# An Exploratory Study of Ad Hoc Parsers in Python

https://arxiv.org/abs/2304.09733

**Michael Schröder**
TU Wien
Vienna, Austria
michael.schroeder@tuwien.ac.at

**Marc Goritschnig**
TU Wien
Vienna, Austria

**Jürgen Cito**
TU Wien
Vienna, Austria
juergen.cito@tuwien.ac.at

**MSR 2023 Registered Report**
**Melbourne, Australia**

TU WIEN Informatics

```python
693        """
694    if string_network.count('/') == 1:
695        try:
696            mask = int(string_network.split('/')[1])
697        except ValueError:
698            return False
699
```

```objc
118        } else {
119            NSRange range = [component rangeOfString:@"max-age="];
120            if (range.location != NSNotFound) {
121                NSInteger seconds = [[component substringFromIndex:range.location + range.length] integerValue];
122                staleTime = [originalDate dateByAddingTimeInterval:(NSTimeInterval)seconds];
123            }
124        }
```

```javascript
169    * @param {string} path
170    */
171    function extractNameFromPath (path) {
172      return path && path !== 'root' ? path.split('/').slice(-2, -1)[0] : 'Root'
173    }
174
175    /**
```

```typescript
21
22    export default function parseUserAgent(ua = navigator.userAgent) {
23        ua = ua.toLowerCase()
24        const osRe = OS_REGEXPS.find((re) => re.test(ua))
25        let [, os = 'other', os_version = '0'] = (osRe && ua.match(osRe)) || []
26        if (os === 'iphone os' || os === 'ipad os') os = 'ios'
27        const browserRe = BROWSER_REGEXPS.find((re) => re.test(ua))
```

# Parsers are everywhere!

```cpp
480                while(str[i] == '\033') {
481                    do {
482                        i++;
483                    } while(i < str.length() && str[i] != 'm');
484                    i++;
485                    if(i >= str.length()) break;
486                }
```

```python
372            if "@" not in value or value.startswith("@") or value.endswith("@"):
373                return False
374    try:
375            p1, p2 = value.split("@")
376        except ValueError:
377            # value contains more than one @.
378            return False
```

```python
35                Actor
36        """
37        if re.search(r'<.+>', string):
38            m = re.search(r'(.*) <(.+?)>', string)
39            name, email = m.groups()
40            return Actor(name, email)
41        else:
```

```javascript
1    const React = window.React;
2
3    function csv(string) {
4      return string.split(/\s*,\s*/);
5    }
6
7    export default function IssueList({issues}) {
```

# Step 1: Program Slicing

- identify string variables via known functions

```python
669   def get_compiler_version(env):
670       """
671       Returns an array of version numbers as ints: [major, minor, patch].
672       The return array should have at least two values (major, minor).
673       """
674       if not env.msvc:
675           # Not using -dumpversion as some GCC distros only return major, and
676           # Clang used to return hardcoded 4.2.1: # https://reviews.llvm.org/D56803
677           try:
678               version = subprocess.check_output([env.subst(env["CXX"]), "--version"]).strip().decode("utf-8")
679           except (subprocess.CalledProcessError, OSError):
680               print("Couldn't parse CXX environment variable to infer compiler version.")
681               return None
682       else:  # TODO: Implement for MSVC
683           return None
684       match = re.search("[0-9]+\.[0-9.]+", version)
685       if match is not None:
686           return list(map(int, match.group().split(".")))
687       else:
688           return None
```

