

1 Error Location in Python: Where the 2 Mutants Hide

3 Joshua Charles Campbell¹, Abram Hindle¹, and José Nelson Amaral¹

4 ¹Department of Computing Science, University of Alberta, Edmonton, Canada,

5 ABSTRACT

Dynamic scripting programming languages present a unique challenge to software engineering tools that depend on static analysis. Dynamic languages do not benefit from the full lexical and syntax analysis provided by compilers and static analysis tools. Prior work exploited a statically typed language (Java) and a simple n -gram language model to find syntax-error locations in programs. This work investigates whether n -gram-based error location on source code written in a dynamic language is effective without static analysis or compilation. UnnaturalCode.py is a syntax-error locator developed for the Python programming language. The UnnaturalCode.py approach is effective on Python code, but faces significantly more challenges than its Java counterpart did. UnnaturalCode.py generalizes the success of previous statically-typed approaches to a dynamically-typed language.

7 Keywords: Software Engineering, Fault Location, Syntax, Dynamic Languages, Python, Language Modelling, Tool

8 1 INTRODUCTION

9 This paper seeks to help Python programmers find coding mistakes by creating an error
10 location tool, UnnaturalCode.py, for use with Python.

11 Campbell *et al.* [5] describe a prototype system, UnnaturalCode, for locating syntax
12 errors in programs written in the Java programming language. UnnaturalCode produced
13 additional diagnostic information for Java source files that failed to compile. Unnatu-
14 ralCode used an n -gram language model to identify snippets of source code that it finds
15 unlikely to be correct. UnnaturalCode was shown to be effective by locating some
16 errors that the Java compiler was not able to locate. Syntax errors and syntax-error
17 reporting are important because researchers have found that syntax errors are serious
18 roadblocks for introductory programming students. Garner *et al.* [9], corroborated by
19 numerous others [17, 18, 16, 38] show that “students very persistently keep seeking
20 assistance for problems with basic syntactic details.”

21 In this paper, a new UnnaturalCode.py system is applied to a new programming
22 language. In the interests of diversity and generalizability, another popular program-
23 ming language, which is very different from Java, was chosen: Python. Besides being
24 a dynamically typed scripting language, Python also offers additional challenges for
25 the localization of errors, and thus is a good candidate to evaluate UnnaturalCode in
26 a different point in the space of programming languages. For instance, Python is an
27 interactive scripting language, while Java is not. Python has different syntactical style

28 from Java: it uses white-space to define the end of statements and the extent of code
29 blocks.

30 Python has taken hold as an instructional language [22, 3]. It is popular to teach
31 Python and Python is now quite popular among scientists [26, 30]. Many new program-
32 mers are being introduced to programming with Python because Python lacks excessive
33 syntax and forces structured and modular programming via forced white-space block
34 structures [37].

35 However, the approach used in Campbell *et al.* [5] does not work with Python
36 because of the differences between Java and Python. A new approach was developed,
37 to generalize the previous approach.

38 Differences between languages have implications for UnnaturalCode, the practice
39 of software development, and the use of Python as a language for experimentation.
40 Experiments and practices that involve generating Python code that may or may not
41 execute are the focus of this paper. These type of experiments are heavily impacted by
42 Python's features, especially its lack of compilation and static analysis.

43 The contributions of this paper are: 1) a corpus-based tool, UnnaturalCode.py, that
44 can help users locate the causes of errors in Python software; 2) a novel methodology
45 and tool for mutation testing of Python tools without restricting the forms those mu-
46 tations can take; 3) an analysis of the differences in approach required for dynamic
47 scripting languages as compared to approaches that relies on statically typed and com-
48 piled languages.

49 **2 BACKGROUND**

50 UnnaturalCode was introduced in Campbell *et al.* [5]. It is a prototype system that
51 effectively finds syntax errors by training itself on working Java source code and then
52 trying to find phrases of source code that it has not seen before. This paper differs
53 from the prior work by introducing a new system with some of the same core ideas.
54 UnnaturalCode.py includes a mutation-testing framework for Python implementations,
55 software engineering tools, and test suites.

56 The UnnaturalCode prototype for Java depended on language features such as declar-
57 ative scope, static typing, context-free syntax and pre-compiled libraries. Additionally,
58 source-code compliance with the rules and requirements of these language features is
59 testable with a single tool: the Java compiler. It is also reasonable to expect the Java
60 compiler to halt (not halting would be a bug). In this paper, however, arbitrary Python
61 programs are used instead of a compiler, and cannot be reasonably assumed to behave
62 in any particular way. This work discusses many of the issues faced with error location
63 in a dynamic language, as well as mutation testing in a dynamic scripting language.

64 **2.1 n -Grams in Software Engineering**

65 Hindle *et al.* [12] proposed the application of n -gram language models to source code.
66 These models were classically applied to natural language text, but Hindle *et al.* [12]
67 showed that software had properties (low entropy) that made it compatible with these
68 kinds of models.

69 N -grams are a phrase formed by at most n words or tokens. An n -gram language
70 model is simply a collection of counts that represent the number of times a phrase
71 appears in a corpus. Empirically the probability of an n -gram is simply the frequency
72 of occurrence of the phrase in the original corpus. For example, if the for loop “for
73 `x in l :`” occurred 10 times in a corpus consisting of 1000 4-grams, its probability
74 would be 0.01. One could consider the probability of a token in its surrounding context.
75 For example, the token “in” might have a probability of 1 given the context of “for
76 `x ---- l :`” because “in” may be the only token that had been observed in that
77 context.

78 Performance of n -gram models tends to increase as n increases, but so does the
79 memory required to store the model. Larger values of n cause sparsity — not all n -
80 grams of size n will be observed if n is large. Empirically, unless the corpus is large
81 enough, there will be unobserved and legitimate n -grams that are not in the corpus. If
82 missing n -grams have a probability of zero, the model cannot estimate any probabilities
83 on the unseen text. This short-coming is addressed by smoothing. Smoothing estimates
84 the probability of unseen n -grams from the probabilities of the largest m -grams (where
85 $m < n$) that are part of the unseen n -gram and that exist in the corpus. For example, if
86 the corpus does not contain the phrase “for `x in l`” but it does contain the phrases
87 “for `x in`” and “`l`” it would estimate the probability of “for `x in l`” using a
88 function of the two probabilities it does know. Specifically, UnnaturalCode.py uses
89 Modified Kneser-Ney smoothing [15].

90 UnnaturalCode.py works by finding n -grams that are improbable – they have high
91 entropy. Thus UnnaturalCode.py exploits the fact that syntactically invalid source code
92 often has a higher entropy. The entropy of source code is computed with respect to a
93 corpus of code that is known to be valid. Unnatural Code looks for source code that is
94 *unnatural* to a n -gram model trained on compilable source code.

95 **2.2 Programming Errors**

96 There is much empirical research on the effect of syntax errors on novices and experi-
97 enced programmers [18, 17, 16, 38, 9, 25]. The consensus is that programming errors
98 consume time, occur frequently for new programmers, and persistently for experienced
99 programmers.

100 Kummerfeld *et al.* [23] studied the effect of syntax errors on programming experi-
101 ence and program comprehension. They found that inexperienced students made rapid
102 and erratic modifications to their source code, in the hurried hope of achieving a com-
103 piling program. These furious modifications do not exhibit much strategy or reasoning,
104 but represent a brute-force approach to programming. When studying experienced pro-
105 grammers, they found that after the failure of a set of syntax-error solving strategies,
106 the experienced programmers would revert to this brute-force random edit strategy as
107 well.

108 Student performance (grades) has been found to negatively correlate with the fre-
109 quency of syntax errors. Tabanao *et al.* [38, 39] studied more than 120 students and
110 found that errors and compilation frequency were negatively correlated with midterm
111 test performance.

112 Jadud *et al.* [18, 17] studied first-year students' programming errors by taking snap-
113 shots of code before compilation. Each failure was manually labelled. Of relevant
114 failures, 5% took students longer than 10 minutes to resolve while 60% only took 20
115 seconds. Thus, addressing and fixing syntax errors is time consuming for inexperienced
116 programmers.

117 **2.3 Technical approaches to Syntax Errors**

118 Two main methods used to identify and report coding errors are *parser-based* and *type-*
119 *based*. *Parser-based* methods augment the parser to enable better reporting or to skip-
120 ping over ambiguous sections in hopes of providing more feedback. Many of these
121 techniques seek to create more error messages that might be more helpful. Burke's
122 parse action deferral [4] backs the parser down the parse stack when an error is encoun-
123 tered and then discards problematic tokens. This approach allows the rest of the file to
124 be checked as well, assuming that the state of the parser is recoverable. In a similar vein,
125 heuristics and cost-based approaches, such as production rule prioritization resumption
126 were combined by Graham *et al.* [10] in order to provide better parse error messages.
127 Parr *et al.* [28] discuss the LL(*) parsing strategy used in parser-generators such as
128 ANTLR. LL(*) parsers dynamically attempt to increase the look-ahead at parse time
129 up to a constant maximum look-head, finally failing over to a backtracking algorithm.
130 This method gains contextual error reporting from a top-down perspective enabling
131 source-level debugging while still performing well enough to be a usable and popular
132 parser.

