

## **Tokenization-and-Categorization-using-OCaml-A-Lexical-Analysis-Approach**

### **Introduction**

The most important phase in the compilation process is lexical analysis, which divides source code into tokens that the compiler or interpreter can utilize for additional processing. This work provides an OCaml-based Lexical Analyzer that can identify and classify tokens, including identifiers, operators, punctuation symbols, integer literals, and keywords.

### **Process steps:**

**Tokenization:** Tokenization involves dividing the input text into discrete tokens, which include operators, keywords, punctuation, integer literals, identifiers, and unknown tokens.

**Categorization:** Based on its type, each token can be categorized as either a keyword, operator, punctuation symbol, integer literal, identifier, or unknown.

**Printing:** The tokens are organized and printed to the standard output, where each token is shown with its associated given category.

### **Explanation:**

A token type which represents different kinds of tokens identified during lexical analysis is specified. Keywords (Keyword), operators (Operator), integer literals (IntLiteral), punctuation symbols (Punctuation), identifiers (Identifier), and unknown tokens (Unknown) are the given types.

The `is_alnum` function can be useful in identifying keywords and IDs since it detects if a given character is alphanumeric. In comparison, a character's status as a punctuation sign is identified by the `is_punctuation` function.

The lexical analyzer's key component is the `tokenize` function. It receives a string as input and separates it into tokens in accordance with the given criteria. The function handle operators, punctuation symbols, whitespace, and alphanumeric characters using secondary functions. It also attempts to use exception handling to convert alphanumeric substrings into integer literals.

The categorized tokens are printed to the standard output by the `print_tokens` function. It prints each token along with its category following looping over the list of tokens.

The main function asks the user to enter a string to be tokenized, reads the input string, utilizes the `tokenize` function to tokenize it, and then uses the `print_tokens` function to print the tokens that are generated.

## Code:

```

1  (* Lexical Analyzer in OCaml *)
2
3  (* Define types for tokens *)
4  type token =
5    | Keyword of string
6    | Operator of string
7    | Punctuation of string
8    | Intliteral of int
9    | Identifier of string
10   | Unknown of string
11
12  (* Function to check if a character is alphanumeric *)
13  let is_alnum c =
14    (c >= 'a' && c <= 'z') ||
15    (c >= 'A' && c <= 'Z') ||
16    (c >= '0' && c <= '9')
17
18  (* Function to check if a character is a punctuation symbol *)
19  let is_punctuation c =
20    List.mem c ['('; ')'; '{'; '}'; ';'; ',']
21
22  (* Function to tokenize a string *)
23  let rec tokenize str =
24    let rec consume_whitespace i =
25      if i < String.length str && (str.[i] = ' ' || str.[i] = '\t') then
26        consume_whitespace (i + 1)
27      else
28        i
29    in
30    let rec tokenize_helper i =
31      if i >= String.length str then
32        []
33      else
34        let next_char = str.[i] in
35        if next_char = ' ' || next_char = '\t' then
36          tokenize_helper (consume_whitespace i)
37        else if is_alnum next_char then
38          let j = ref i in
39          while !j < String.length str && is_alnum str.[!j] do
40            incr j
41          done;
42          let token_str = String.sub str i (j - i) in
43          if token_str = "if" || token_str = "else" || token_str = "for" || token_str = "if else" || token_str = "while" || token_str = "let" || token_str = "in" || token_str = "then" then
44            Keyword token_str :: tokenize_helper !j
45          else
46            begin
47              try
48                Intliteral (int_of_string token_str) :: tokenize_helper !j
49              with Failure _ ->
50                Identifier token_str :: tokenize_helper !j
51            end
52        else if is_punctuation next_char then
53          Punctuation (Char.escaped next_char) :: tokenize_helper (i + 1)
54        else if next_char = '+' || next_char = '-' || next_char = '*' || next_char = '/' || next_char = '=' || next_char = '!' then
55          let lookahead = if i + 1 < String.length str then Some str.[i + 1] else None in
56          (match lookahead with
57           | Some '-' ->
58             Operator (Char.escaped next_char ^ Char.escaped '-') :: tokenize_helper (i + 2)
59           | _ ->
60             Operator (Char.escaped next_char) :: tokenize_helper (i + 1))
61        else
62          Unknown (Char.escaped next_char) :: tokenize_helper (i + 1)
63    in
64    tokenize_helper 0
65
66  (* Function to print tokens *)
67  let print_tokens tokens =
68    List.iter (fun token ->
69      match token with
70      | Keyword s -> Printf.printf "Keyword: %s\n" s
71      | Operator s -> Printf.printf "Operator: %s\n" s
72      | Punctuation s -> Printf.printf "Punctuation: %s\n" s
73      | Intliteral i -> Printf.printf "Intliteral: %d\n" i
74      | Identifier s -> Printf.printf "Identifier: %s\n" s
75      | Unknown s -> Printf.printf "Unknown: %s\n" s
76    ) tokens
77
78  (* Main function *)
79  let () =
80    Printf.printf "Enter a string to tokenize:\n";
81    let input_str = read_line () in
82    let tokens = tokenize input_str in
83    print_tokens tokens

```

### Test case 1:

```
camel@sm3094: ~  
camel@sm3094:~$ ./Lexical_Analyser  
Enter a string to tokenize:  
if x==35 then {a=b+9;} else {b=a-1;}  
Keyword: if  
Identifier: x  
Operator: ==  
IntLiteral: 35  
Keyword: then  
Punctuation: {  
Identifier: a  
Operator: =  
Identifier: b  
Operator: +  
IntLiteral: 9  
Punctuation: ;  
Punctuation: }  
Keyword: else  
Punctuation: {  
Identifier: b  
Operator: =  
Identifier: a  
Operator: -  
IntLiteral: 1  
Punctuation: ;  
Punctuation: }
```

### Test case 2:

```
camel@sm3094:~$ ./Lexical_Analyser  
Enter a string to tokenize:  
$$$ invalid token  
Unknown: $  
Unknown: $  
Unknown: $  
Identifier: invalid  
Identifier: token
```

### Random Test case:

```
camel@sm3094: ~  
camel@sm3094:~$ ./Lexical_Analyser  
Enter a string to tokenize:  
let x=42;  
Keyword: let  
Identifier: x  
Operator: =  
IntLiteral: 42  
Punctuation: ;
```

```
camel@sm3094:~$ ./Lexical_Analyser
Enter a string to tokenize:
if x==10 then {y=x+5;}else{y=x-5;}
Keyword: if
Identifier: x
Operator: ==
IntLiteral: 10
Keyword: then
Punctuation: {
Identifier: y
Operator: =
Identifier: x
Operator: +
IntLiteral: 5
Punctuation: ;
Punctuation: }
Keyword: else
Punctuation: {
Identifier: y
Operator: =
Identifier: x
Operator: -
IntLiteral: 5
Punctuation: ;
Punctuation: }
```

### References:

<https://v2.ocaml.org/manual/lex yacc.html>

[https://caml.inria.fr/pub/old\\_caml\\_site/ocaml/htmlman/manual026.html](https://caml.inria.fr/pub/old_caml_site/ocaml/htmlman/manual026.html)