478>
- [23] Robert Dyer, Hridesh Rajan, Hoan Anh Nguyen, and Tien N. Nguyen. 2014. Mining Billions of AST Nodes to Study Actual and Potential Usage of Java Language Features. In *ICSE* (Hyderabad, India) (*ICSE 2014*). Association for Computing Machinery, New York, NY, USA, 779–790. <https://doi.org/10.1145/2568225.2568295>
- [24] Mattia Fazzini, Qi Xin, and Alessandro Orso. 2019. Automated API-Usage Update for Android Apps. In *ISSTA 2019* (Beijing, China) (*ISSTA 2019*). Association for Computing Machinery, New York, NY, USA, 204–215. <https://doi.org/10.1145/3293882.3330571>
- [25] Mattia Fazzini, Qi Xin, and Alessandro Orso. 2020. APIMigrator: An API-Usage Migration Tool for Android Apps. In *MOBILESoft '20* (Seoul, Republic of Korea) (*MOBILESoft '20*). Association for Computing Machinery, New York, NY, USA, 77–80. <https://doi.org/10.1145/3387905.3388608>
- [26] Michael Feathers. 2004. *Working Effectively with Legacy Code: WORK EFFECT LEG CODE_p1*. Prentice Hall Professional.
- [27] Eclipse foundation. 2021. *JDT Core Component*. Eclipse. https://www.eclipse.org/jdt/core/#JDT_CORE Accessed: 2021-03-31.
- [28] Lyle Franklin, Alex Gyori, Jan Lahoda, and Danny Dig. 2013. LAMBDAFICATOR: From Imperative to Functional Programming through Automated Refactoring. In *ICSE* (San Francisco, CA, USA) (*ICSE '13*). IEEE Press, 1287–1290. <https://doi.org/10.1109/ICSE.2013.6606699>
- [29] Yaroslav Golubev, Jiawei Li, Viacheslav Blushev, Timofey Bryksin, and Iftekhar Ahmed. 2021. Changes from the trenches: Should we automate them? *arXiv preprint arXiv:2105.10157* (2021).
- [30] Google. 2021. *PyType*. <https://github.com/google/pytype> Accessed: 2021-03-31.
- [31] Fabrice Harel-Canada, Lingxiao Wang, Muhammad Ali Gulzar, Quanquan Gu, and Miryung Kim. 2020. Is Neuron Coverage a Meaningful Measure for Testing Deep Neural Networks?. In *FSE* (Virtual Event, USA) (*ESEC/FSE 2020*). ACM, New York, NY, USA, 851–862. <https://doi.org/10.1145/3368089.3409754>
- [32] Xincheng He, Lei Xu, Xiangyu Zhang, Rui Hao, Yang Feng, and Baowen Xu. 2021. PyART: Python API Recommendation in Real-Time. In *ICSE 2021*. 1634–1645. <https://doi.org/10.1109/ICSE43902.2021.00145>
- [33] Abram Hindle, Earl T. Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. 2016. On the Naturalness of Software. *Commun. ACM* 59, 5 (April 2016), 122–131. <https://doi.org/10.1145/2902362>
- [34] Reid Holmes, Robert J. Walker, and Gail C. Murphy. 2006. Approximate Structural Context Matching: An Approach to Recommend Relevant Examples. *IEEE Trans. Softw. Eng.* 32, 12 (Dec. 2006), 952–970. <https://doi.org/10.1109/TSE.2006.117>
- [35] Nargiz Humbatova, Gunel Jahangirova, Gabriele Bavota, Vincenzo Riccio, Andrea Stocco, and Paolo Tonella. 2020. Taxonomy of Real Faults in Deep Learning Systems. In *ICSE '20* (Seoul, South Korea) (*ICSE '20*). ACM, New York, NY, USA, 1110–1121. <https://doi.org/10.1145/3377811.3380395>
- [36] Nargiz Humbatova, Gunel Jahangirova, and Paolo Tonella. 2021. DeepCrime: Mutation Testing of Deep Learning Systems Based on Real Faults. In *ISSTA-2021* (Virtual, Denmark) (*ISSTA 2021*). ACM, New York, NY, USA, 67–78. <https://doi.org/10.1145/3460319.3464825>