133 Other parse-based methods avoid reporting spurious parse errors. They often work
134 by attempting to repair the source code or by resuming the parse after an error. Kim
135 *et al.* [21] and Corchuelo *et al.* [7] apply search methods to find repairs to source code
136 to enable a parser to continue, often without user intervention. These are merely two
137 examples of such techniques: research on parser error recovery and diagnostics spans
138 many decades and is far too voluminous to list here.

139 *Type-based* static analysis is leveraged as an alternative method of parse/syntax
140 error fixing. Python employs a liberal type system and parser which can be viewed
141 as a barrier to *Type-based* static analysis. Cannon concluded that one barrier to type
142 inference and type-based static analysis was that "Object-oriented programming also
143 affects the effectiveness of type inference." [6]. Many attempts at type inference, such
144 as Starkiller by Salib [35] and Psycho by Rigo *et al.* [34, 32] and Python ignore part of
145 Python's semantics or just handle subsets of Python's language. PyPy [33], is the most
146 modern effort at static analysis of Python code by enabling just in time (JIT) compi-
147 lation of Python code. RPython by Ancona *et al.* [1] attempted to limit the semantics
148 of Python to work "naturally" within statically typed domains such the JVM or CLI.
149 Thus there has been some work on Python, static analysis and *type-based* analysis, but
150 much of it either chooses to work on safe subsets of the languages or restrict language
151 semantics.

152 *Type-based* analysis focuses on reconciling the types of the identifier and functions
153 called, rather than focusing on grammatical correctness. Heeren [11] leverages types
154 to implement a constraint-based framework within the compiler.

155 Some *type-based* techniques leverage corpora either to aid their type-based reason-
156 ing with heuristics or to improve search performance. Hristova *et al.* [14] leverage pre-
157 defined heuristic rules to address common mistakes observed in code. Lerner *et al.* [24]
158 use a corpus of compilable software and its types to improve type error messages for
159 statically typed languages.

160 One *mutation-based* technique was described by Weimer *et al.* [40]. They use
161 Genetic Algorithms to mutate parse trees to fix defects. In comparison to previous
162 mutation-based approaches, the mutations presented in this paper are not guaranteed to
163 produce either parsable text or a valid parse tree.

164 3 IMPLEMENTATION

165 UnnaturalCode.py is intended to augment the Python run-time environment's error mes-
166 sages with additional information about the probable location of a coding error. It is
167 comprised of two major parts. The first is a piece of Python software, and the second
168 is a modified version of the MIT Language Model (MITLM). MITLM is a software
169 implementation of the n -gram language model written in C++.

170 UnnaturalCode.py wraps around the system's Python interpreter in order to add the
171 desired functionality: UnnaturalCode.py is invoked to run Python programs instead
172 of the system's Python interpreter. First, UnnaturalCode.py runs the desired Python
173 program in exactly the same way as the system's Python interpreter. Then, Unnatural-
174 Code.py checks the result of the execution.

175 There are two possible outcomes:

- 176 • In the case that the Python program exited successfully, or more preferably, the
177 Python program's test suite *passed*, UnnaturalCode.py can add it to its corpus of
178 known-good Python code.
- 179 • In the case that the Python program exits with an *error*, UnnaturalCode.py at-
180 tempts to locate the source of that error and to give the user a suggestion of where
181 to look for coding mistakes along with the standard Python error information.

182 UnnaturalCode.py could be integrated into Beck's test-driven-development [2] pro-
183 cess in two places: as a method to help locate test-failing code in the case of an error
184 inducing a test failure; after tests pass UnnaturalCode.py could be updated with the
185 new working code so that it becomes more familiar with the successful test-compliant
186 system.

187 UnnaturalCode.py does not interfere with the usual process of executing Python
188 software. UnnaturalCode.py's goal is only to augment the usual Python diagnostics
189 and error output with its own suggestions. In order to achieve this goal it performs the
190 following functions: first, it lexically analyzes the Python program with its own custom
191 lexical analyzer. This lexical analyzer is based on Python's standard lexical analyzer,
192 but modified to continue in the case of an error instead of stopping.

193 Second, it breaks the Python file that caused the error into sequences of 20 contigu-
194 ous tokens using a sliding window. A 20-token sliding window is used to ensure that
195 MITLM has sufficient context to employ smoothing. Then, it sends each sliding win-
196 dow to MITLM. MITLM can either be running locally or on a server. MITLM returns,

197 for each sequence, a single value representing the cross entropy (or, log probability) of
198 that sequence versus the corpus of known-good Python code.

199 Finally, `UnnaturalCode.py` reports to the user the sequence that has the highest
200 cross-entropy, along with its line number, file name, and the usual Python diagnos-
201 tics and error messages. The entire `UnnaturalCode.py` process usually takes less than
202 1 second and 300-500MB of RAM on a modern 64-bit PC for a single large Python
203 project.

204 MITLM is configured to use a 10-gram model, and the modified Kneser-Ney smooth-
205 ing, when estimating entropy values. This configuration allows `UnnaturalCode.py` to
206 receive reasonable entropy estimates even for token sequences that have never been
207 seen before.

208 `UnnaturalCode.py` also includes a program mutation tool, which it uses to test itself.
209 This tool can be used to test test suites, error handling, and Python implementations.
210 The mutation tool includes rules for 14 different types of mutations, all of which are
211 designed to be extremely general.

212 11 of those 14 types of random mutations were studied in this paper: token deletion,
213 token insertion, token replacement, digit deletion, digit insertion, letter deletion, letter
214 insertion, symbol deletion, symbol insertion, line dedenting, line indenting.

215 These 11 types of random-edit-mutations are intended to simulate mistakes that a
216 typical developer may make when writing Python code, such as misspelling identifiers,
217 typos, unbalanced parentheses, braces and brackets, bad indentation, missing charac-
218 ters, and using incorrect operators.

219 In all three token mutations, a token is chosen at random. Deletion simply removes
220 the chosen token. Insertion inserts a copy of the chosen token at a random location.
221 Replacement writes a copy of the chosen token over another randomly selected token.
222 In digit, letter and symbol mutations, a single digit, letter, or symbol character is deleted
223 or inserted randomly in the file. In the indentation mutations a line in the file is selected
224 at random and its indentation is decreased or increased. Digit, letter, symbol, and
225 indentation mutations may mutate any part of the file, including code and comments.
226 Token mutations can only affect the code and never comments.

227 Other techniques for mutating Python code are far more limited, such as the opera-
228 tors presented by Derezińska *et al.* [8] or the mutations employed by Pester [27]. These
229 mutators are designed to produce reasonably valid and executable Python programs,
230 whereas the mutations used in the experiments here are not. `UnnaturalCode.py` muta-
231 tions are designed to be as general as possible, they do not guarantee an executable
232 program after mutation. Thus the set of possible of programs we generate is larger than
233 the set that Jester [27] produces. In this paper these mutations are only applied once
234 so that they produce general text files that are similar to known good Python programs.
235 Two types of mutation rules are available: rules that are guaranteed to produce text that
236 Python can successfully lexically analyze, and rules that do not have this guarantee.

237 In order to obtain as much information on whether a Python program is valid or not,
238 while also preventing that program from affecting the operation of `UnnaturalCode.py`,
239 `UnnaturalCode.py` runs code in a separate process. Execution is limited to 10 sec-
240 onds, though this limit was never needed. The mutation testing system has several
241 features that manage the execution of random, unpredictable, Python code and extract

242 error-report data produced by Python and Python programs. Python types may change
243 at runtime, Python does not enforce encapsulation, and Python does not have a stan-
244 dard external debugger. Thus, UnnaturalCode.py executes Python code in a separate
245 process to prevent unknown programs and mutants from changing types, global and
246 local variables used by UnnaturalCode.py. UnnaturalCode.py also obtains debugging
247 information from the process under test, and exports that information back to the main
248 UnnaturalCode.py process for examination. UnnaturalCode.py ensures that all testing
249 processes halt by forcefully killing processes in case they exceed a preset amount of
250 execution time.

251 UnnaturalCode.py depends only on its own code, the implementation of the Python
252 language, and MITLM. MITLM is the only component of UnnaturalCode.py that was
253 preserved from the prototype Java implementation of the tool. UnnaturalCode.py uses
254 a slightly modified version of Python’s own lexical analyzer. It does not use a lexical
255 analyzer generated by a parser-generator from a grammar.

256 Additionally, UnnaturalCode.py is intended not only for research but also for prac-
257 tical use by Python developers. One can download and experiment with Unnatural-
258 Code.py as it is distributed freely on Github: <https://github.com/orezpraw/unnaturalcode>.
259

260 4 EXPERIMENTAL VALIDATION PROCEDURE

261 UnnaturalCode.py is designed to help the Python developer locate simple programming
262 mistakes such as typos. Most of these mistakes are syntactic in nature, although some
263 may be semantic errors, such as misspelled identifiers. Consider the following example
264 Python program:

```
def functionWhichExpectsTwoArguments(a, b):  
    return True  
def testA():  
    functionWhichExpectsTwoArguments("a" "b")  
def testB():  
    functionWhichExpectsTwoArguments("a",-"b")  
def testC():  
    functionWhichExpectsTwoArguments["a", "b"]
```

265 The program listed above executes without error in the Python interpreter and loads
266 without error if imported as a module, indicating that it has basic syntactic validity.
267 However, importing this Python module and running any of the three test functions
268 would result in a `TypeError`. A Python programmer could quickly identify the sim-
269 ple mistakes: there is a comma missing in `testA`; in `testB` there is a stray “-” and
270 in `testC` square brackets were used instead of parentheses.