https://github.com/neshume/godot/blob/e43c867/methods.py
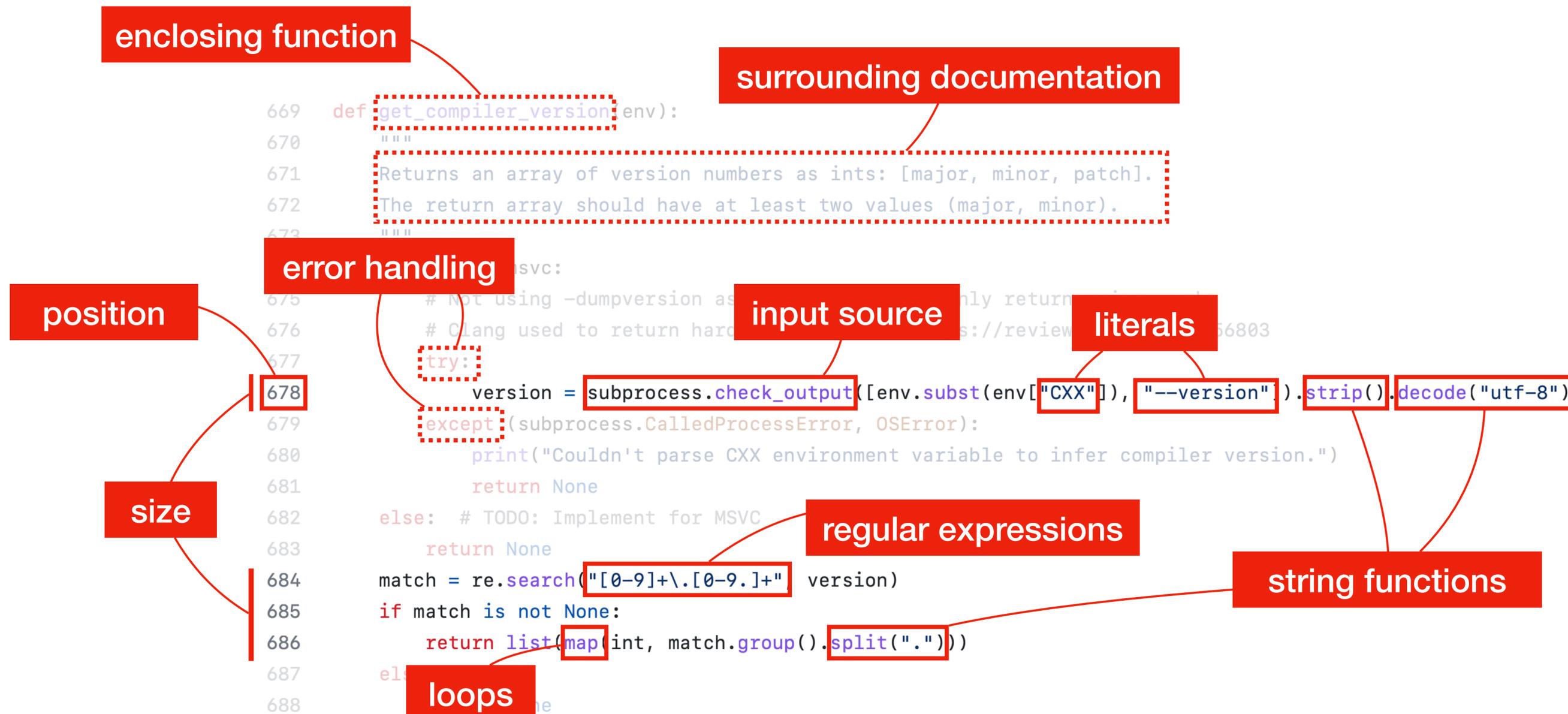
3

# Step 1: Program Slicing

- identify string variables via known functions

- construct forward slice based on data dependencies

```python
669    def get_compiler_version(env):
670        """
671        Returns an array of version numbers as ints: [major, minor, patch].
672        The return array should have at least two values (major, minor).
673        """
674        if not env.msvc:
675            # Not using --dumpversion as some GCC distros only return major, and
676            # Clang used to return hardcoded 4.2.1: # https://reviews.llvm.org/D56803
677            try:
678                version = subprocess.check_output([env.subst(env["CXX"]), "--version"]).strip().decode("utf-8")
679            except (subprocess.CalledProcessError, OSError):
680                print("Couldn't parse CXX environment variable to infer compiler version.")
681                return None
682        else:  # TODO: Implement for MSVC
683            return None
684    match = re.search("[0-9]+\.[0-9.]+", version)
685    if match is not None:
686        return list(map(int, match.group().split(".")))
687    else:
688        return None
```

