



Python Lex and Yacc

Prepared by

Zeyad moneer abo El-Zahab

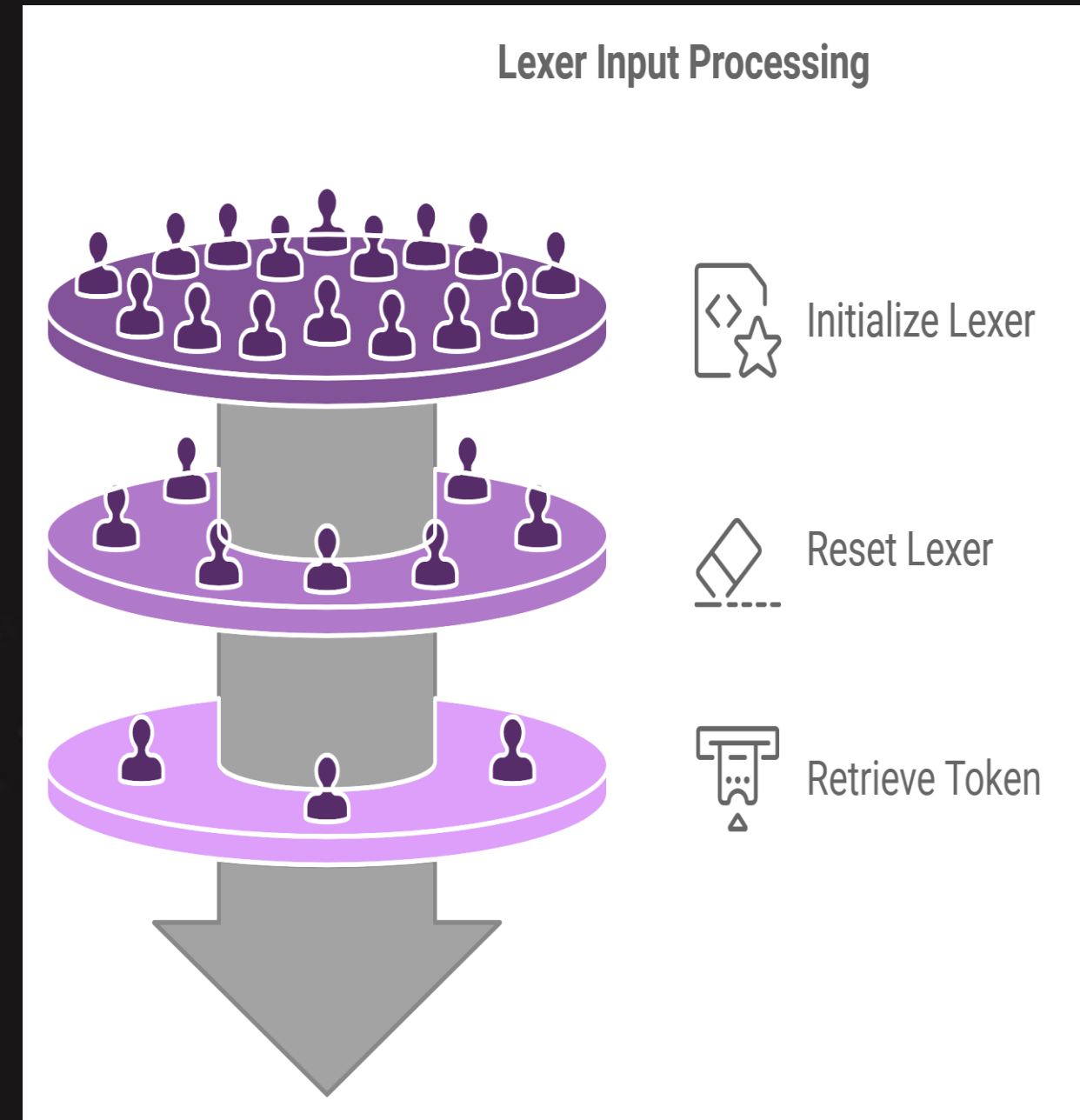
Agenda

- ❑ Building and using the lexer
- ❑ The @TOKEN decorator
- ❑ Optimized mode
- ❑ Debugging
- ❑ Alternative specification of lexers
- ❑ Lexer cloning
- ❑ Internal lexer state
- ❑ Conditional lexing and start conditions
- ❑ How to use conditional analysis ?



