

Creating a lexical analyser

Marcin Kuta

Theory of Compilation
Laboratory 1

Basic concepts

Basic concepts

Pattern [0-9] +

Token INTNUM

Lexem 1920

Example of specification

Example of specification

```
a      { Action1 }
abb    { Action2 }
a*b+  { Action3 }
```

Input: abb

Possible tokenizations

```
a|bb|  { Action1, Action3 }
abb|   { Action2 }
abb|   { Action3 }
```

Rules of lexical specification

Two rules

- ① Principle of maximal match
- ② Detailed specifications before general specification
 - If an input string matches two patterns, the pattern which appears earlier in the specification list is chosen

Example of specification

```
a      { Action1 }
abb    { Action2 }
a*b+  { Action3 }
```

Input: abb

Possible tokenizations

```
a|bb| Action1, Action3
abb| Action2
abb| Action3
```

Example of specification

```
a      { Action1 }
abb    { Action2 }
a*b+  { Action3 }
```

Input: abb

Possible tokenizations

a|bb| Action1, Action3

abb| Action2

abb| Action3

Example of specification

```
a      { Action1 }
abb    { Action2 }
a*b+  { Action3 }
```

Input: abb

Possible tokenizations

~~a|bb| Action1, Action3~~

~~abb| Action2~~

~~abb| Action3~~

Scanner specification

In practice, we distinguish three types of tokens:

- literals
- reserved keywords
- general tokens

Scanner specification in SLY or PLY

Literals:

- Lexems are one-character
- Token can be represented by its one-character lexem

```
literals = [ '+', '-' , '*' , '/' ]  
literals = "+-*/"
```

Scanner specification in PLY

General tokens

- One token matches many lexems
- Specified with regular expressions

Examples:

NUM - matches many numbers

```
def t_NUM(t):
    r"\d+"
    return t
```

ID - matches many identifiers

```
def t_ID(self,t):
    r"[a-zA-Z_]\w*"
    return t
```

Scanner specification in PLY

Reserved keywords:

- One token corresponds to exactly one lexem
- Lexems are longer than one character
- Their specification matches also specification of an identifier, so they should appear earlier on the specification list

```
reserved = {    'break':           'BREAK',
                 'continue':        'CONTINUE',
                 'if':              'IF',
                 'else':            'ELSE',
             }
tokens = [ "ID", "EQ", "NEQ", "LE", "GE" ] + list(
    reserved.values())

def t_ID(t):
    r"[a-zA-Z_]\w*"
    t.type = reserved.get(t.value, 'ID')
    return t
```

Scanner specification in PLY

Pattern to be avoided - individual rules for reserved keywords:

```
t_BREAK = r'break'  
t_CONTINUE = r'continue'  
t_IF = r'if'  
t_ELSE = r'else'
```

Scanner specification in SLY

```
tokens = [ "ID", "EQ", "NEQ", "LE", "GE" ] + list(  
    reserved.values())  
  
ID = r'[a-zA-Z_][a-zA-Z0-9_]*'  
  
ID['break'] = 'BREAK'  
ID['continue'] = 'CONTINUE'  
ID['if'] = 'IF'  
ID['else'] = 'ELSE'
```

References

- ① <https://sly.readthedocs.io/en/latest/sly.html>, Sect. Writing a lexer
- ② <https://github.com/dabeaz/sly>
- ③ <http://www.dabeaz.com/ply/ply.html>, Sect. 4, Lex