271 All three of these mistakes would be quickly caught by a compiler at compile time.
272 However, Python must load this file and actually run one of the three broken lines
273 of code in order to discover this mistake. The experimental validation that follows
274 attempts to evaluate UnnaturalCode.py’s ability to locate simple coding mistakes such
275 as these, including both mistakes that Python can catch and mistakes that it cannot.

276 First, known-good Python source files were collected from a variety of Python
277 projects including Django, pip, setuptools and Zope. These are popular, well-known
278 projects that are used in real-world production Python environments. The files collected
279 are assumed to be in good condition, and free of syntax errors. Some files were mod-
280 ified to substitute relative import paths for absolute import paths in order to run. Not
281 every file from these projects was used because some of them required configuration
282 or configuration files, or relied on external libraries not available to the authors, such
283 as Oracle database libraries and could not run without such dependencies. Files with
284 less than 21 tokens were also excluded because they are too short to produce mean-
285 ingful results. 936 files remained after removing files that were inappropriate for the
286 mutation experiments. UnnaturalCode.py will always be able to locate the error in a
287 file shorter than its sliding-window length. Including such files is akin to asking which
288 line a syntax error is on in a file with only one line. The Python files used are available
289 at <https://github.com/orezpraw/pythonCorpus>.

290 Collecting these 936 files required each to be tested, library dependencies installed,
291 run-time dependencies configured, module paths specified or corrected. Programs with
292 errors may not terminate; therefore a limit of 10 seconds was imposed on the execution
293 of any single file. Among the collected test cases no known-good file would run for
294 more than 10 seconds. These known-good files were then used to build a corpus, which
295 was then used to build a 10-gram language model with MITLM. UnnaturalCode.py
296 updates its own corpus as soon as a valid version of a new or changed file is discov-
297 ered. Therefore, UnnaturalCode.py has all known-good versions in the corpus at the
298 beginning of the validation procedure, including files from the current project. Un-
299 naturalCode.py is designed to have the most-recent working versions of all files in its
300 corpus during regular usage (that is the entire set of 936 files). Thus, starting Unnat-
301 uralCode.py with a corpus of all known good files is done to test UnnaturalCode.py's
302 intended use case of running after modifying a Python source file. Each file was then
303 repeatedly subjected to a random-edit mutation and tested against both Python and Un-
304 naturalCode.py.

305 Once a file was mutated, the location of the mutation was noted, and the mutant
306 file was 'required', imported and executed, using *Python 2.7* because the original
307 files were all designed to be run with this Python version. This is similar to running the
308 file from the commandline: `python2.7 filename.py`. This tests if the file parses
309 and executes. One of two possible outcomes was recorded: 1) the file ran successfully;
310 or 2) the file exited with an error. The exact type of error and location were recorded
311 and compared to the location of the mutation, see Table 1. For the sake of brevity, only
312 common errors reported by Python are shown in the results in Table 5.

313 5 RESULTS

314 The data in this section is presented as the fraction of experiments for which the first
315 result returned by Python or UnnaturalCode.py is near the location of the mutation,
316 which is denoted *precision*. Precision is a measure of the performance of an informa-
317 tion retrieval system. In this paper, precision measures how often Python and Unnat-
318 uralCode.py *locate* a mutation. False positive rate is irrelevant to UnnaturalCode.py
319 because it only runs in the presence of an error indicated by Python. Only a single

320 result is considered. Therefore, precision is equal to *recall*, 1-precision, precision at
 321 1 result, and mean reciprocal rank (MRR) with only 1 result. This metric was cho-
 322 sen because Python only produces at most a single result. Therefore, the comparison
 323 would be unfair if UnnaturalCode.py was allowed to produce more than one location
 324 for examination by the user, although it is capable of doing so.

Table 1. Experimental Data Summary

Python Source Files	936
Source Files from django	623
Source Files from pip	159
Source Files from zope	75
Source Files from setuptools	52
Source Files from Python	25
Source Files from markerlib	2
DeleteToken mutations	560400
InsertToken mutations	560400
ReplaceToken mutations	532372
DeleteDigit mutations	79225
InsertDigit mutations	93500
DeleteSymbol mutations	93550
InsertSymbol mutations	93400
DeleteLetter mutations	93550
InsertLetter mutations	93523
Dedent mutations	93450
Indent mutations	93550
Total data points	2386920

325 Table 1 shows some summary statistics about the experimental data gathered. Each
 326 file was subjected to many different mutations of each type in order to obtain a mean
 327 precision value.

328 Table 2 shows the overall performance of Python and UnnaturalCode.py on the 11
 329 types of mutations tested. Each number in the table represents the fraction of injected
 330 errors that were detected. The baseline for all fractions is the total number of errors
 331 injected. `Py` is the Python interpreter. `Py Only` are the errors detected by Python
 332 but not detected by UnnaturalCode.py. Similarly, the fraction of errors detected by
 333 UnnaturalCode.py appears in the `UC` column and the errors exclusively detected by
 334 UnnaturalCode.py are in the `UC Only` column. Then the table shows errors that were
 335 detected by both. `Either` is the union of detection by both methods and `None` are
 336 errors that are not detected.

337 One example of a mutation that occurred during the experiment — that appears in
 338 the `UC`, `UC Only`, and `Either` amounts on the `DeleteToken` row in Table 2 — is the
 339 following code, which is missing the dot operator between `self` and `discard`:

```
def __isub__(self, it):
    if it is self:
        self.clear()
```

Table 2. Fraction of Mutations Located by Mutation Type

	Py	Py Only	UC	UC Only	Both	Either	None
DeleteToken	0.64	0.14	0.65	0.15	0.50	0.79	0.21
InsertToken	0.64	0.09	0.77	0.23	0.55	0.86	0.14
ReplaceToken	0.63	0.13	0.74	0.23	0.51	0.86	0.14
DeleteDigit	0.25	0.01	0.52	0.28	0.24	0.53	0.47
InsertDigit	0.33	0.02	0.62	0.31	0.31	0.64	0.36
DeleteSymbol	0.43	0.15	0.49	0.22	0.27	0.65	0.35
InsertSymbol	0.46	0.14	0.50	0.18	0.32	0.64	0.36
DeleteLetter	0.19	0.03	0.52	0.36	0.17	0.55	0.45
InsertLetter	0.31	0.03	0.58	0.30	0.28	0.61	0.39
Dedent	0.00	0.00	0.09	0.09	0.00	0.09	0.91
Indent	0.33	0.10	0.38	0.15	0.24	0.48	0.52

```

else:
    for value in it:
        selfdiscard(value)
    return self
MutableSet.register(set)

```

340 Python does not report an error when running this code because the block con-
341 taining the mutation is never reached, while UnnaturalCode.py reports the 20-token
342 window indicated by the bold text above.

343 Assuming that UnnaturalCode.py is used in conjunction with the Python interpreter
344 to improve error detection, the important data in Table 2 appear in the `Py` and `Either`
345 column. For instance, for the set of programs used in the evaluation and for the random-
346 edit insertions used, combining UnnaturalCode.py with the Python interpreter would
347 improve the detection of token-replacement errors from 63% to 86%.

348 For all three types of token mutations, Python and UnnaturalCode.py perform com-
349 parably, with the combination locating 9-23% more of the total number of mutations
350 than either Python or UnnaturalCode.py alone. This result is similar to the result ob-
351 tained in [5] where interleaving UnnaturalCode.py and JavaC error messages always
352 improved the score. Though the single-character and indentation mutations are harder
353 for both Python and UnnaturalCode.py to detect, the combination of Python and Unnat-
354 uralCode.py detects the most mutations. Surprisingly, most indentation mutations did
355 not cause errors on execution.

356 Another similarity between these experimental results and the previous results in
357 Campbell *et al.* [5] is that UnnaturalCode.py struggles more with deletion mutations
358 than any other mutation type.