Building and using the lexer

- This function uses Python reflection to read the regular expression rules out of the calling context and build the lexer. Once the lexer has been built, two methods can be used to control the lexer.
 - `lexer.input(data)`. Reset the lexer and store a new input string.
 - `lexer.token()`. Return the next token. Returns special `LexToken` instance on success or `None` if the end of the input text has been reached.



Building and using the lexer

```
LexToken(WORD, 'hello', 1, 0)  
LexToken(NUMBER, '123', 1, 6)  
LexToken(WORD, 'world', 1, 10)  
LexToken(NUMBER, '456', 1, 16)
```

```
lex.py > ...  
1  import ply.lex as lex  
2  
3  tokens = (  
4      'NUMBER',  
5      'WORD',  
6  )  
7  
8  t_WORD = r'[a-zA-Z]+'  
9  t_NUMBER = r'\d+'  
10  
11 t_ignore = ' \t'  
12  
13 def t_error(t):  
14     print(f"illegal char :{t.value[0]}")  
15     t.lexer.skip(1)  
16  
17 lexer = lex.lex()  
18  
19 data = "hello 123 world 456"  
20  
21 lexer.input(data)  
22  
23 while True:  
24     tok = lexer.token()  
25     if not tok:  
26         break  
27     print(tok)  
28
```

The @TOKEN decorator

- in some applications, you may want to define build tokens from as a series of more complex regular expression rules.

The '@TOKEN' decorator

provides a more convenient way to define tokens compared to traditional methods.

Improved Readability

Clearer and more concise code, making it easier to understand the lexer's structure.