https://github.com/neshume/godot/blob/e43c867/methods.py

4

# Step 2: Feature Extraction

- 25 metrics across 9 research questions (initial list)



```python
669   def get_compiler_version(env):
670       """
671       Returns an array of version numbers as ints: [major, minor, patch].
672       The return array should have at least two values (major, minor).
673       """
674           msvc:
675           # Not using -dumpversion as        nly return    
676           # Clang used to return hard                s://review    56803
677           try:
678           version = subprocess.check_output([env.subst(env["CXX"]), "--version"]).strip().decode("utf-8")
679           except (subprocess.CalledProcessError, OSError):
680               print("Couldn't parse CXX environment variable to infer compiler version.")
681               return None
682       else:  # TODO: Implement for MSVC
683           return None
684   match = re.search("[0-9]+\.[0-9.]+"  version)
685   if match is not None:
686       return list(map(int, match.group().split(".")))
687   els
688              ne
```

**enclosing function**

**surrounding documentation**

**error handling**

**position**

**input source**

**literals**

**size**

**regular expressions**

**string functions**

**loops**

https://github.com/neshume/godot/blob/e43c867/methods.py

5

# Step 2: Feature Extraction

- 25 metrics across 9 research questions (initial list)

**Table 1: Initial list of metrics extracted for each ad hoc parser.**

| RQs | Metric | Description |
|---|---|---|
| 1 2 | Project Name | name of the project containing the ad hoc parser |
| 1 | Project LOC | total lines of code in the containing project |
| 2 | Module Name | name of the enclosing module/file |
| 2 | EF Name | name of the enclosing function |
| 2 | EF LOC | total lines of code in the enclosing function |
| 2 | Position | position of the ad hoc parser within the enclosing function |
| 1 2 3 | LOC | lines of code in the ad hoc parser |
| 3   6 | CYCLO | cyclomatic complexity of the ad hoc parser |
| 2   4 | Input Source | source of the input string: EF argument, global variable, function call, etc. |
| 2   4 | Input Origin | origin of the input string: command-line, file, environment variable, etc. |
| 3 | Expression Count | number of expressions in the ad hoc parser |
| 3 | Variable Count | number of variables in the ad hoc parser |
| 3 | Function Count | number of function calls in the ad hoc parser |
| 5 6 7 8 | Function Names | names of all functions called in the ad hoc parser |
| 5 | Function Origins | origin of each called function: user-defined or from a library |
| 5 6 | Function Positions | position of all function calls within the ad hoc parser |
| 5 | Function Arguments | arguments with which each function is called, besides the input string |
| 5   8 | Syntactic Sugar | special syntax used in the ad hoc parser: subscript notation, tuples, list comprehensions, etc. |
| 6 | Regular Expressions | arguments to known regex functions or regex literals used in the ad hoc parser |
| 6 7 | Loop Bounds | constant, linear on input string, complex, or unbounded |
| 7 | Loop Types | for, while, functional (map, split, etc.), or recursive |
| 6 7 | Loop Nesting Depth | how deeply nested loops in the ad hoc parser are |
| 8 | Caught Exceptions | all exceptions caught by the ad hoc parser or the enclosing function |
| 8 | Uncaught Exceptions | all uncaught exceptions (excluding explicitly raised...) |
| 8 | Raised Exceptions | |

# An Exploratory Study of Ad Hoc Parsers in Python

https://arxiv.org/abs/2304.09733

- reveal common syntactic and semantic characterstics of ad hoc parsing code

- understand the nature of ad hoc parsers

- inform future program analysis efforts

**SNEAK PREVIEW**

### automatic grammar inference

| Inferred Grammar | Inferred Inputs |
|---|---|
| $s \rightarrow int \mid int \text{ , } s$ | ✖ (empty) |
| $int \rightarrow space^* \ (+ \mid -)^? \ digit \ (\_^? \ digit)^* \ space^*$ | ✔ `1,2,3` |
| $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$ | ✔ `10_000,4` |
| $space \rightarrow \text{␣} \mid \text{\t} \mid \text{\n} \mid \text{\v} \mid \text{\f} \mid \text{\r}$ | ✔ `+01_2,␣␣3␣` |

```python
xs = map(int, s.split(","))
```