359 Table 2 shows the performance of UnnaturalCode.py and Python under the assump-
360 tion that every mutation is an error. However, this is clearly not the case for some
361 mutations. This provides an upper bound on performance. In order to provide a lower
362 bound, Table 3 shows the performance of Python and UnnaturalCode.py on the 11
363 types of mutations tested, while only counting mutations known to cause an error. By

Table 3. Fraction of Error-Generating Mutations Located by Mutation Type

	Py	Py Only	UC	UC Only	Both	Either	None
DeleteToken	0.85	0.18	0.71	0.04	0.67	0.89	0.11
InsertToken	0.70	0.10	0.81	0.21	0.60	0.91	0.09
ReplaceToken	0.70	0.14	0.75	0.19	0.56	0.89	0.11
DeleteDigit	0.74	0.04	0.83	0.13	0.71	0.87	0.13
InsertDigit	0.75	0.05	0.86	0.16	0.70	0.91	0.09
DeleteSymbol	0.71	0.25	0.56	0.11	0.45	0.82	0.18
InsertSymbol	0.77	0.23	0.63	0.09	0.54	0.86	0.14
DeleteLetter	0.67	0.09	0.73	0.15	0.58	0.82	0.18
InsertLetter	0.72	0.07	0.81	0.16	0.65	0.88	0.12
Dedent	0.00	0.00	0.03	0.03	0.00	0.03	0.97
Indent	0.71	0.21	0.60	0.09	0.50	0.80	0.20

364 removing some data points that were counted against Python, Python’s precision im-
365 proves across the board and especially for deletion mutations. Python errors which
366 are not near the location of the mutation are still counted against Python’s precision.
367 UnnaturalCode.py performs similarly either way for the token mutations, but its extra
368 contribution to precision when combined with Python is reduced for token deletion
369 mutations.

Table 4. Fraction of Error-Generating Mutations Located by Token Type

	Py	Py Only	UC	UC Only	Both	Either	None
NAME	0.75	0.08	0.87	0.20	0.67	0.95	0.05
OP	0.68	0.21	0.62	0.15	0.47	0.83	0.17
NEWLINE	0.17	0.01	0.51	0.34	0.16	0.52	0.48
STRING	0.61	0.05	0.83	0.27	0.56	0.88	0.12
INDENT	0.47	0.05	0.73	0.31	0.41	0.78	0.22
DEDENT	0.16	0.01	0.46	0.30	0.15	0.47	0.53
NUMBER	0.71	0.06	0.87	0.22	0.65	0.93	0.07

370 Table 4 shows the performance of Python and UnnaturalCode.py in relation to the
371 lexical type of the token mutated, for token mutations. For token-replacement mu-
372 tations, the type of the replacement token is considered, not the type of the original
373 token.

374 One interesting result is that Python struggles the most with mutations involving to-
375 ken newlines and indentation tokens. Additionally, both Python and UnnaturalCode.py
376 struggle the least with mutations involving names (identifier tokens).

377 Table 5 shows the frequency at which Python generates different types of errors
378 based on which type of token mutation is performed. The “None” row indicates that
379 Python did not detect a problem in the mutant file.

Table 5. Python Exception Frequency by Mutation Type

	DeleteToken	InsertToken	ReplaceToken
SyntaxError	0.60 (336117)	0.68 (383365)	0.61 (323125)
IndentationError	0.11 (62263)	0.23 (128356)	0.29 (152778)
None	0.25 (139131)	0.09 (47771)	0.09 (47522)
ImportError	0.02 (12338)	0.00 (23)	0.01 (2789)
NameError	0.01 (3177)	0.00 (355)	0.01 (4250)
ValueError	0.01 (4876)	0.00 (100)	0.00 (144)
TypeError	0.00 (1685)	0.00 (339)	0.00 (987)
AttributeError	0.00 (641)	0.00 (32)	0.00 (595)
OptionError	0.00 (5)	0.00 (15)	0.00 (135)
IndexError	0.00 (73)	0.00 (3)	0.00 (12)
error	0.00 (29)	0.00 (14)	0.00 (9)
InvalidInterface	0.00 (42)	0.00 (2)	0.00 (0)
KeyError	0.00 (3)	0.00 (12)	0.00 (5)
HaltingError	0.00 (3)	0.00 (8)	0.00 (6)
ImproperlyConfigured	0.00 (7)	0.00 (0)	0.00 (5)
ArgumentError	0.00 (1)	0.00 (1)	0.00 (4)
UnboundLocalError	0.00 (2)	0.00 (2)	0.00 (1)
OGRException	0.00 (3)	0.00 (0)	0.00 (2)
AssertionError	0.00 (1)	0.00 (0)	0.00 (2)
GEOSException	0.00 (2)	0.00 (0)	0.00 (0)
TemplateDoesNotExist	0.00 (0)	0.00 (1)	0.00 (1)
OSError	0.00 (1)	0.00 (0)	0.00 (0)
FieldError	0.00 (0)	0.00 (1)	0.00 (0)

380 Many token mutations change the indentation of the code in Python. Indentation-
381 Error is a type of SyntaxError, and therefore a large number of mutations result in some
382 type of SyntaxError. The third most common outcome of running a mutant Python file
383 is that no error is raised, and it is followed by relatively rare ImportError, NameEr-
384 rors, and ValueError that are to be expected from mutations affecting library loading,
385 identifiers, and literals. Other types of error occur, however they are extremely rare.

386 Some deletions lead to mutant Python programs that do not contain syntax or seman-
387 tics errors that could be raised by an automatic tool. For instance, deleting a unary “not”
388 operator changes the semantics of the program, but this change does not cause identi-
389 fier, syntax, argument or type errors. The results in Table 5 indicate that **25%** of token
390 deletion mutations yield Python programs that can be successfully ran or imported.
391 For most Python files this means that they successfully defined classes, variables and
392 functions. However, some of the functions defined may not have been executed. It is
393 unlikely that all 25% of the deletions that lead to no error detection actually resulted in
394 error-free programs.

395 In comparison, the digit, symbol, letter and indentation mutations do not always
396 lead to a change in the semantics of the python program. These mutations can occur in

397 comments or literals. Additionally, the indentation mutations, `indent` and `dedent`, may
 398 change a line of code but not affect the block structure if their change is to a single-line
 399 block of code. The results in Table 3 do not include any mutations that did not change
 400 the semantics of the Python program.

Table 6. Fraction of Mutations Located by Exception Type

	Py	Py Only	UC	UC Only	Both	Either	None
SyntaxError	0.78	0.15	0.73	0.11	0.62	0.88	0.12
IndentationError	0.63	0.10	0.83	0.30	0.53	0.93	0.07
None	0.00	0.00	0.49	0.49	0.00	0.49	0.51
ImportError	0.98	0.07	0.93	0.02	0.91	1.00	0.00
NameError	0.95	0.04	0.96	0.05	0.90	1.00	0.00
ValueError	0.99	0.08	0.92	0.00	0.92	1.00	0.00
TypeError	0.45	0.06	0.92	0.53	0.39	0.98	0.02
AttributeError	0.49	0.04	0.87	0.42	0.45	0.91	0.09
OptionError	0.00	0.00	0.85	0.85	0.00	0.85	0.15
IndexError	0.02	0.00	0.90	0.88	0.02	0.90	0.10

401 Table 6 shows the performance of Python and `UnnaturalCode.py` in relation to the
 402 type of error seen by Python. For this measurement, a file for which Python does not
 403 detect an injected error is counted as an irrelevant result. Thus, the “None” reports zero
 404 mutations as detected by Python — it is not the fraction of files that do not contain
 405 mutations.

406 The results in Table 6 indicate that Python’s ability to detect indentation error is
 407 rather poor: this is unsurprising and mirrors the examples shown in `UnnaturalCode` on
 408 Java [5]. While it is difficult to determine whether there is a missing indentation or an
 409 extra indentation (comparable to a missing `{` or an extra `}` in Java) using static analysis,
 410 it is easier for `UnnaturalCode.py` to locate such errors because `UnnaturalCode.py` has
 411 information about the author’s or project’s coding style and history.

412 Unfortunately, it is very difficult to determine why `UnnaturalCode.py` reports a lo-
 413 cation that is not the location of the mutation. This could be because of poor perfor-
 414 mance or because `UnnaturalCode.py` is reporting actual faulty code that is present in
 415 the code as it was released by project authors, before mutation. The experimental re-
 416 sults are computed under the assumption that the code, as it was released by its authors,
 417 is syntax-error-free. Determining whether the code, as released by the authors, is actu-
 418 ally syntax-error-free, or if `UnnaturalCode.py` is reporting locations near bugs that were
 419 shipped would require expert auditing of the locations suggested by `UnnaturalCode.py`.

420 The cumulative proportion of results falling into several range of distances is shown
 421 in Table 7. Each range is a column. Each report that counts toward the “0” proportion
 422 counts also toward the “0-1” proportion, “0-2” proportion, and so on. Python reports a
 423 line number with its error messages, and `UnnaturalCode.py` reports a 20-token window
 424 that may be one line, less than one line, or multiple lines long. The results for `Unnat-`
 425 `uralCode.py` in Table 7 are computed by taking the distance between the mutation and
 426 the beginning of the window reported by `UnnaturalCode.py`.

