

Grammars for Free

Toward Grammar Inference for Ad Hoc Parsers

<https://arxiv.org/abs/2202.01021>



Michael Schröder

TU Wien

Vienna, Austria

michael.schroeder@tuwien.ac.at



Jürgen Cito

TU Wien and Meta Platforms, Inc.

Vienna, Austria

juergen.cito@tuwien.ac.at

New Ideas and Emerging Results, ICSE 2022

Pittsburgh, PA, USA



Informatics

```
psf/requests > requests/utils.py 47 matches Python main
693 """
694 if string_network.count('/') == 1:
695     try:
696         mask = int(string_network.split('/')[1])
697     except ValueError:
698         return False
699
```

```
cortinico/react-native > Libraries/Image/RCTImageCache.m 2 matches Objective-C main
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 }
```

```
vuejs/vuex > src/plugins/devtool.js 6 matches JavaScript main
169 * @param {string} path
170 */
171 function extractNameFromPath (path) {
172     return path && path !== 'root' ? path.split('/').slice(-2, -1)[0] : 'Root'
173 }
174
175 /**
```

```
x10an14/docs > javascripts/user-agent.ts 6 matches TypeScript main
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!

```
coolreader18/libqalculate > src/qalc.cc 86 matches C++ master
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 }
```

```
django/django > django/utils/html.py 1 match Python main
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
```

```
packages/GitPython/lib/git/actor.py 3 matches 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:
```

```
fixtures/dom/src/components/IssueList.js 1 match 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}) {
```

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

"1,2,3"

[1,2,3]

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

The image shows a code snippet with two annotations. A red 'X' is positioned to the left of the variable 'xs', with an arrow pointing from 'xs' to the 'X'. Above the second double quote in the split function argument, there are two more double quotes, with an arrow pointing from them to the original double quote.

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

↑
"+01_2,3_0_4 "

↓
[12,304]

```
class int([x])
class int(x, base=10)
```

Return an integer object constructed from a number or string *x*, or return 0 if no argument given. If *x* defines `__int__()`, `int(x)` returns `x.__int__()`. If *x* defines `__index__()`, it returns `x.__index__()`. If *x* defines `__trunc__()`, it returns `x.__trunc__()`. For floating point numbers, this truncates towards zero.

If *x* is not a number or if *base* is given, then *x* must be a string, `bytes`, or `bytearray` instance representing an **integer literal** in radix *base*. Optionally, the literal can be preceded by + or - (space in between) and surrounded by whitespace. A base-*n* literal consists of the digits 0 to *n*-1 (or A to Z for bases above 10) having values 0 to *n*-1. The default *base* is 10. The allowed values for *base* are 2, 8, 10, or 16, and so that `int('010', 0)` is not legal, while `int('010')` is as well.

```
map(function, iterable, ...)
Return an iterator that applies function to the items from all iterables in parallel. With multiple iterables, the iterator stops when the shortest iterable is exhausted. For cases where the function inputs are already arranged into argument tuples, see itertools.starmap().
```

```
str.split(sep=None, maxsplit=-1)
```

Return a list of the words in the string, using *sep* as the delimiter string. If *maxsplit* is given, at most *maxsplit* splits are done (thus, the list will have at most `maxsplit+1` elements). If *maxsplit* is not specified or -1, then there is no limit on the number of splits (all possible splits are made).

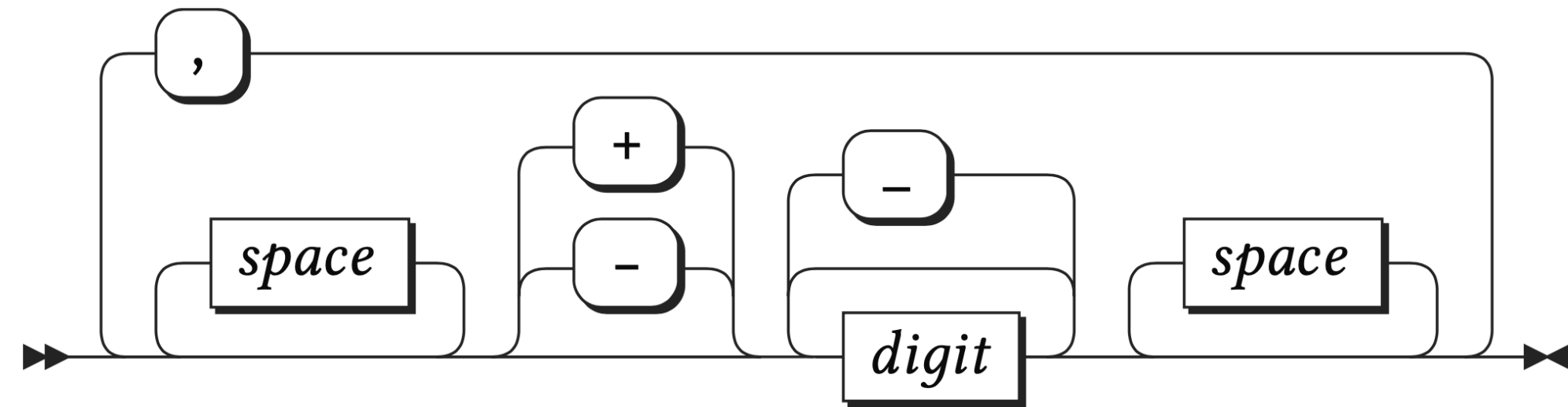
If *sep* is given, consecutive delimiters are not grouped together and are deemed to delimit empty strings (for example, `'1,,2'.split(',')` returns `['1', '', '2']`). The *sep* argument may consist of multiple characters (for example, `'1<>2<>3'.split('<>')` returns `['1', '2', '3']`). Splitting an empty string with a specified separator returns `['']`.

For example:

```
>>> '1,2,3'.split(',')
['1', '2', '3']
>>> '1,2,3'.split(',', maxsplit=1)
['1', '2,3']
>>> '1,2,,3'.split(',')
['1', '2', '', '3']
```

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

"1,2,3"	→	[1,2,3]	✓
" +01_2,3_0_4 "	→	[12,304]	✓
" "	→	✗	
	⋮		



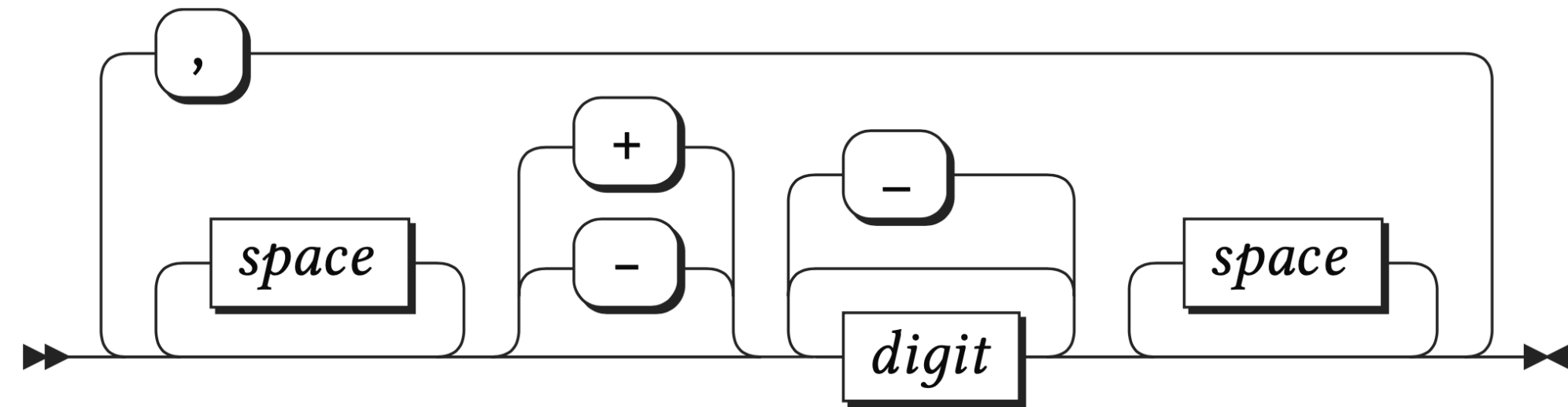
$s \rightarrow int \mid int , s$
 $int \rightarrow space^* sign^? digit (_? digit)^* space^*$
 $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
 $sign \rightarrow + \mid -$
 $space \rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r$

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

"1,2,3"	→	[1,2,3]	✓
" +01_2,3_0_4 "	→	[12,304]	✓
" "	→	✗	
	⋮		

The Need for Grammars

- Program Comprehension & Documentation
 - high-level data-centric perspective
 - diverse notations: regex, ABNF, PEG, state machine, railroad diagram,...
- Fuzzing
 - systematically generate random test inputs
 - grammar-based fuzzers can penetrate into deep program states (past syntactic checks)
- Language-Theoretic Reasoning
 - safety of parsers is connected to theoretic properties (e.g., computability bounds)
 - input languages should be minimally powerful