- [37] Md Johirul Islam, Giang Nguyen, Rangeet Pan, and Hridesh Rajan. 2019. A Comprehensive Study on Deep Learning Bug Characteristics. In *ESEC/FSE 2019* (Tallinn, Estonia) (*ESEC/FSE 2019*). Association for Computing Machinery, New York, NY, USA, 510–520. <https://doi.org/10.1145/3338906.3338955>
- [38] Md Johirul Islam, Rangeet Pan, Giang Nguyen, and Hridesh Rajan. 2020. Repairing Deep Neural Networks: Fix Patterns and Challenges. In *ICSE '20* (Seoul, Republic of Korea) (*ICSE '20*). ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/1122445.1122456>
- [39] Suhas Kabinna, Cor-Paul Bezemer, Weiyei Shang, and Ahmed E. Hassan. 2016. Logging Library Migrations: A Case Study for the Apache Software Foundation Projects. In *MSR '16* (Austin, Texas) (*MSR '16*). ACM, New York, NY, USA, 154–164. <https://doi.org/10.1145/2901739.2901769>
- [40] Ameya Ketkar, Ali Mesbah, Davood Mazinanian, Danny Dig, and Edward Af-tandilian. 2019. Type Migration in Ultra-Large-Scale Codebases. In *ICSE '19* (Montreal, Quebec, Canada) (*ICSE '19*). IEEE Press, 1142–1153. <https://doi.org/10.1109/ICSE.2019.00117>
- [41] Ameya Ketkar, Oleg Smirnov, Nikolaos Tsantalis, Danny Dig, and Timofey Bryksin. 2022. Inferring and Applying Type Changes. In *44th International Conference on Software Engineering (ICSE '22)* (Pittsburgh, United States) (*ICSE '22*). ACM. <https://doi.org/10.1145/3510003.3510115>
- [42] Ameya Ketkar, Nikolaos Tsantalis, and Danny Dig. 2020. Understanding Type Changes in Java. In *FSE* (Virtual Event, USA) (*ESEC/FSE 2020*). Association for Computing Machinery, New York, NY, USA, 629–641. <https://doi.org/10.1145/3368089.3409725>
- [43] Dongsun Kim, Jaechang Nam, Jaewoo Song, and Sunghun Kim. 2013. Automatic Patch Generation Learned from Human-Written Patches. In *ICSE* (San Francisco, CA, USA) (*ICSE '13*). IEEE Press, 802–811.

- [44] M. Lamothe, W. Shang, and T. Chen. 2020. A3: Assisting Android API Migrations Using Code Examples. *TSE* 2020 01 (apr 2020), 1–1. <https://doi.org/10.1109/TSE.2020.2988396>
- [45] Jiakun Liu, Qiao Huang, Xin Xia, Emad Shihab, David Lo, and Shanping Li. 2020. Is Using Deep Learning Frameworks Free? Characterizing Technical Debt in Deep Learning Frameworks. In *ICSE* (Seoul, South Korea) (*ICSE-SEIS '20*). ACM, New York, NY, USA, 1–10. <https://doi.org/10.1145/3377815.3381377>
- [46] Lucy Ellen Lwakatare, Aiswarya Raj, Jan Bosch, Helena Holmström Olsson, and Ivica Crnkovic. 2019. A Taxonomy of Software Engineering Challenges for Machine Learning Systems: An Empirical Investigation. In *Agile Processes in Software Engineering and Extreme Programming*, Philippe Kruchten, Steven Fraser, and François Coallier (Eds.). Springer International Publishing, Cham, 227–243.
- [47] Robert C Martin. 2009. *Clean code: a handbook of agile software craftsmanship*. Pearson Education.
- [48] Matplotlib. 2021. *Matplotlib*. <https://matplotlib.org> Accessed: 2021-05-05.