Table 7. Distance in Lines of Code by Mutation Type

		0	0-1	0-2	0-5	0-10	0-20	>20
DeleteToken	Py	0.85	0.90	0.92	0.93	0.94	0.95	0.05
	UC	0.28	0.57	0.73	0.93	0.98	0.99	0.01
InsertToken	Py	0.70	0.70	0.70	0.71	0.74	0.77	0.23
	UC	0.35	0.65	0.80	0.95	0.99	1.00	0.00
ReplaceToken	Py	0.70	0.73	0.74	0.75	0.77	0.80	0.20
	UC	0.32	0.63	0.78	0.95	0.99	1.00	0.00
DeleteDigit	Py	0.74	0.75	0.76	0.77	0.79	0.81	0.19
	UC	0.20	0.43	0.59	0.77	0.88	0.93	0.07
InsertDigit	Py	0.75	0.82	0.83	0.85	0.86	0.87	0.13
	UC	0.17	0.45	0.63	0.85	0.93	0.97	0.03
DeleteSymbol	Py	0.71	0.84	0.88	0.91	0.91	0.92	0.08
	UC	0.30	0.53	0.67	0.84	0.92	0.95	0.05
InsertSymbol	Py	0.77	0.86	0.88	0.90	0.91	0.92	0.08
	UC	0.24	0.51	0.67	0.86	0.93	0.97	0.03
DeleteLetter	Py	0.67	0.70	0.71	0.73	0.74	0.77	0.23
	UC	0.17	0.44	0.62	0.84	0.92	0.96	0.04
InsertLetter	Py	0.72	0.81	0.82	0.83	0.85	0.86	0.14
	UC	0.18	0.46	0.64	0.85	0.93	0.97	0.03
Dedent	Py	0.00	0.01	0.02	0.05	0.11	0.16	0.84
	UC	0.25	0.48	0.66	0.92	0.98	1.00	0.00
Indent	Py	0.71	0.79	0.82	0.86	0.88	0.90	0.10
	UC	0.08	0.36	0.57	0.84	0.93	0.97	0.03

427 If we only consider the exact line of the error, Python usually outperforms Unnatu-
428 ralCode.py in terms of the location of the error, according to Table 7. However if we
429 consider up to five lines before and after the reported error, UnnaturalCode indicates a
430 line near the error more often than Python.

431 5.1 Comparison to UnnaturalCode with Java

Table 8. MRR Comparison

	Python	Java	UC.py	UC.java
DeleteToken	0.85	0.92	0.74	0.87
InsertToken	0.70	0.93	0.83	0.99
ReplaceToken	0.70	0.93	0.77	0.98

432 Table 8 shows the performance of Python and UnnaturalCode.py compared to the
433 results obtained previously for Java and the prototype version of UnnaturalCode. Only
434 mutants that produced an error in Python are considered. Mean reciprocal ranks are re-
435 ported for Java and both versions of UnnaturalCode. Precision is computed for Python
436 as the proportion of times that Python reports the error on the same line as the mutation.
437 This metric is very similar to *mean reciprocal rank* (MRR) because MRR places as
438 much weight on the first result as it does on the second through last results combined.
439 Thus the results here differ in the methodology used to present results in Campbell *et*
440 *al.* [5] where MRR was used.

441 Given that Python produces at most one result, this metric is somewhat comparable
442 to the mean reciprocal rank. The precision is computed for Python in the exact same
443 way as the MRR was computed for the Java compiler in the case that the Java compiler
444 only returned a single result. Though the Java compiler can produce up to 100 results,
445 it does produce single results occasionally.

446 The results show that the new UnnaturalCode.py system does not perform as well
447 as the prototype version of UnnaturalCode, and that Python’s own error detection mech-
448 anism does not perform as well as the Java compiler error reporting. The distribution
449 of MRR scores for the Java compiler and UnnaturalCode is shown in Figure 2.

450 As was found with the UnnaturalCode prototype for Java, some source files simply
451 seem to be more difficult for UnnaturalCode.py than others. Files that seem to be diffi-
452 cult for UnnaturalCode.py have long sequences of tokens that do not appear elsewhere,
453 such as lists of string literals. Figure 1 shows the distribution of UnnaturalCode.py,
454 Python, and combined “either” precision over the files tested. Python also seems to
455 have difficulties with specific files. Typically, random-edit mutations tend to produce
456 random results from Python’s own parser: regardless of whether Python’s error mes-
457 sages identify the correct location, different mutations tend to cause Python to identify
458 different lines as the source of the location.

459 On some files, however, such as `posixpath.py`, from the Python standard li-
460 brary, the Python parser often reports a parse error at the same location. In `posixpath.py`,
461 this location is on line 182 and 72% of random-edit mutations cause Python to report
462 an error at this location. This effect occurs only when the actual mutant line is after the
463 line that Python’s parser tends to point to. In the case of `posixpath.py`, the line that
464 the Python parser consistently reports is shown below. This may indicate that Python’s
465 own parser struggles with multi-line statements when the line-continuation character \
466 is used.

```
467         return s1.st_ino == s2.st_ino and \
```

468 6 DISCUSSION

469 6.1 UnnaturalCode.py Performance

470 The results described in the previous section show that both UnnaturalCode.py and
471 Python are able to locate mutations made to Python source files most of the time for
472 most mutation types. UnnaturalCode.py is able to identify most mutations that cause
473 a program to fail to execute. UnnaturalCode.py is able to locate some mutants that
474 Python misses. Python is able to locate some mutants that UnnaturalCode.py misses.

475 When used alongside Python, UnnaturalCode.py improves the chances of the correct
476 location being reported to the user. UnnaturalCode.py struggles more with random
477 deletion mutations than any other type of mutation.

478 Despite the fact that the Python version of UnnaturalCode.py has a lower MRR
479 score than the Java version did, the MRR score still indicates that most of the time the
480 correct result is in the top few. *These results indicate that the usefulness of natural*
481 *language techniques used with programming languages includes scripting languages*
482 *like Python.*

483 The Java version of UnnaturalCode was shown to perform much worse on code it
484 had never seen before in Campbell *et al.* [5]. It is safe to assume that UnnaturalCode.py
485 also performs worse on new code that it has not seen before. However, this is not the
486 intended usage of UnnaturalCode.py. UnnaturalCode.py is designed to have good code
487 added automatically to its corpus as soon as it is compiled or run. UnnaturalCode.py's
488 corpus is therefore updated much more often than the project's source code repository.
489 This allows UnnaturalCode.py to follow code evolution very closely. *Testing Unnat-*
490 *uralCode.py with completely unrelated test and training code would not relate to its*
491 *real-world use case.*

492 **6.2 Properties of Python**

493 Python is a very flexible language with many useful and ultimately powerful features.
494 But some of this power limits other aspects and properties of the language. The fol-
495 lowing differences between Python and Java needed to be accounted for in Unnatural-
496 Code.py: 1) Python is not compiled; 2) Python programs may remove or add identifiers
497 to or from any scope during execution; 3) Python types may change during execution;
498 4) Python has indentation-based syntax; 5) Python's lexical analyzer does not produce
499 white-space tokens; 6) run-time dependencies and requirements exceed compile-time
500 dependencies and requirements; and 7) Python only produces at most a single syntax
501 error.

502 **No Compilation** Python is not a compiled language. This means that, in contrast with
503 compiled languages such as Java, there is no oracle for basic Python program validity.
504 The absence of an oracle creates several challenges for UnnaturalCode.py, and requires
505 that both UnnaturalCode.py's goals and implementation be generalized.

506 The main challenge is that there is no way to tell if a Python program is valid.
507 Executing the program will check it for basic syntactic validity, but some issues such
508 as misspelled identifiers can not be checked merely by parsing the program. While
509 there is a fine line between syntactic and semantic errors, due to their simple and typo-
510 driven nature, UnnaturalCode.py is able to assist in locating semantic errors caused by
511 typographical mistakes made when working with identifiers, expressions, values, and
512 indentation. UnnaturalCode.py has no ability to discern between a syntactical Python
513 error and a semantic Python error.

514 Determining if a Python program is valid is an undecidable question because mak-
515 ing such determination requires running the program. So, the approach taken by Un-
516 naturalCode.py is to run the program. Besides being an interactive scripting language,
517 in Python any scope, name-space, type, constant, or library may change at any time
518 during program execution.

519 A Python script may not necessarily halt. Furthermore, a Python program may
520 execute without error even though it contains misspelled identifiers, broken imports, or
521 other problems with the code that are easily caught in any compiled language.

522 The results in the previous section indicate that this is not merely a theoretical foot-
523 note, but that random-edit mutations produce cases in which it is difficult to determine
524 the validity of a Python program a significant fraction ($\frac{1}{4}$) of the time.

525 The interactive, dynamic nature of Python has implications not just for the experi-
526 ments presented in this paper, but also any other experiment that depends on the execu-
527 tion of Python code of unknown quality. Techniques such as genetic programming and
528 mutation would clearly be impacted: in order to be certain that programs produced by
529 a genetic algorithm were valid in basic ways, the programs produced would have to be
530 severely limited in terms of what mutations and combinations were allowed. Indeed,
531 this is the approach taken in Derezińska *et al.* [8], and Jester [27]. For example, one
532 would have to ensure that no programs that referenced a variable before assignment
533 were produced by limiting the set of possible output programs, because use before def-
534 inition can not be checked easily after program generation. Another consideration is
535 that generating Python code with correct indentation, or repairing Python indentation,
536 requires more memory and time than checking pairs of braces.

537 **Dynamic Scope** Even if a given Python program seems well-formed at one point
538 during execution, it may not be at another point. For example, if a program uses an
539 identifier, not only is that identifier not added to any scope until execution, but it may
540 be removed or changed during execution. Holkner *et al.* [13] found in Python software
541 advertised as production-stable quality that variables are created and removed, and their
542 types changed *after* program initialization. Politz *et al.* [31] discuss the differences
543 between Python scope and traditional compiled languages.