 - <https://langsec.org> [LANGSEC] regards the Internet insecurity epidemic as a consequence of *ad hoc* programming of input handling [...]



$$s \rightarrow int \mid int , s$$

$$int \rightarrow space^* sign^? digit (_? digit)^* space^*$$

$$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

$$sign \rightarrow + \mid -$$

$$space \rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r$$

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

Parser : Grammar \approx Function : Type

Type Inference

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

[Int]

String

✨ Grammar Inference ✨

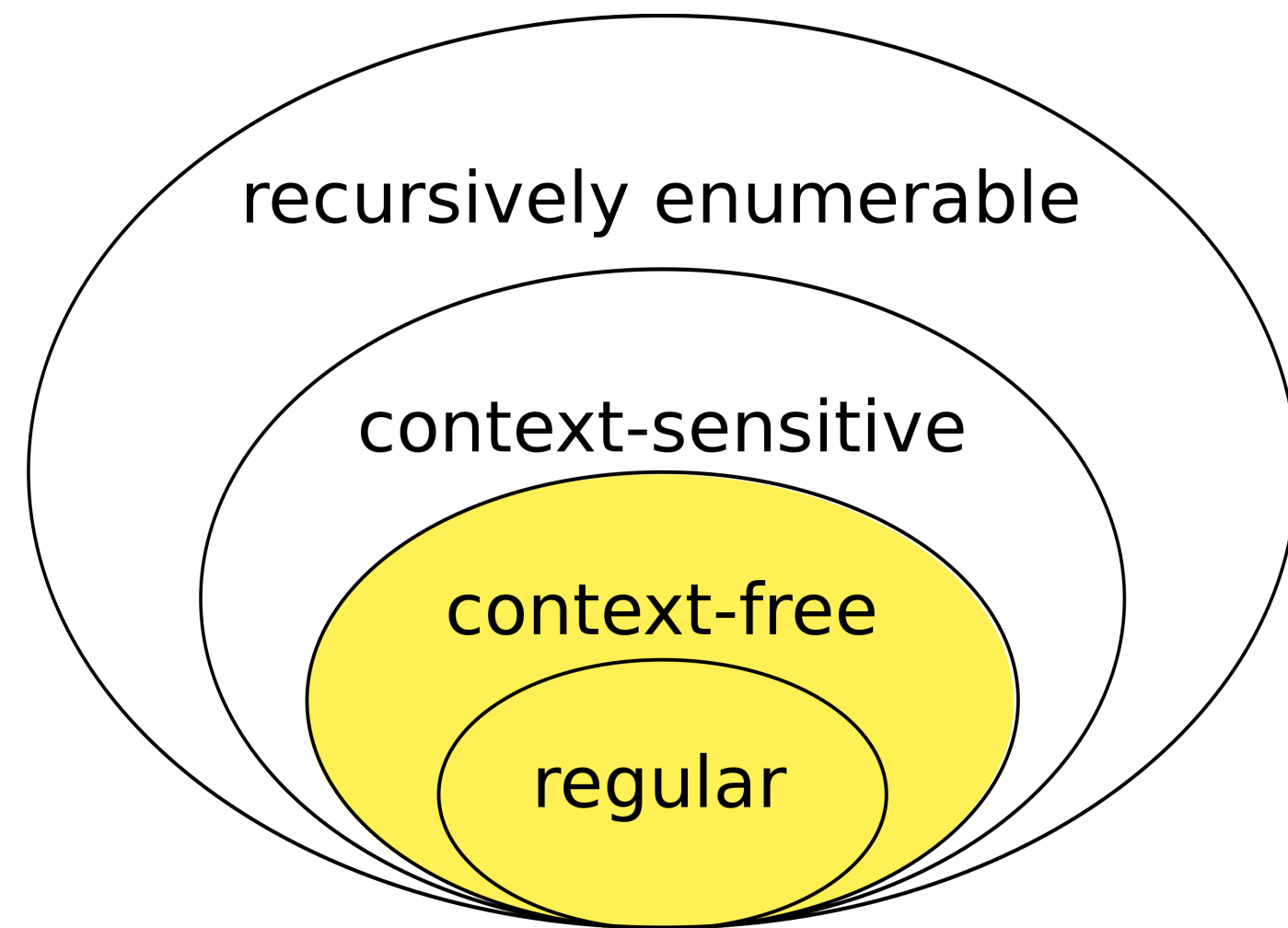
```
 $s \rightarrow int \mid int , s$   
 $int \rightarrow space^* sign^? digit (_? digit)^* space^*$   
 $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$   
 $sign \rightarrow + \mid -$   
 $space \rightarrow \_ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r$ 
```

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

[Int]

String { }

Toward Grammar Inference



The Chomsky hierarchy

$$L = \{a^n b^n \mid n > 0\}$$

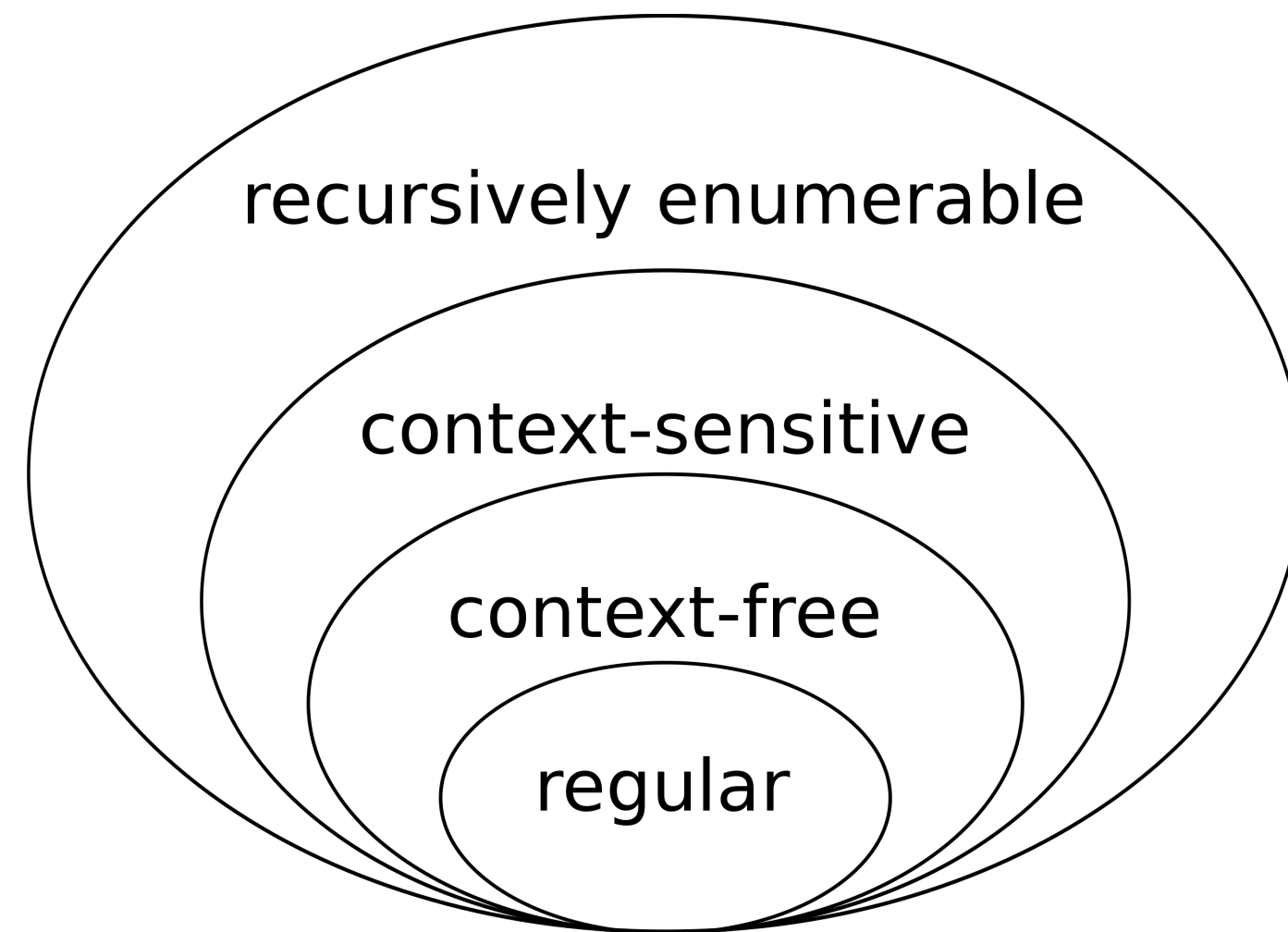
Language

Grammar

$S \rightarrow aSb$
 $S \rightarrow \epsilon$

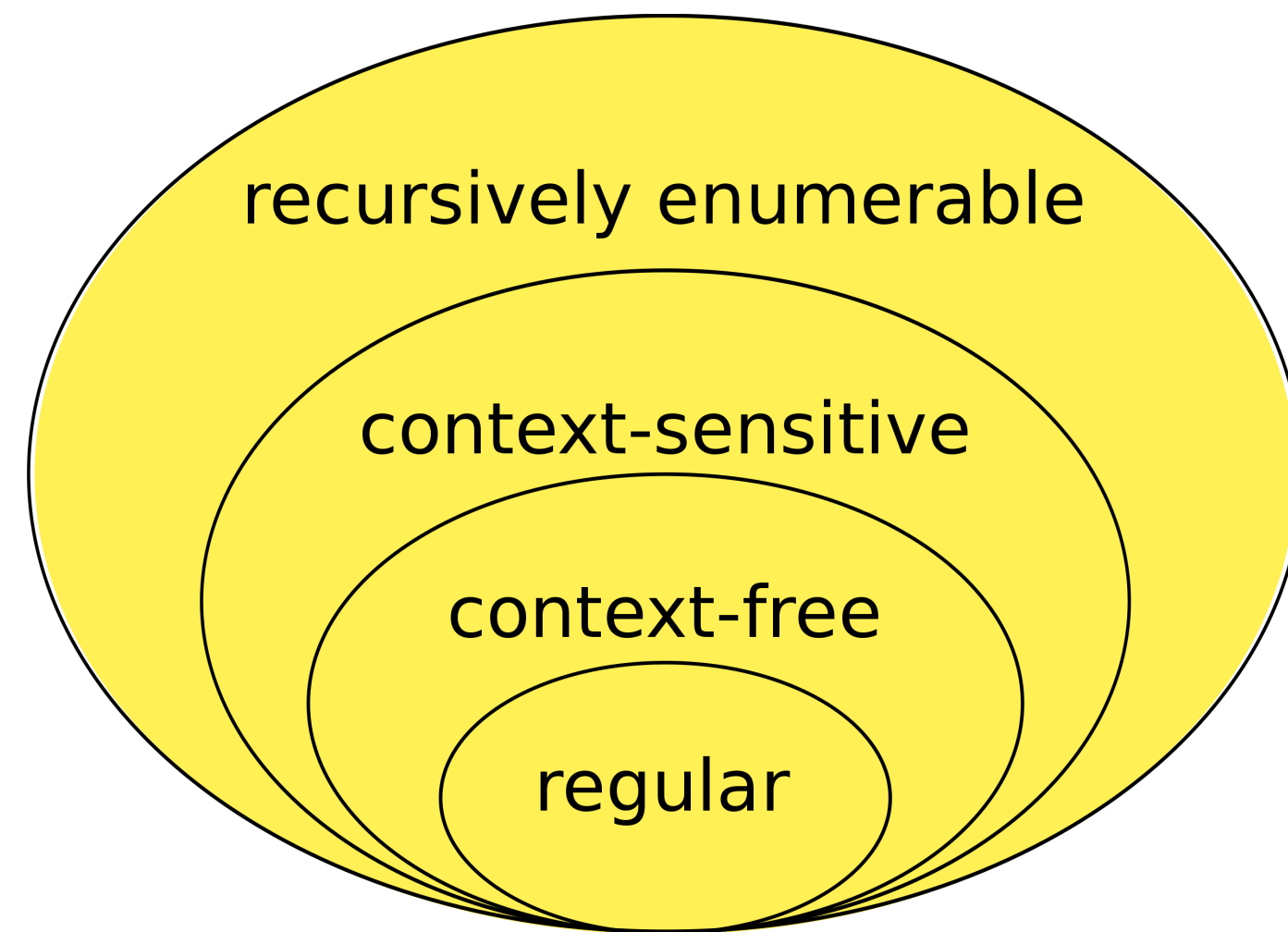
Machine

