

TP Analyse lexicale.

L'objectif de cette séance est de réaliser un analyseur lexical¹ pour une partie du langage Java² qui sera utilisé dans l'Unité d'Enseignement de Technologies Objets (UE TOB) au second semestre. Cette partie du langage sera étudiée dans une future séance de TD sur les grammaires. La partie lexicale du langage est spécifiée par le document <http://docs.oracle.com/javase/specs/jls/se8/html/jls-3.html>.

1 Analyseur lexical

Un analyseur lexical est la première étape de l'analyse d'un langage. Celle-ci est suivie par l'analyse syntaxique³ qui sera l'objet des futures séances de TP et l'analyse sémantique⁴ qui sera l'objet de l'UE Sémantique et Traduction des Langages de la majeure Informatique en deuxième année.

Un analyseur lexical est spécifié sous la forme d'une séquence d'association entre une expression régulière et une action. L'action est effectuée lorsque l'analyseur lexical reconnaît un mot qui fait parti du langage décrit par l'expression régulière associée. Si plusieurs expressions régulières reconnaissent le même mot, seule l'action associée est la première expression régulière est exécutée. Si plusieurs mots reconnus par des expressions régulières contiennent un même préfixe, alors l'action exécutée est celle associée à l'expression régulière qui reconnaît le mot le plus long.

2 L'outil camllex

Pour réaliser cette séance, nous allons utiliser l'outil `camllex`⁵ qui traduit la spécification d'un analyseur lexical en un programme `caml`⁶ qui implante un analyseur lexical sous la forme d'un automate fini.

3 Etude d'un exemple

L'objectif est d'étudier l'exemple suivant en s'appuyant sur la documentation de l'outil `camllex` : <http://caml.inria.fr/pub/docs/manual-ocaml-4.02-1/lexyacc.html>. Pour compiler le programme, utilisez la commande `make` qui produit l'exécutable `mainJava`. Cet exécutable prend en paramètre le nom du fichier que l'on veut analyser.

4 Construction d'un analyseur lexical

L'objectif de cet exercice est de compléter le fichier fourni en exemple `lexerJava.mll` pour traiter l'ensemble des unités lexicales de Java telles qu'elles sont décrites dans la documentation : <http://docs.oracle.com/javase/specs/jls/se8/html/jls-3.html>.

Vous vous attacherez à traiter les différentes formes de constantes (entiers, réels, caractères) ainsi que les différentes formes de commentaires.

¹http://en.wikipedia.org/wiki/Lexical_analysis

²[http://en.wikipedia.org/wiki/Java_\(programming_language\)](http://en.wikipedia.org/wiki/Java_(programming_language))

³<http://en.wikipedia.org/wiki/Parsing>

⁴<http://en.wikipedia.org/wiki/Semantics>

⁵<http://caml.inria.fr/pub/docs/manual-ocaml-4.02-1/lexyacc.html>

⁶<http://caml.inria.fr/pub/docs/manual-ocaml-4.02-1/>

Listing 1: Fichier `lexerJava.ml`

```
1 {
2
3   open TokenJava
4   exception LexicalError
5
6 }
7
8 let chiffre = ['0'-'9']
9 let alpha = ['a'-'z' 'A'-'Z']
10 let alphanum = alpha | chiffre | '_'
11
12 rule lexer = parse
13   | ['\n' '\t' ' ' ' ']+      { lexer lexbuf }
14   | ", "                      { VIRG }
15   | "int"                     { INT }
16   | "while"                   { TANTIQUE }
17   | chiffre+ as texte        { ENTIER (int_of_string texte) }
18   | eof                       { FIN }
19   | _ as texte               { raise LexicalError }
20
21 {
22
23 }
```

Listing 2: Fichier `mainJava.ml`

```
1 open TokenJava;;
2
3 if (Array.length Sys.argv > 1)
4 then
5   let lexbuf = (Lexing.from_channel (open_in Sys.argv.(1))) in
6   let token = ref (LexerJava.lexer lexbuf) in
7   while ((!token) != FIN) do
8     (printToken (!token));
9     (token := (LexerJava.lexer lexbuf))
10  done
11 else
12  (print_endline "MainJava_fichier")
```