- [49] Davood Mazinanian, Ameya Ketkar, Nikolaos Tsantalis, and Danny Dig. 2017. Understanding the Use of Lambda Expressions in Java. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 85 (Oct. 2017), 31 pages. <https://doi.org/10.1145/3133909>
- [50] Ali Mesbah, Andrew Rice, Emily Johnston, Nick Glorioso, and Edward Aftandilian. 2019. DeepDelta: Learning to Repair Compilation Errors. In *ESEC/FSE*. 925–936. <https://doi.org/10.1145/3338906.3340455>
- [51] Emerson Murphy-Hill, Chris Parnin, and Andrew P. Black. 2009. How We Refactor, and How We Know It. In *ICSE '09* (*ICSE '09*). Association for Computing Machinery, New York, NY, USA, 287–297. <https://doi.org/10.1109/ICSE.2009.5070529>
- [52] Stas Negara, Nicholas Chen, Mohsen Vakilian, Ralph E. Johnson, and Danny Dig. 2013. A Comparative Study of Manual and Automated Refactorings. In *ECOOP'13* (Montpellier, France) (*ECOOP'13*). Springer-Verlag, Berlin, Heidelberg, 552–576. https://doi.org/10.1007/978-3-642-39038-8_23
- [53] Stas Negara, Mihai Codoban, Danny Dig, and Ralph E. Johnson. 2014. Mining Fine-Grained Code Changes to Detect Unknown Change Patterns. In *ICSE* (Hyderabad, India) (*ICSE 2014*). Association for Computing Machinery, New York, NY, USA, 803–813. <https://doi.org/10.1145/2568225.2568317>
- [54] Stas Negara, Mohsen Vakilian, Nicholas Chen, Ralph E. Johnson, and Danny Dig. 2012. Is It Dangerous to Use Version Control Histories to Study Source Code Evolution?. In *ECOOP'12* (Beijing, China) (*ECOOP'12*). Springer-Verlag, Berlin, Heidelberg, 79–103. https://doi.org/10.1007/978-3-642-31057-7_5
- [55] Anh Tuan Nguyen, Michael Hilton, Mihai Codoban, Hoan Anh Nguyen, Lily Mast, Eli Rademacher, Tien N. Nguyen, and Danny Dig. 2016. *API Code Recommendation Using Statistical Learning from Fine-Grained Changes*. ACM, New York, NY, USA, 511–522. <https://doi.org/10.1145/2950290.2950333>
- [56] Anh Tuan Nguyen, Tung Thanh Nguyen, Hoan Anh Nguyen, Ahmed Tamrawi, Hung Viet Nguyen, Jafar Al-Kofahi, and Tien N. Nguyen. 2012. Graph-Based Pattern-Oriented, Context-Sensitive Source Code Completion. In *ICSE* (Zurich, Switzerland) (*ICSE '12*). IEEE Press, 69–79.
- [57] Giang Nguyen, Johir Islam, Rangeet Pan, and Hridesh Rajan. 2022. Manas: Mining Software Repositories to Assist AutoML. In *ICSE'22: The 44th International Conference on Software Engineering* (Pittsburgh, PA, USA). <https://doi.org/10.1145/3510003.3510052> To appear.
- [58] H. Nguyen, T. N. Nguyen, D. Dig, S. Nguyen, H. Tran, and M. Hilton. 2019. Graph-Based Mining of In-the-Wild, Fine-Grained, Semantic Code Change Patterns. In *ICSE 2019*. IEEE Computer Society, Los Alamitos, CA, USA, 819–830. <https://doi.org/10.1109/ICSE.2019.00089>
- [59] Hoan Anh Nguyen, Anh Tuan Nguyen, Tung Thanh Nguyen, Tien N. Nguyen, and Hridesh Rajan. 2013. A Study of Repetitiveness of Code Changes in Software Evolution. In *ASE '13* (Silicon Valley, CA, USA) (*ASE'13*). IEEE Press, 180–190. <https://doi.org/10.1109/ASE.2013.6693078>
- [60] Hoan Anh Nguyen, Tung Thanh Nguyen, Gary Wilson, Anh Tuan Nguyen, Miryung Kim, and Tien N. Nguyen. 2010. A Graph-Based Approach to API Usage Adaptation. *SIGPLAN Not.* 45, 10 (Oct. 2010), 302–321. <https://doi.org/10.1145/1932682.1869486>
- [61] Numpy. 2021. *NumPy*. <https://numpy.org> Accessed: 2021-05-05.