```
i = 0
while s[i] == 'a':
    i = i + 1
n = i
while i < len(s):
    assert s[i] == 'b'
    i = i + 1
    n = n - 1
assert n == 0
```



The Chomsky hierarchy

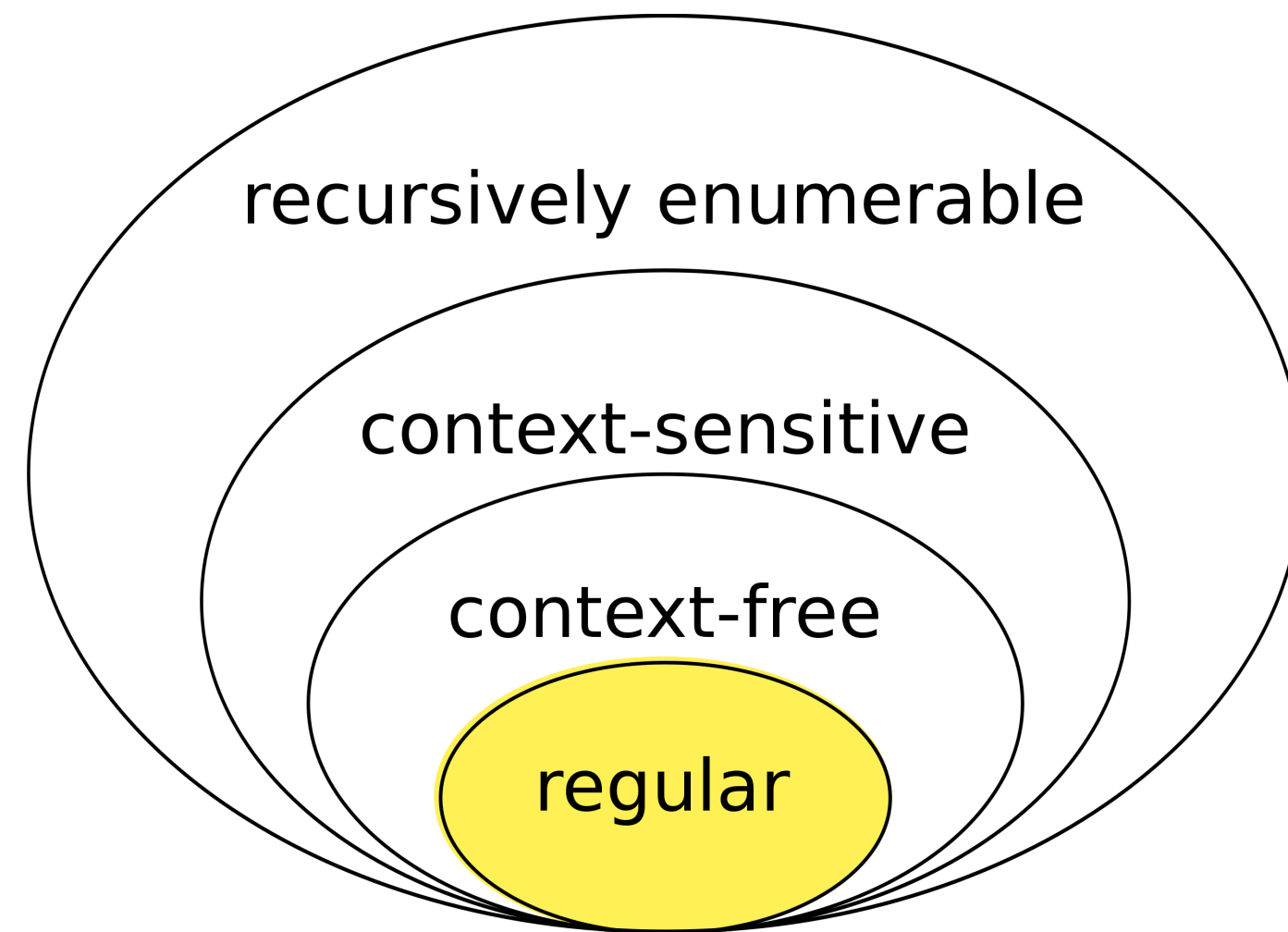
```
01 import math
02
03 while True:
04     line = input("enter 3d vector: ")
05     xs = line.split(",")
06     if len(xs) == 3: break
07     else: print("please try again.")
08
09 [x,y,z] = map(int, xs)
10 n = math.sqrt(x**2 + y**2 + z**2)
11 print(f"vector length: {n}")
```



The Chomsky hierarchy