544 Though the assumption that types, scopes and names will not change at runtime
545 is tempting to make, it is not a given. For example, Django, a popular Python web-
546 framework, makes use of this flexibility. This flexibility is used to enable developers to
547 inject their own code and types into exceptions when debugging Python code.

548 The presence of runtime changes to built-in exception types in Python programs
549 created a challenge for the implementation of UnnaturalCode.py. UnnaturalCode.py's
550 built-in mutation testing system uses Python's standard `multiprocessing` library
551 to run mutant Python code in a separate process, to implement the 10-second timeout,
552 and to communicate the result of running the mutant Python code back to the main
553 UnnaturalCode.py process. Sending the raw exception caused by the mutant code, if
554 any, is impossible because the type of the exception generated by the mutant code may
555 not exist in the master process, or may exist in a different form. Relying on serialization
556 and de-serialization for inter-process communication is unreliable because the types
557 may be different. They may be different even when both processes are running the
558 same program. Thus, UnnaturalCode.py must convert exceptions and other debugging
559 information to plain strings.

560 **No Constants** The ability to change semantics at run time is not limited to identifiers.
561 Python has no real concept of constants: even `math.pi` (π) may be changed during
562 execution. Furthermore, even the type system is not constant at run time. For exam-

563 ple, because libraries and modules are imported during execution, a Python program
564 with two threads may lose the ability to communicate between its threads if one thread
565 imports something that overwrites a built-in Python type. This is not merely a theo-
566 retical consequence, but one that actually occurs in some of the software used for the
567 experiments presented in this paper.

568 Any changes that the program made to the state of the Python run-time environ-
569 ment — including changes to local scopes, global scope, and types — are discarded by
570 terminating the process after successful execution, exception, or failure to halt. This
571 allows UnnaturalCode.py to protect itself from having its semantics disrupted by other
572 software.

573 UnnaturalCode.py cannot improve the detection of coding mistakes, it only at-
574 tempts to locate them once a problem occurs. Thus, it can be completely bypassed
575 in several ways. The existence of a misspelled identifier in any relevant scope may
576 not be verified if said identifier is used in a conditional block that does not run. Fur-
577 thermore, SyntaxErrors produce exceptions that can be caught and ignored, which will
578 prevent both Python and UnnaturalCode.py from reporting these errors.

579 **Indented Blocks** Python’s use of indentation to define code blocks creates a very
580 different kind of lexical analyzer than most programming languages. While it is con-
581 sidered usual to specify tokens with regular expressions and parse trees with context-
582 free grammars using parser generators such as ANTLR [29], this is not the case for
583 Python. In order to track indentation levels, even the lexical analyzer must be context-
584 sensitive [19].

585 During lexical analysis, Python produces tokens of type “NEWLINE,” “INDENT,”
586 and “DEDENT.” These are roughly equivalent to semi-colons, opening braces and clos-
587 ing braces in Java, respectively. However, they do not necessarily correspond to white-
588 space in a Python source file. In fact, “DEDENT” corresponds to a lack of white-space,
589 and “NEWLINE” is only present once at the end of a multi-line statement.

590 UnnaturalCode.py employs a modified version of the official Python lexical ana-
591 lyzer, `tokenize.py`, which produces token streams even in the following cases: 1)
592 bad indentation; 2) unterminated multi-line statements; and 3) unterminated multi-line
593 literals. The modified lexical analyzer is only used for querying the language model.
594 It is never used to parse or execute Python code. Changes made to the Python lexical
595 analyzer were minimal: it merely continues in whatever state it was in before encoun-
596 tering an error, or assumes that the end of the file terminates multi-line statements and
597 literals.

598 In comparison with other lexical analyzers, such as ANTLR [29], Python does not
599 produce white-space tokens. Instead, white-space is tracked by computing the differ-
600 ence between the end of one token and the beginning of the next. White spaces con-
601 tain semantic information, even in languages where they are syntactically irrelevant.
602 Though this information may be irrelevant to Python, it is relevant to the human devel-
603 opers. Since UnnaturalCode.py cannot see white spaces in Python, its efficacy both in
604 practice and in the experimental validation procedure described in may be negatively
605 impacted.

606 **Run Time Dependencies** Python software may have dependencies at run time. One
607 example are libraries that work with database engines such as MySQL and that require
608 not only MySQL libraries and development files, but also require MySQL to actually
609 be running when the library is loaded with Python's `import` statement. Such require-
610 ments complicates the use of real-world software in validation experiments. Further-
611 more, the order in which libraries are loaded is occasionally important for successful
612 execution: if they are loaded in the wrong order they may break.

613 Validation used only source files from open-source software that could actually run
614 during the experiments. This is a significantly higher bar that source code must meet
615 than merely compiling. One example that was encountered was that some Python mod-
616 ules depend on a configuration module, which must be written by the user, in order
617 to be imported. This configuration model defines run-time behaviour such as database
618 authentication.

619 **One Error** Python only produces a single syntax error at most, and this syntax error
620 is usually without diagnostics. In contrast, most other compilers, such as C and Java
621 compilers, can produce many syntax errors for a single failed compilation. The report-
622 ing of multiple errors employs techniques such as those described by Kim *et al.* [21]
623 and Corchuelo *et al.* [7].

624 To be consistent with Python's behaviour this study reports only one result from
625 `UnnaturalCode.py` to the user. `UnnaturalCode.py` is capable of producing as many re-
626 sults as there are contiguous 20-token sequences in the source file. However, given that
627 Python only reports at most one syntax or other error at a time, while `javac` reported
628 50 or more, `UnnaturalCode.py` is limited to report only one result so that it can be more
629 easily used and compared with Python.

630 Python, could, perhaps be served well by implementing several features to miti-
631 gate the difficulties listed above. First, Python could add a mode similar to Perl's `use`
632 `strict` that forces all variables to be declared with their scope. Second, Python could
633 add variables, similar to Java's `final` variables, which can not be changed once set.
634 Third, Python could prevent the run-time mutation of types, especially built-in and
635 standard library types. The Python project has recently, as of 2014, announced future
636 improvements to Python typing, such as type hinting, which will address this difficulty.
637 Fourth, Python could employ the standard extensions of basic parser technology to pro-
638 duce better parser diagnostics, error recovery, and multiple-error messages in a similar
639 fashion to C, C++, Java and many other languages, as recommended by Cannon [6].
640 Sippu *et al.* [36] describe many of these techniques.

641 7 THREATS TO VALIDITY

642 When comparing these results to the results obtained in `UnnaturalCode Java`, the valid-
643 ity of the comparison is threatened by several factors. First, the metric used in Campbell
644 *et al.* [5], MRR, is not directly comparable or convertible to the metric used in this pa-
645 per: the precision of the only result. The conversion of Python's precision at 1 to an
646 MRR in Table 8 is biased against Python because it produces only one result. Second,
647 the Python lexer produces different token types than Java's lexer. For example, Python's
648 lexer does not produce white-space tokens. This implies that the distribution of seman-

649 tic meanings of the mutations generated for Python code differs from the distribution
650 of semantic meanings of the mutations performed for Java code. Third, Campbell *et*
651 *al.* [5] did not use mutations that passed compilation: rather, any mutation that com-
652 piled was discarded. However, those mutations were included here. Table 3 shows the
653 performance of UnnaturalCode.py on Python programs that are known to be broken.

654 Python provides no mechanism to perform static analysis such as type checking,
655 scope checking, and constant immutability checking. Thus, the experimental evaluation
656 of UnnaturalCode.py provides results at both extremes: results assuming that every
657 mutant would fail the checks in table 2, and results assuming that every mutant that
658 executed successfully would successfully pass the checks in Table 3.

659 As an analysis of random-edit mutations in Python source files, the validity of these
660 results is threatened mainly by the possibility that the mutations made uniformly ran-
661 domly do not reflect the distribution of mistakes made by humans when authoring
662 Python software, which may not be uniformly random. In order to address this concern,
663 both whole-token and single-character mutation experiments were performed. Whole-
664 token mutations represent changes to program structure and single-character mutations
665 represent typographical errors that a human may make.

666 **8 FUTURE WORK**

667 There are many other software engineering tools available for Python and their per-
668 formance or some aspects of their performance may be characterized by using the
669 random-edit-mutation testing tools presented in this paper. These include tools such
670 as `pyflakes`, and `pylint`. UnnaturalCode.py could be integrated into such tools as
671 another method of critiquing code. A scientific comparison of the use of random-edit-
672 mutation tool as opposed to a tool with fewer, specialized mutations such as Pester [27],
673 to characterize the coverage and effectiveness of test suites would also be useful. Just *et*
674 *al.* [20] has found that mutations are representative of errors introduced by human soft-
675 ware developers.

676 An extension of UnnaturalCode.py intended to provide feedback on possible errors
677 while an author types would be useful. This feedback could then be compared to other
678 systems that give similar feedback such as the Eclipse integrated development environ-
679 ment.