- [62] Oracle. 2021. *Java SE specification*. Oracle. <https://docs.oracle.com/javase/specs/> Accessed: 2021-03-31.
- [63] Purit Phan-udom, Naruedon Wattanakul, Tattiya Sakulniwat, Chaiyong Ragkhitwetsagul, Thanwadee Sunetnanta, Morakot Choetkiertikul, and Raula Gaikovina Kula. 2020. Teddy: Automatic Recommendation of Pythonic Idiom Usage For Pull-Based Software Projects. In *ICSME-2020*. 806–809. <https://doi.org/10.1109/ICSME46990.2020.00098>
- [64] Python. 2021. *Context Manager*. <https://docs.python.org/3/reference/datamodel.html#context-managers> Accessed: 2021-03-31.
- [65] Python. 2021. *Functions*. Python. <https://docs.python.org/3/library/functions.html> Accessed: 2021-03-31.
- [66] Python. 2021. *list-comprehensions*. <https://docs.python.org/3/tutorial/datastructures.html#list-comprehensions> Accessed: 2021-05-05.
- [67] Python. 2021. *Python AST*. Python. <https://docs.python.org/3/library/ast.html> Accessed: 2021-03-31.
- [68] Python. 2021. *With Statement*. Python. https://docs.python.org/3/reference/compound_stmts.html#the-with-statement Accessed: 2021-03-31.
- [69] PythonTypeInformation. 2021. *PythonTypeInformation*. GitHub. <https://github.com/mlcodepatterns/PythonTypeInformation> Accessed: 2021-05-05.
- [70] Sebastian Raschka and Vahid Mirjalili. 2017. *Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow*.
- [71] Martin P. Robillard and Robert Deline. 2011. A Field Study of API Learning Obstacles. *Empirical Softw. Engg.* 16, 6 (Dec. 2011), 703–732. <https://doi.org/10.1007/s10664-010-9150-8>
- [72] Reudimam Rolim, Gustavo Soares, Rohit Gheyi, and Loris D'Antoni. 2018. Learning Quick Fixes from Code Repositories. (2018). <http://arxiv.org/abs/1803.03806>
- [73] RuleOfThree. 2021. *RuleOfThree*. wicik2. <http://wiki.c2.com/?RuleOfThree> Accessed: 2021-05-05.
- [74] Tattiya Sakulniwat, Raula Gaikovina Kula, Chaiyong Ragkhitwetsagul, Morakot Choetkiertikul, Thanwadee Sunetnanta, Dong Wang, Takashi Ishio, and Kenichi Matsumoto. 2019. Visualizing the Usage of Pythonic Idioms over Time: A Case Study of the with open Idiom. In *IWESEP-2019*. IEEE, 43–435.
- [75] D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. 2015. Hidden Technical Debt in Machine Learning Systems (*NIPS'15*). MIT Press, Cambridge, MA, USA, 2503–2511.
- [76] Qingchao Shen, Haoyang Ma, Junjie Chen, Yongqiang Tian, Shing-Chi Cheung, and Xiang Chen. 2021. A Comprehensive Study of Deep Learning Compiler Bugs. In *FSE* (Athens, Greece) (*ESEC/FSE 2021*). Association for Computing Machinery, New York, NY, USA, 968–980. <https://doi.org/10.1145/3468264.3468591>
- [77] D. Silva, J. Silva, G. De Souza Santos, R. Terra, and M. O. Valente. 5555. RefDiff 2.0: A Multi-language Refactoring Detection Tool. *TSE* 2020 01 (jan 5555), 1–1. <https://doi.org/10.1109/TSE.2020.2968072>
- [78] Danilo Silva, Nikolaos Tsantalis, and Marco Tulio Valente. 2016. Why We Refactor? Confessions of GitHub Contributors. In *FSE* (Seattle, WA, USA) (*FSE 2016*). Association for Computing Machinery, New York, NY, USA, 858–870. <https://doi.org/10.1145/2950290.2950305>
- [79] Janice Singer, Susan E Sim, and Timothy C Lethbridge. 2008. Software engineering data collection for field studies. In *Guide to Advanced Empirical Software Engineering*. Springer, 9–34.