```
01 import math
02
03 while True:
04     line = input("enter 3d vector: ")
05     xs = line.split(",")
06     if len(xs) == 3: break
07     else: print("please try again.")
08
09 [x,y,z] = map(int, xs)
10 n = math.sqrt(x**2 + y**2 + z**2)
11 print(f"vector length: {n}")
```

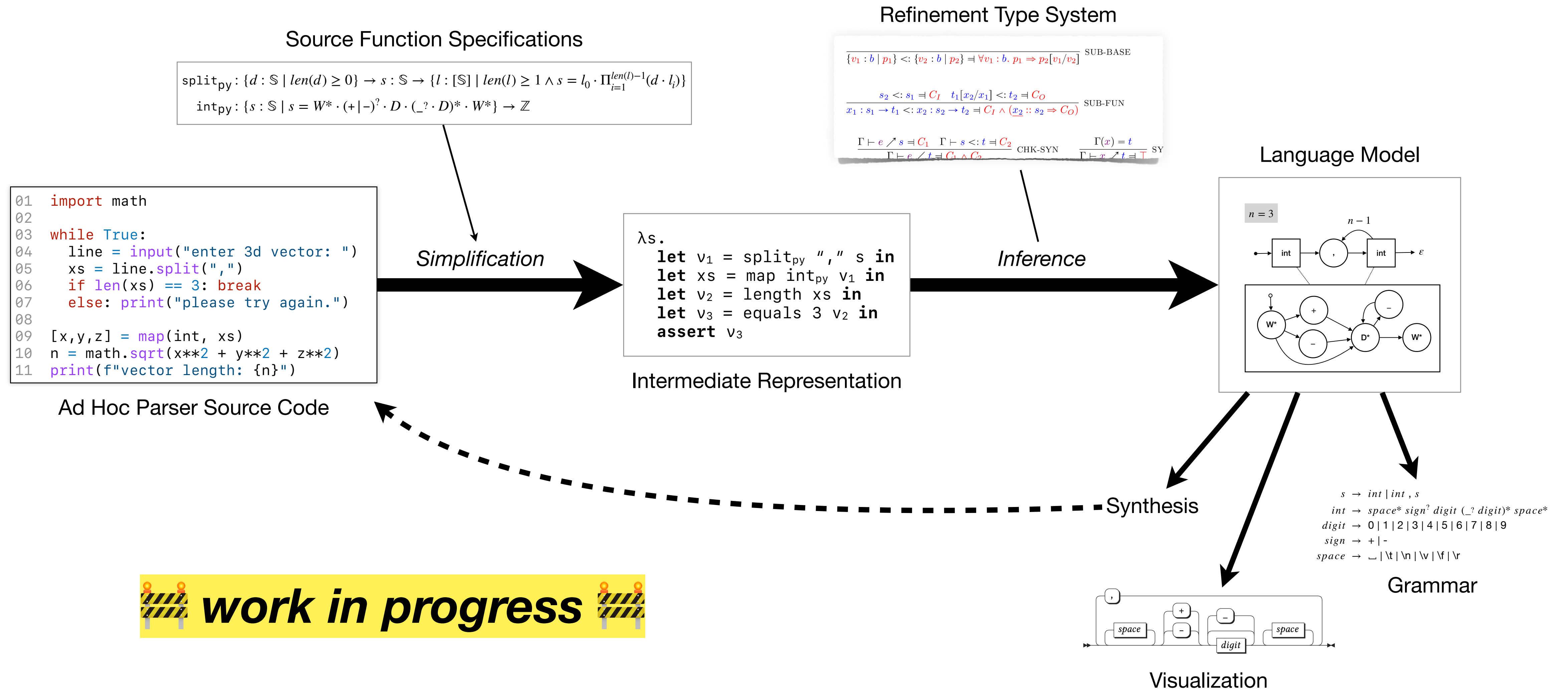
Intuition: Parsers are Embedded Machines



The Chomsky hierarchy