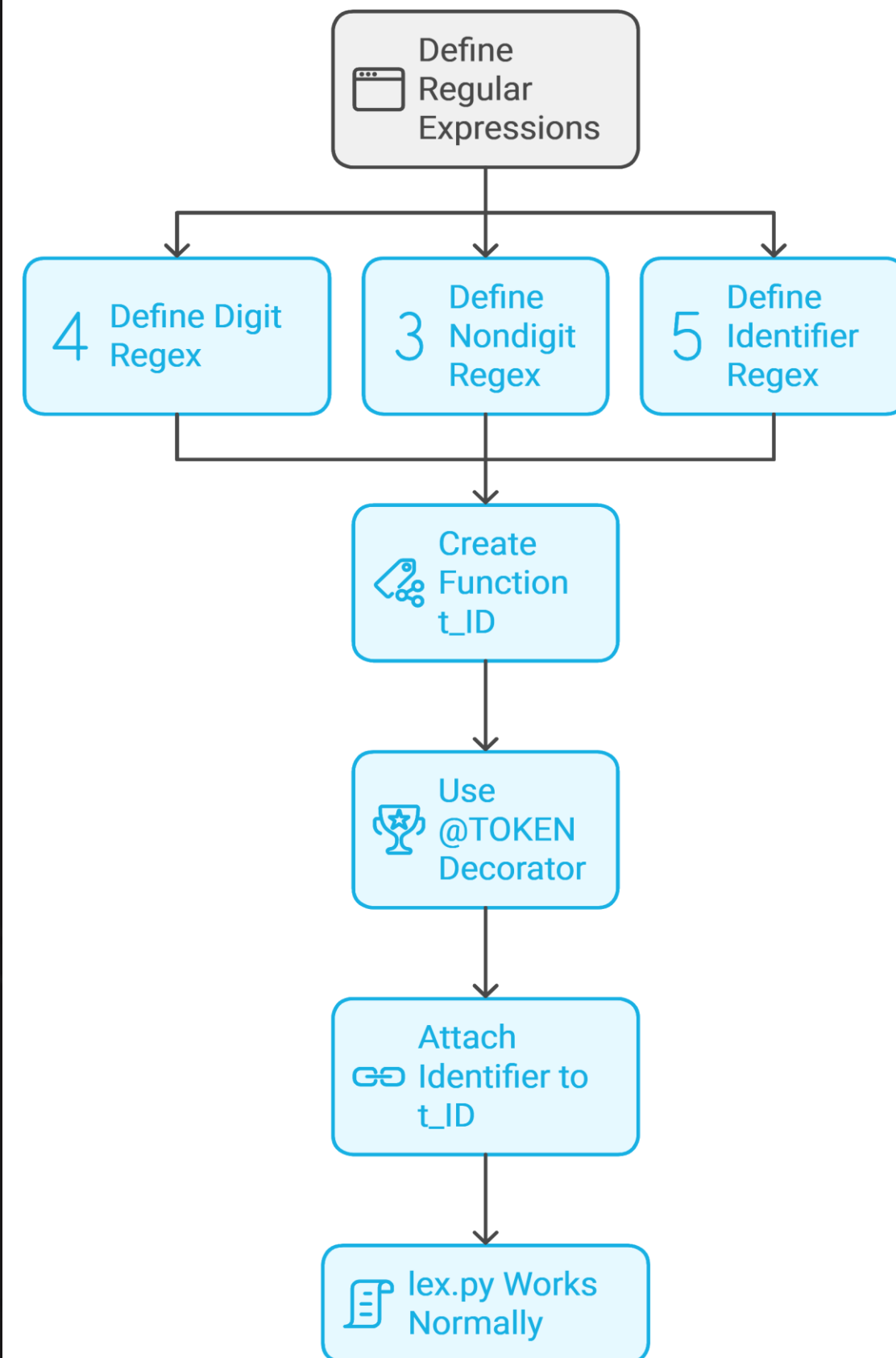
Reduced Boilerplate

The decorator takes care of common lexer setup, reducing the amount of code you need to write.

The @TOKEN decorator

@token.py > ...

```
1  from ply.lex import TOKEN
2
3  digit = r'([0-9])'
4  nondigit = r'([_A-Za-z])'
5  identifier = r'(' + nondigit + r'(' + digit + r'|' + nondigit + r')*)'
6
7  @TOKEN(identifier)
8  def t_ID(t):
9      print("it's identifier : ", t.value)
10     return t
11
12     # -----
13     # Explain the difference
14
15     import ply.lex as lex
16
17     def t_ID(t):
18         r'([_A-Za-z])([_A-Za-z0-9])*'
19         print("it's identifier : ", t.value)
20         return t
21
```



Optimized mode

- **Purpose of Optimizations:**







Reducing memory usage and processing time during lexing.

- **Enabling Optimized Mode:**

Setting `optimize=1` in the `lex()` function.

- **Implications:**

Trade-off between detailed error reporting and performance gains in optimized mode.

Pros	VS	Cons
 Improved performance		 Ignores documentation strings
 Reduced startup time		 Disables error checking
 Works in optimized mode		 Requires confidence in code

Optimized mode

optimized mode.py > ...

```
1 import ply.lex as lex
2
3 tokens = ('PRINT', 'IF', 'WHILE')
4
5 t_PRINT = r'print'
6 t_IF = r'if'
7 t_WHILE = r'while'
```

```
9 # optimize mode
10 lexer = lex.lex(optimize=1)
11
12 # change the file name, It is automatically imported on subsequent runs.
13 lexer = lex.lex(optimize=1, lextab="footab")
```

> .venv

@token.py

conditional 2.py

Conditional lexing and...

conditional.py

lex.py

lextab.py

optimized mode.py

@token.py

conditional 2.py

Conditional lexing and...

conditional.py

footab.py

lex.py

lextab.py

optimized mode.py

Debugging

1

Enabling Debugging

Set (debug=1) to track lexer behavior.

2

Print Statements

Insert print statements in your lexer rules to track the tokenization process.

3

Logging

Utilize a logging framework to record token information for later analysis.