- [80] Brett Slatkin. 2019. *Effective python: 90 specific ways to write better python*. Addison-Wesley Professional.
- [81] Oleg Smirnov, Ameya Ketkar, Timofey Bryksin, Nikolaos Tsantalis, and Danny Dig. 2022. IntelliTC: Automating Type Changes in IntelliJ IDEA. In *44th International Conference on Software Engineering Companion* (*ICSE '22 Companion*) (Pittsburgh, United States) (*ICSE '22 Companion*). ACM/IEEE. <https://doi.org/10.1145/3510454.3516851>
- [82] Y. Tang, R. Khatchadourian, M. Bagherzadeh, R. Singh, A. Stewart, and A. Raja. 2021. An Empirical Study of Refactorings and Technical Debt in Machine Learning Systems. In *ICSE 2021*. IEEE Computer Society, Los Alamitos, CA, USA, 238–250. <https://doi.org/10.1109/ICSE43902.2021.00033>
- [83] TensorFlow. 2022. *TensorFlow Core v2.7.0*. https://www.tensorflow.org/api_docs/python/tf/io/gfile/GFile Accessed: 2022-01-05.
- [84] Cédric Teyton, Jean-Rémy Fallier, Marc Palyart, and Xavier Blanc. 2014. A Study of Library Migrations in Java. *J. Softw. Evol. Process* 26, 11 (Nov. 2014), 1030–1052. <https://doi.org/10.1002/smr.1660>
- [85] N. Tsantalis, A. Ketkar, and D. Dig. 5555. RefactoringMiner 2.0. *TSE* 2020 01 (jul 5555), 1–1. <https://doi.org/10.1109/TSE.2020.3007722>
- [86] TypeShed. 2021. *TypeShed*. <https://github.com/python/typeshed> Accessed: 2021-09-03.
- [87] C. Wan, S. Liu, H. Hoffmann, M. Maire, and S. Lu. 2021. Are Machine Learning Cloud APIs Used Correctly?. In *ICSE 2021*. IEEE Computer Society, Los Alamitos, CA, USA, 125–137. <https://doi.org/10.1109/ICSE43902.2021.00024>
- [88] David Wicks. 2017. The coding manual for qualitative researchers. *Qualitative research in organizations and management: an international journal* (2017). <https://doi.org/10.1108/QROM-08-2016-1408>
- [89] Shengzhe Xu, Ziqi Dong, and Na Meng. 2019. Meditor: Inference and Application of API Migration Edits. In *Proceedings of the 27th International Conference on Program Comprehension* (Montreal, Quebec, Canada) (*ICPC '19*). IEEE Press, Piscataway, NJ, USA, 335–346. <https://doi.org/10.1109/ICPC.2019.00052>
- [90] Ming Yan, Junjie Chen, Xiangyu Zhang, Lin Tan, Gan Wang, and Zan Wang. 2021. Exposing Numerical Bugs in Deep Learning via Gradient Back-Propagation. In *FSE* (Athens, Greece) (*ESEC/FSE 2021*). ACM, New York, NY, USA, 627–638.
- [91] 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 *ICSE* (Seoul, South Korea) (*ICSE '20*). ACM, New York, NY, USA, 1159–1170. <https://doi.org/10.1145/3377811.3380362>
- [92] T. Zhang, C. Gao, L. Ma, M. Lyu, and M. Kim. 2019. An Empirical Study of Common Challenges in Developing Deep Learning Applications. In *ISSRE*. IEEE Computer Society, Los Alamitos, CA, USA, 104–115. <https://doi.org/10.1109/ISSRE.2019.00020>