```
01 import math
02
03 while True:
04     line = input("enter 3d vector: ")
05     xs = line.split(",")
06     if len(xs) == 3: break
07     else: print("please try again.")
08
09     [x,y,z] = map(int, xs)
10     n = math.sqrt(x**2 + y**2 + z**2)
11     print(f"vector length: {n}")
```


Vision: Automatic Grammar Inference



New Possibilities

Interactive Documentation



- inferred grammar is always up-to-date
- closely linked to underlying source code

Inferred grammar:

$$\text{pred} \rightarrow a = v_1 \mid a o v_2$$
$$a \rightarrow (\Sigma \backslash o)^*$$
$$v_1 \rightarrow \Sigma^*$$
$$v_2 \rightarrow \text{eval}$$
$$o \rightarrow < \mid <= \mid > \mid >=$$

```
01 def pred(s):
02     ops = r"(=|<|<=|>|>=)"
03     a,o,v = re.split(ops, s)
04     if o[0] == "<" or o[0] == ">":
05         v = eval(v)
06     return p(a,o,v)
```

Inferred grammar:

$$\text{pred} \rightarrow a = v_1 \mid a o v_2$$
$$a \rightarrow (\Sigma \backslash o)^*$$
$$v_1 \rightarrow \Sigma^*$$
$$v_2 \rightarrow \text{eval}$$
$$o \rightarrow < \mid <= \mid > \mid >=$$