**Initiate Lexer
Creation**



**Set Debug
Mode to 1**



**Create Lexer
with
Debugging**



**Process
Complete**

Debugging

lexer.log

```
1  DEBUG:root:Token NUMBER: 3 at position 0
2  INFO:root:Token: LexToken(NUMBER,3,1,0)
3  ERROR:root:Illegal character ' ' at position 1
4  INFO:root:Token: LexToken(PLUS,'+',1,2)
5  ERROR:root:Illegal character ' ' at position 3
6  DEBUG:root:Token NUMBER: 5 at position 4
7  INFO:root:Token: LexToken(NUMBER,5,1,4)
8  ERROR:root:Illegal character ' ' at position 5
9  INFO:root:Token: LexToken(TIMES,'*',1,6)
10 ERROR:root:Illegal character ' ' at position 7
11 DEBUG:root:Token NUMBER: 10 at position 8
12 INFO:root:Token: LexToken(NUMBER,10,1,8)
13 ERROR:root:Illegal character ' ' at position 10
14 INFO:root:Token: LexToken(MINUS,'-',1,11)
15 ERROR:root:Illegal character ' ' at position 12
16 DEBUG:root:Token NUMBER: 4 at position 13
17 INFO:root:Token: LexToken(NUMBER,4,1,13)
18 ERROR:root:Illegal character ' ' at position 14
19 INFO:root:Token: LexToken(DIVIDE,'/',1,15)
20 ERROR:root:Illegal character ' ' at position 16
21 DEBUG:root:Token NUMBER: 2 at position 17
22 INFO:root:Token: LexToken(NUMBER,2,1,17)
23
```

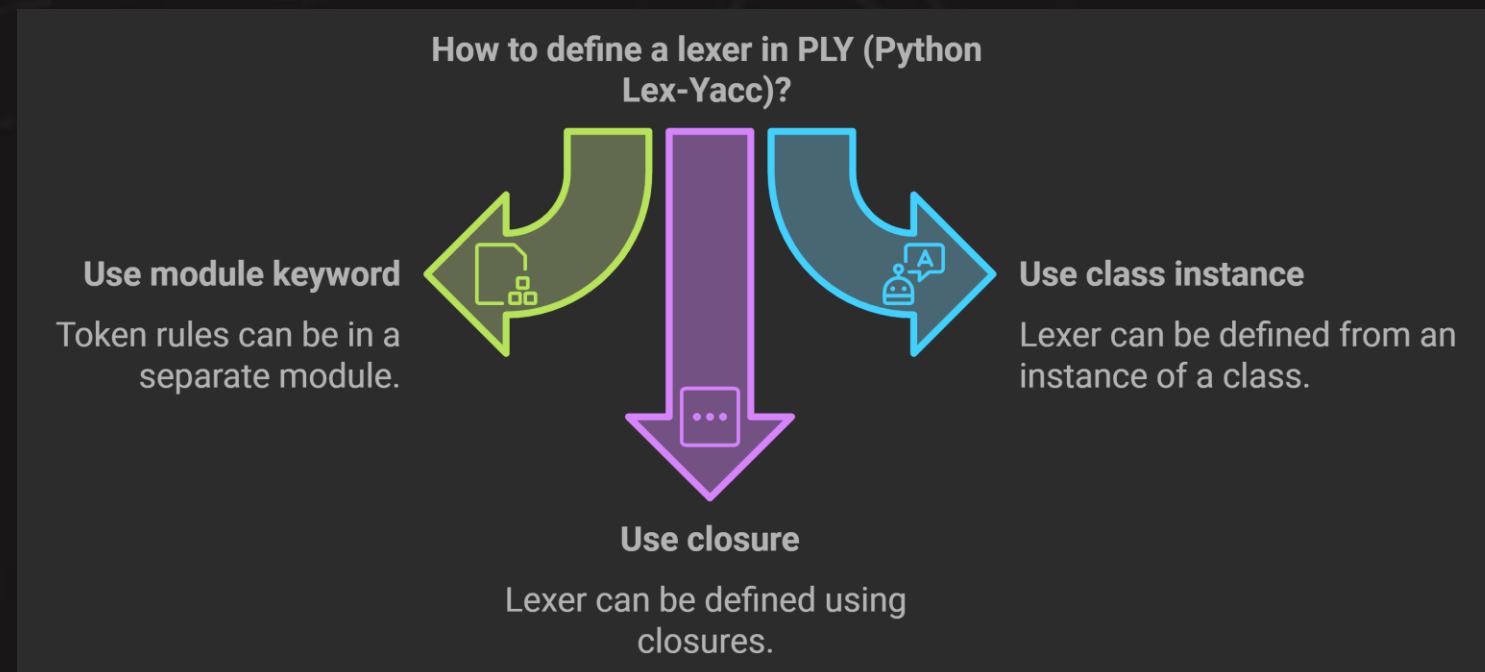
debug.py > ...

```
1  import ply.lex as lex
2  import logging
3
4  # logging
5  logging.basicConfig(filename='lexer.log', level=logging.DEBUG)
6
7  tokens = ('NUMBER', 'PLUS', 'MINUS', 'TIMES', 'DIVIDE')
8
9  t_PLUS = r'\+'
10 t_MINUS = r'\-'
11 t_TIMES = r'\*'
12 t_DIVIDE = r'\/'
13
14 def t_NUMBER(t):
15     r'\d+'
16     t.value = int(t.value)
17     logging.debug(f"Token NUMBER: {t.value} at position {t.lexpos}")
18     return t
19
20 def t_error(t):
21     logging.error(f"Illegal character '{t.value[0]}' at position {t.lexpos}")
22     t.lexer.skip(1)
23
24 # Debug mode
25 lexer = lex.lex(debug=1)
26
27 data = "3 + 5 * 10 - 4 / 2"
28 lexer.input(data)
29
30 for token in lexer:
31     logging.info(f"Token: {token}")
32
```