680 **9 CONCLUSIONS**

681 Python, and other dynamically typed, interactive, scripting languages, pose many chal-
682 lenges to the Software Engineering experimenter. This paper provides a tool, and some
683 techniques, for addressing these challenges. This tool may also be used when experi-
684 menting with genetic programming or other types of automatic code generation. These
685 techniques may be applied when languages with a robust compiler is not available or
686 desired. If a robust compiler is available, it should be employed. Because such a com-
687 piler can be used as an oracle for basic program validity including syntactic validity
688 and identifier validity.

689 UnnaturalCode.py can augment Python's own error messages by suggesting a loca-
690 tion for the developer to check. Often that is the location of a mistake, as shown in the

691 “UC” column of Table 3. Using UnnaturalCode.py with Python improves the chances
692 of locating a mutant piece, or mistaken piece, of code. Therefore UnnaturalCode.py
693 can be a valuable tool for debugging language implementations such as Python be-
694 cause syntax errors appear unnatural to language models trained with working source
695 code. The n -gram language models can help find syntax errors when the model is built
696 from a corpus of code known to be correct. Additionally, UnnaturalCode.py is able to
697 locate some semantic errors, such as type errors induced by typos.

698 Python is a language that raises challenges for mutant-based experimentation: for
699 instance, Python does not report as faulty 25% of programs with a single missing token.
700 Thus, when performing mutation-based experiments with “scripting” languages such as
701 Python, researchers must be aware of the issues discussed in this paper. Typical code
702 errors that would be caught quickly and automatically by a compiler for a language
703 such as Java can be difficult to automatically discover and report in Python.

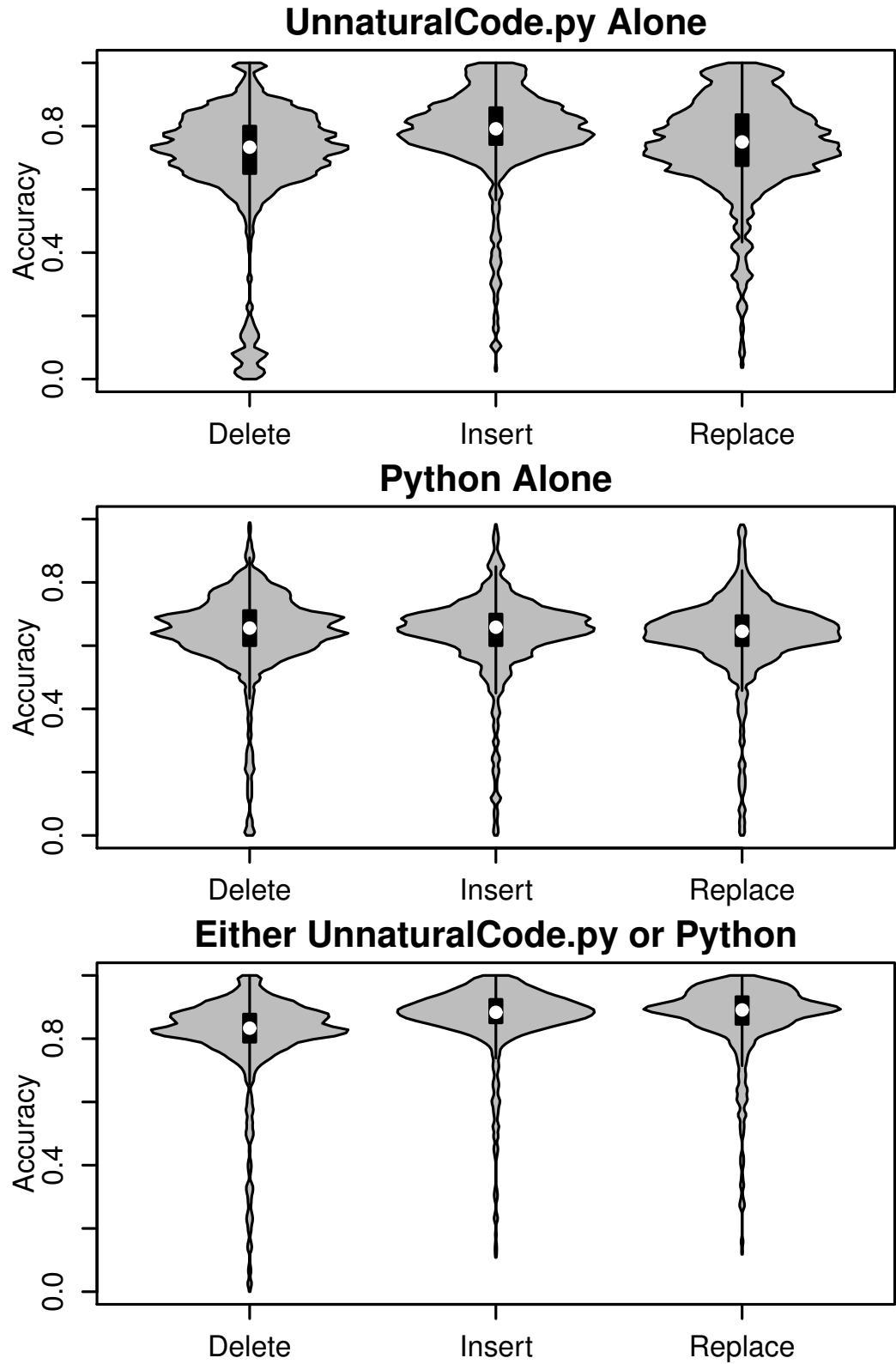


Figure 1. Independent and combined performance of UnnaturalCode.py and Python. Width indicates the relative number of source files on which the precision level is achieved.

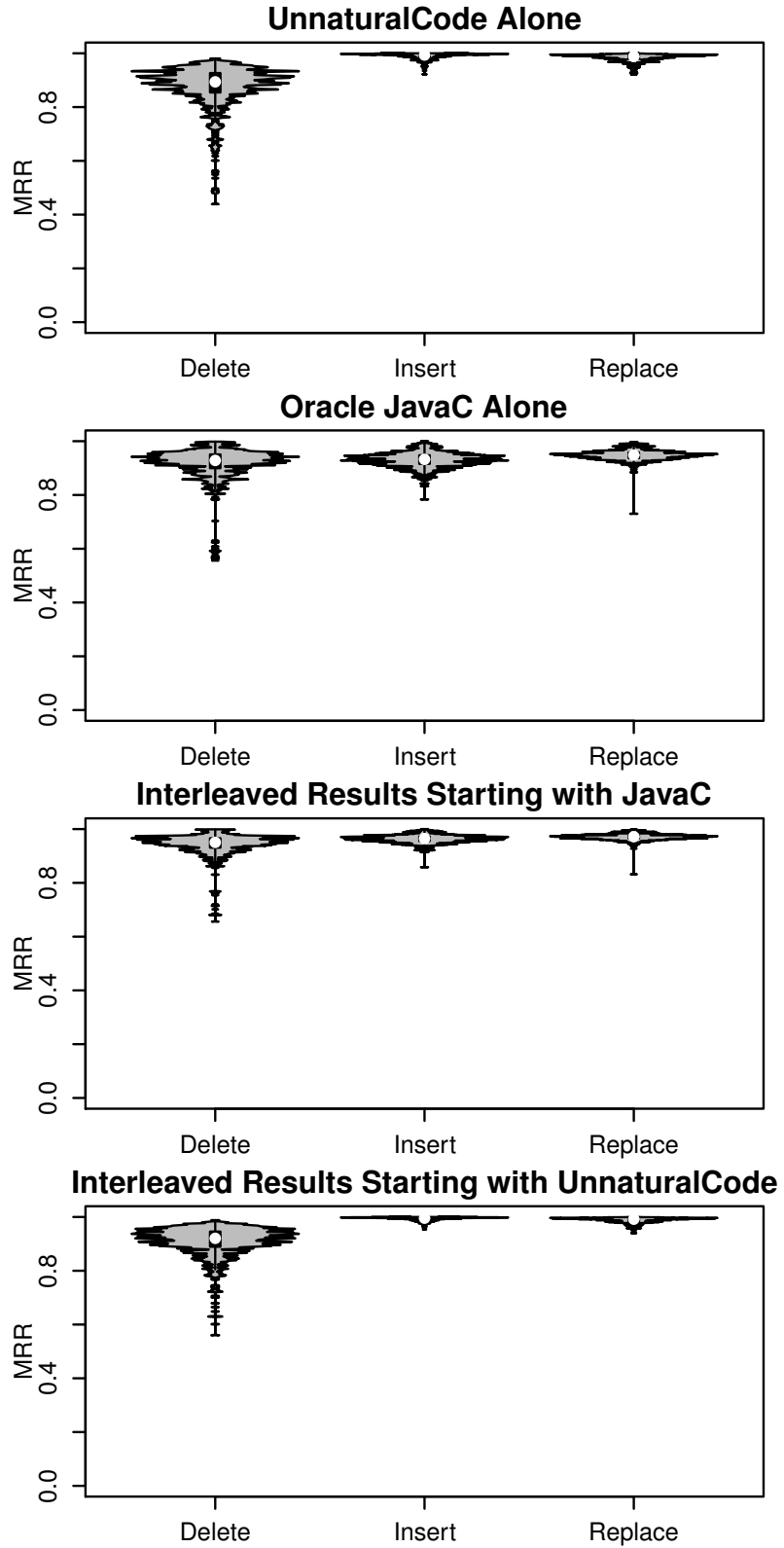


Figure 2. Independent and combined performance of UnnaturalCode Java prototype and JavaC, from Campbell *et al.* [5]. Width indicates the relative number of source files on which the MRR score is achieved.