```
01 def pred(s):
02     ops = r"(=|<|<=|>|>=)"
03     a,o,v = re.split(ops, s)
04     if o[0] == "<" or o[0] == ">":
05         v = eval(v)
06     return p(a,o,v)
```

Inferred grammar:

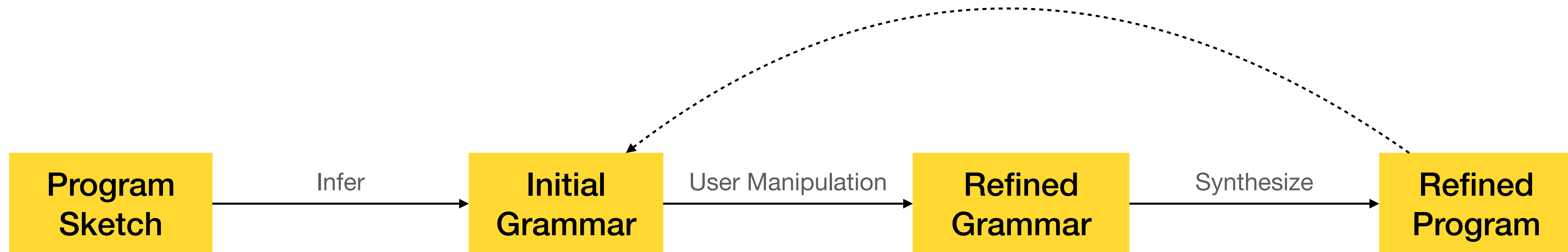
$$\text{pred} \rightarrow a = v_1 \mid a o v_2$$
$$a \rightarrow (\Sigma \backslash o)^*$$
$$v_1 \rightarrow \Sigma^*$$
$$v_2 \rightarrow \text{eval}$$
$$o \rightarrow < \mid <= \mid > \mid >=$$