lex.py
lexer.log
optimized mode.py

Alternative Specification of Lexers

- Alternative specification of lexers refers to different methods or approaches used to define and implement lexers, Unlike traditional lexer generators that use regular expressions, alternative specifications may utilize different frameworks(class , module , closure, etc.)
- purpose :** These alternatives can improve flexibility, error handling, and performance depending on the specific needs of the programming language or application."



Alternative Specification of Lexers

```
def MyLexer():  
    # Regular expression rules for simple tokens  
    t_PLUS = r'\+'  
    t_MINUS = r'\-'  
    t_TIMES = r'\*'  
    t_DIVIDE = r'\/'  
    t_LPAREN = r'\('  
    t_RPAREN = r'\)'
```

```
class MyLexer(object):  
    # List of token names. This is always required  
    tokens = (  
        'NUMBER',  
        'PLUS',  
        'MINUS',  
        'TIMES',  
        'DIVIDE',  
        'LPAREN',  
        'RPAREN',  
    )
```

- lex.py
- optimized mode.py
- tokrules.py

```
ALT.py > ...  
1 import ply.lex as lex  
2 import tokrules  
3 lexer = lex.lex(module=tokrules)  
4 lexer.input("3 + 4 * 1")  
5 while True:  
6     tok = lexer.token()  
7     if not tok:  
8         break  
9     print(tok)
```

Lexer cloning

- Duplicate Lexers

- ✓ Create multiple copies of the original lexical analyzer object so that these copies can process multiple text streams independently.

- Independent State

- ✓ Duplicate copies of the lexical analyzer maintain their own internal state, preventing information or data from interfering between different copies.

- Parallel Processing

- ✓ The cloning feature can be used to perform parallel processing operations, so that multiple inputs are analyzed simultaneously using cloned copies.