REFERENCES

- [1] D. Ancona, M. Ancona, A. Cuni, and N. D. Matsakis. Rpython: a step towards reconciling dynamically and statically typed oo languages. In *Proceedings of the 2007 symposium on Dynamic languages*, pages 53–64. ACM, 2007.
- [2] K. Beck. *Test-driven development : by example*. Addison-Wesley, Boston, 2003.
- [3] T. Bell, P. Andreae, and A. Robins. Computer science in nz high schools: The first year of the new standards. In *Proceedings of the 43rd ACM Technical Symposium on Computer Science Education, SIGCSE '12*, pages 343–348, New York, NY, USA, 2012. ACM.
- [4] M. G. Burke and G. A. Fisher. A practical method for LR and LL syntactic error diagnosis and recovery. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 9(2):164–197, Mar. 1987.
- [5] J. C. Campbell, A. Hindle, and J. N. Amaral. Syntax Errors Just Aren't Natural: Improving Error Reporting with Language Models. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, page 10, 2014.
- [6] B. Cannon. Localized type inference of atomic types in python. Master's thesis, California Polytechnic State University, 2005.
- [7] R. Corchuelo, J. A. Pérez, A. Ruiz, and M. Toro. Repairing syntax errors in LR parsers. *ACM Trans. Program. Lang. Syst.*, 24(6):698–710, Nov. 2002.
- [8] A. Derezińska and K. Hałas. Operators for mutation testing of python programs. *ICS Research Report*, 2014.
- [9] S. Garner, P. Haden, and A. Robins. My program is correct but it doesn't run: a preliminary investigation of novice programmers' problems. In *Proceedings of the 7th Australasian conference on Computing education-Volume 42*, pages 173–180. Australian Computer Society, Inc., 2005.
- [10] S. L. Graham, C. B. Haley, and W. N. Joy. Practical LR error recovery. *SIGPLAN Not.*, 14(8):168–175, Aug. 1979.
- [11] B. J. Heeren. *Top Quality Type Error Messages*. PhD thesis, Universiteit Utrecht, Netherlands, 2005.
- [12] A. Hindle, E. Barr, Z. Su, M. Gabel, and P. Devanbu. On the naturalness of software. In *Software Engineering (ICSE), 2012 34th International Conference on*, pages 837–847, June 2012.
- [13] A. Holkner and J. Harland. Evaluating the dynamic behaviour of python applications. In *Proceedings of the Thirty-Second Australasian Conference on Computer Science - Volume 91, ACSC '09*, pages 19–28, Darlinghurst, Australia, Australia, 2009. Australian Computer Society, Inc.
- [14] M. Hristova, A. Misra, M. Rutter, and R. Mercuri. Identifying and correcting Java programming errors for introductory computer science students. *ACM SIGCSE Bulletin*, 35(1):153–156, 2003.
- [15] B. Hsu and J. Glass. Iterative language model estimation: efficient data structure & algorithms. 2008.
- [16] J. Jackson, M. Cobb, and C. Carver. Identifying top java errors for novice programmers. In *Frontiers in Education, 2005. FIE'05. Proceedings 35th Annual Conference*, pages T4C–T4C. IEEE, 2005.

- 748 [17] M. C. Jadud. A first look at novice compilation behaviour using bluej. *Computer*
749 *Science Education*, 15(1):25–40, 2005.
- 750 [18] M. C. Jadud. Methods and tools for exploring novice compilation behaviour. In
751 *Proceedings of the second international workshop on Computing education re-*
752 *search*, pages 73–84. ACM, 2006.
- 753 [19] T. Jim. Python is not context free. [http://trevorjim.com/](http://trevorjim.com/python-is-not-context-free/)
754 [python-is-not-context-free/](http://trevorjim.com/python-is-not-context-free/), 2012.
- 755 [20] R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, and G. Fraser. Are
756 mutants a valid substitute for real faults in software testing. In *22nd International*
757 *Symposium on the Foundations of Software Engineering (FSE 2014)*, 2014.
- 758 [21] I.-S. Kim and K.-M. Choe. Error repair with validation in LR-based parsing. *ACM*
759 *Trans. Program. Lang. Syst.*, 23(4):451–471, July 2001.
- 760 [22] D. Krpan and I. Bilobrk. Introductory programming languages in higher education.
761 In *MIPRO, 2011 Proceedings of the 34th International Convention*, pages 1331–
762 1336. IEEE, 2011.
- 763 [23] S. K. Kummerfeld and J. Kay. The neglected battle fields of syntax errors. In
764 *Proceedings of the fifth Australasian conference on Computing education-Volume*
765 *20*, pages 105–111. Australian Computer Society, Inc., 2003.
- 766 [24] B. S. Lerner, M. Flower, D. Grossman, and C. Chambers. Searching for type-error
767 messages. In *Conference on Programming Language Design and Implementation*
768 *(PLDI)*, pages 425–434, San Diego, CA, USA, 2007.
- 769 [25] L. McIver. The effect of programming language on error rates of novice program-
770 mers. In *12th Annual Workshop of the Psychology of Programming Interest Group*,
771 pages 181–192. Citeseer, 2000.
- 772 [26] K. J. Millman and M. Aivazis. Python for scientists and engineers. *Computing in*
773 *Science & Engineering*, 13(2):9–12, 2011.
- 774 [27] I. R. Moore. Jester: the junit test tester. [http://jester.sourceforge.](http://jester.sourceforge.net/)
775 [net/](http://jester.sourceforge.net/), 2000–2013.
- 776 [28] T. Parr and K. Fisher. Ll(*): The foundation of the antlr parser generator. In
777 *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language*
778 *Design and Implementation, PLDI ’11*, pages 425–436, New York, NY, USA, 2011.
779 ACM.
- 780 [29] T. J. Parr and R. W. Quong. Antlr: A predicated-ll(k) parser generator. *Software:*
781 *Practice and Experience*, 25(7):789–810, 1995.
- 782 [30] F. Pérez, B. E. Granger, and J. D. Hunter. Python: an ecosystem for scientific
783 computing. *Computing in Science & Engineering*, 13(2):13–21, 2011.
- 784 [31] J. G. Politz, A. Martinez, M. Milano, S. Warren, D. Patterson, J. Li, A. Chitipothu,
785 and S. Krishnamurthi. Python: The full monty. *SIGPLAN Not.*, 48(10):217–232,
786 Oct. 2013.
- 787 [32] A. Rigo. Representation-based just-in-time specialization and the psyco prototype
788 for python. In *Proceedings of the 2004 ACM SIGPLAN symposium on Partial*
789 *evaluation and semantics-based program manipulation*, pages 15–26. ACM, 2004.
- 790 [33] A. Rigo and S. Pedroni. Pypy’s approach to virtual machine construction. In *Com-*
791 *panion to the 21st ACM SIGPLAN Symposium on Object-oriented Programming*
792 *Systems, Languages, and Applications, OOPSLA ’06*, pages 944–953, New York,
793 NY, USA, 2006. ACM.

- 794 [34] A. Rigo and C. Tismer. Psycho, the python specializing compiler,. <http://psyco.sourceforge.net/>, 2001.
- 795
- 796 [35] M. Salib. Faster than c: Static type inference with starkiller. In *in PyCon Proceedings, Washington DC*, pages 2–26. SpringerVerlag, 2004.
- 797
- 798 [36] S. Sippu and E. Soisalon-Soininen. *Parsing Theory: LR(k) and LL(k) Parsing*. EATCS monographs on theoretical computer sciences: European Association for
- 799 Theoretical Computer Science. Springer, 1990.
- 800
- 801 [37] F. Stajano. Python in education: Raising a generation of native speakers. In *Proceedings of 8 th International Python Conference*, pages 2000–01, 2000.
- 802
- 803 [38] E. S. Tabanao, M. M. T. Rodrigo, and M. C. Jadud. Identifying at-risk novice java
- 804 programmers through the analysis of online protocols. In *Philippine Computing*
- 805 *Science Congress*, 2008.
- 806 [39] E. S. Tabanao, M. M. T. Rodrigo, and M. C. Jadud. Predicting at-risk novice
- 807 java programmers through the analysis of online protocols. In *Proceedings of*
- 808 *the seventh international workshop on Computing education research, ICER '11*,
- 809 pages 85–92, New York, NY, USA, 2011. ACM.
- 810 [40] W. Weimer, T. Nguyen, C. Le Goues, and S. Forrest. Automatically finding patches
- 811 using genetic programming. In *Proceedings of the 31st International Conference*
- 812 *on Software Engineering*, pages 364–374, Washington, DC, USA, 2009. IEEE
- 813 Computer Society.