```
01 def pred(s):
02     ops = r"(=|<|<=|>|>=)"
03     a,o,v = re.split(ops, s)
04     if o[0] == "<" or o[0] == ">":
05         v = eval(v)
06     return p(a,o,v)
```

Bi-Directional Parser Synthesis



- semantic program transformation via inferred grammar
→ program sketching



```
01 def foo(s):  
02     i = 0  
03     while s[i] == "a":  
04         i += 1
```

$$\text{foo} \rightarrow a^*(\Sigma \setminus a)\Sigma^*$$

$$\text{foo} \rightarrow a^*$$

```
01 def foo(s):  
02     i = 0  
03     while i < len(s):  
04         assert s[i] == "a"
```

Mining & Learning



- grammar as *equivalence class* over concrete implementations
 - detecting code clones of ad hoc parsers
 - grammar-enhanced semantic code search

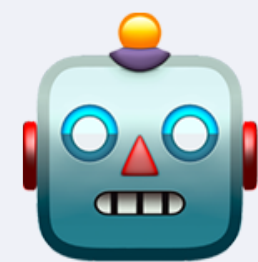
```
🔍 s -> int | int , s
```

Viewer.py	2 matches Python
<pre>640 ranges = self.find_ranges() 641 split = str.split 642 point = map(int, split(self.text.index(CURRENT), ',')) 643 for start, end in ranges: 644 startv = map(int, split(start, ','))</pre>	
transformScriptTags.ts	1 match TypeScript
<pre>122 return null; 123 } 124 return rawValue.split(",").map(item => parseInt(item)); 125 } 126</pre>	

Mining & Learning



- learn how implicit input specifications evolve over time
 - grammar-aware semantic change tracking



⚠ Merging #420 (6a36b23) into main (224b18b) will change the input grammar of a function.

Before:




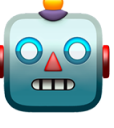

$\text{baz} \rightarrow a*b$

After:

$\text{baz} \rightarrow \Sigma a*b$

```
18 18 def baz(s):
19 - i = 0
19 + i = 1
20 20 while s[i] == "a"
21 21     i += 1
22 22 assert s[i] == "b"
```

Future Plans

-  mining study of ad hoc parsers in the wild ongoing
-  grammar inference via refinement types of simplified parsing IR ongoing
- $\forall x$ proving soundness of simplification and inference steps tbd
-  evaluation on corpus of curated ad hoc parser samples tbd
-  large-scale mining study of inferred grammars tbd
-  user studies on grammar comprehension tbd

Grammars for Free

Toward Grammar Inference for Ad Hoc Parsers

<https://arxiv.org/abs/2202.01021>



Michael Schröder

TU Wien

Vienna, Austria

michael.schroeder@tuwien.ac.at



Jürgen Cito

TU Wien and Meta Platforms, Inc.

Vienna, Austria

juergen.cito@tuwien.ac.at

New Ideas and Emerging Results, ICSE 2022
Pittsburgh, PA, USA



Informatics