Lexer cloning

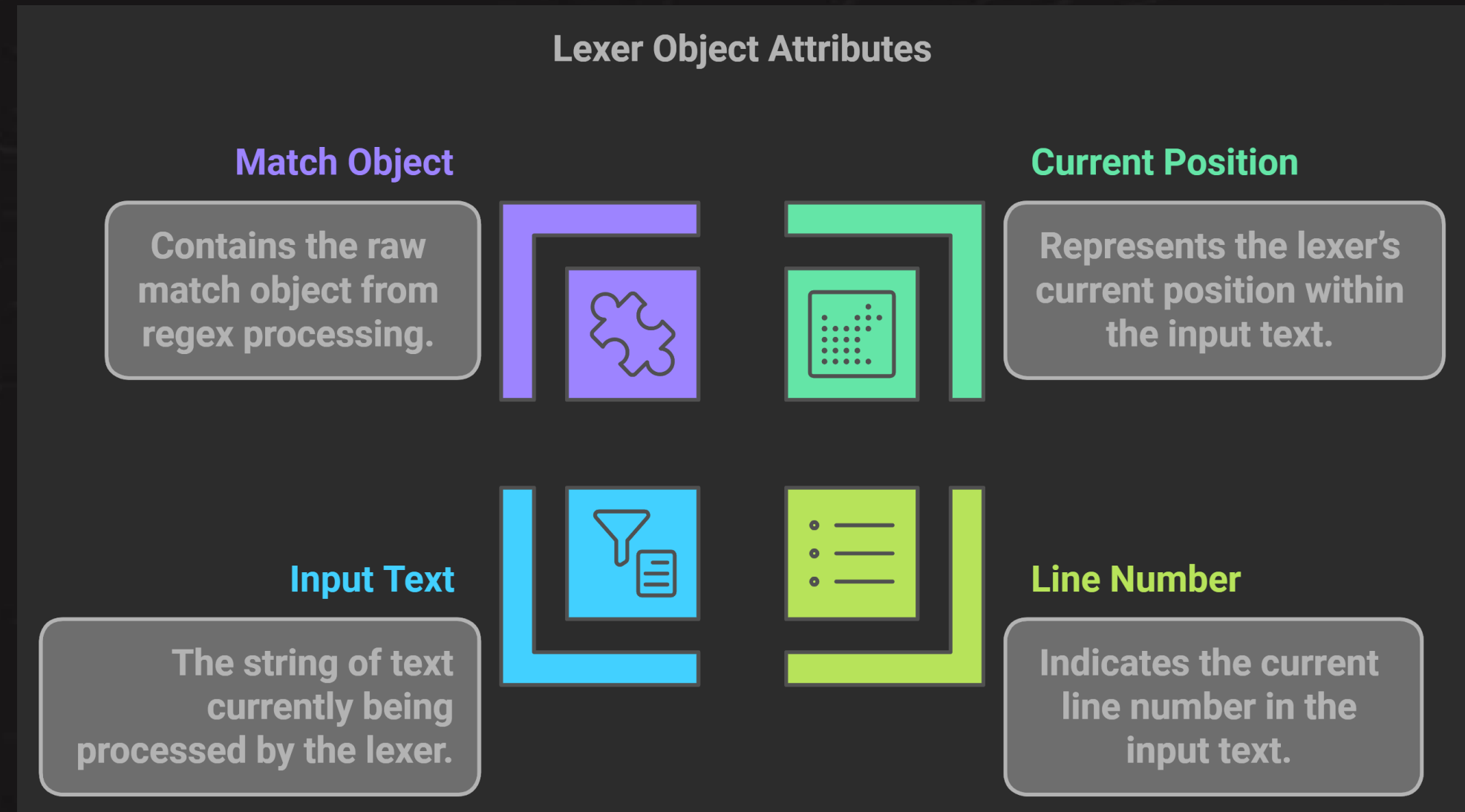
Number of numbers in the expression: 3
The expression has an odd number of numbers.

```
clone.py > ...
28 |
29 |
30 | lexer.input("3 + 4 * 5")
31 |
32 | # lexer clone
33 | lookahead_lexer = lexer.clone()
34 |
35 | # Analyze a portion of text using cloning
36 | # Read all numbers using the cloned version
37 | count_numbers = 0
38 | while True:
39 |     tok = lookahead_lexer.token()
40 |     if not tok:
41 |         break
42 |     if tok.type == 'NUMBER':
43 |         count_numbers += 1
44 |
45 | print("Number of numbers in the expression:", count_numbers)
46 | if count_numbers % 2 == 1:
47 |     print("The expression has an odd number of numbers.")
48 | else:
49 |     print("The expression has an even number of numbers.")
50 |
```

Internal lexer state

- A Lexer object `lexer` has a number of internal attributes that may be useful in certain situations.

- ✓ `lexer.lexpos`
- ✓ `lexer.lineno`
- ✓ `lexer.lexdata`
- ✓ `lexer.lexmatch`



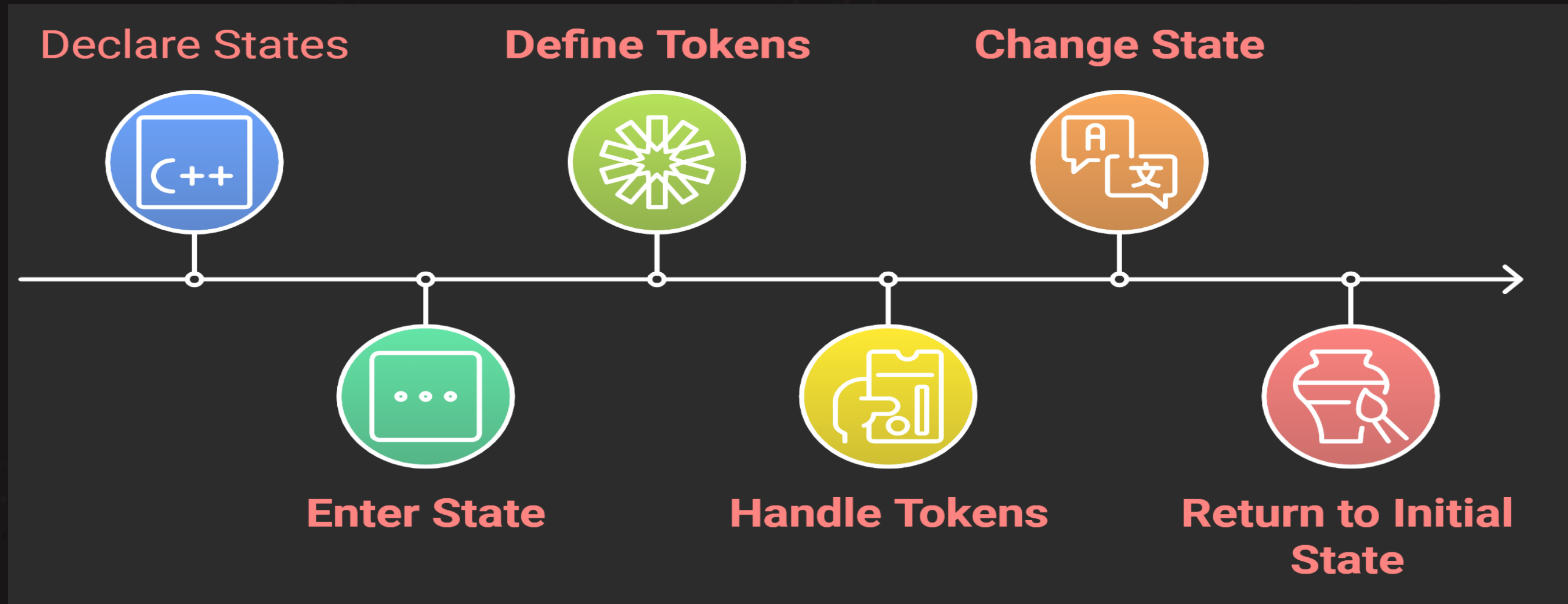
Conditional lexing and start conditions

- In advanced parsing, you may encounter situations where you need to deal with multiple grammar rules within the same text. For example, **In a script that contains comments or special sections enclosed in parentheses, you need dedicated parsing rules to deal with these sections.** To facilitate this, PLY provides the ability to define different parsing cases, so that you can switch between them as needed.
- **Parsing cases can be of two types:**
 - ✓ **Exclusive case:** If the case is exclusive, the default parsing rules will be disabled, and parsing will be limited to the rules defined for that case only.
 - ✓ **Inclusive case:** In this case, the rules defined for that case will be combined with the default rules.

How to use conditional analysis ?

- In order to implement conditional analysis in a symbol parser using the PLY library, states are defined that are switched between based on specific symbols encountered in the text. The default state is **INITIAL**, which is where the parser starts, and can then be switched to other states based on the needs of the analysis.
- **Working steps:**
 - ✓ Define states: Different states are defined in the parser using the states list.
 - ✓ Switch between states: The `begin()`, `push_state()`, and `pop_state()` functions are used to change the state during the analysis.
 - ✓ Return state: You can return to the default state (or any other state) based on the analysis

Conditional lexing and start conditions



Thank you

Any questions?