

Analyse lexicale: Outil Flex

© Compilation / M.MOURCHID,
SMI/2021-2022

1

Analyse lexicale: outil flex

- Description de l'outil Flex
- Format d'un programme Flex
- Les motifs
- Fonctionnement des analyseurs Flex
- Les actions prédéfinies

© Compilation / M.MOURCHID,
SMI/2021-2022

2

Ecriture d'un analyseur lexical

Il y a 3 méthodes pour développer un analyseur lexical :

- Méthode manuelle : on écrit à la main un analyseur lexical comme un programme en utilisant un langage de programmation évoluée comme *C*, *C++* ou *Java*.
- Méthode par automate : on construit un automate fini déterministe qui représente le lexique du langage source, puis on le muni d'actions d'analyse lexicale.
- Méthode automatique : on écrit une spécification d'un analyseur lexical à l'aide d'un programme qui génère automatiquement l'analyseur lexical spécifiée. Cette spécification est faite par un ensemble d'expressions régulières enrichies par des actions d'analyse lexicale.

© Compilation / M.MOURCHID,
SMI/2021-2022

3

Flex: Fast Lexer

- **Flex** est un outil de **génération automatique d'analyseurs lexicaux**.
- Un fichier *Flex* contient la description d'un analyseur lexical à générer.
- Cette description est donnée sous la forme d'expressions régulières étendues et du code écrit en langage *C* (ou *C++*).
- *Flex* génère comme résultat un fichier contenant le code *C* du futur analyseur lexical.
- Ce fichier est nommé **lex.yy.c**.
- La compilation de ce fichier par un compilateur *C*, génère finalement le code exécutable de l'analyseur lexical en question.

© Compilation / M.MOURCHID,
SMI/2021-2022

4

Flex: Fast Lexer

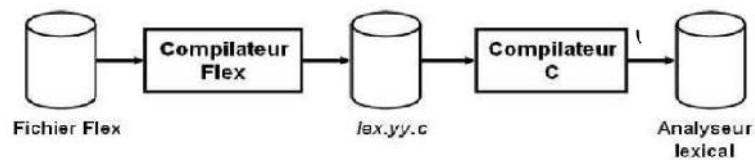


Figure: Processus de création d'un analyseur lexcical avec *Flex*

© Compilation / M.MOURCHID,
SMI/2021-2022

5

Flex: Fast Lexer

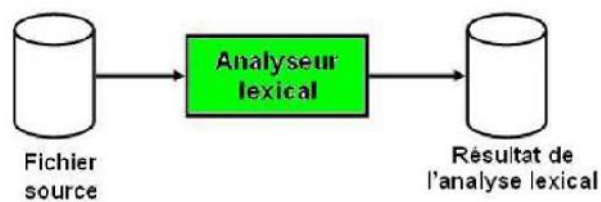


Figure: Utilisation d'un analyseur lexcical généré avec *Flex*

© Compilation / M.MOURCHID,
SMI/2021-2022

6

Flex: Fast Lexer

- Lorsque l'exécutable est mis en œuvre, il analyse le fichier source pour chercher les occurrences d'expressions régulières.
- Lorsqu'une chaîne est reconnue par une expression régulière, il exécute le code C correspondant.

Flex: Fast Lexer

Les étapes à suivre pour obtenir un analyseur avec *Flex* :

1. Ecrire, puis enregistrer votre analyseur dans un fichier portant une extension ".flex" ou ".lex". Par exemple, "scan.flex".
2. Compiler votre analyseur par la commande flex :
flex scan.flex
3. Compiler le fichier par la commande gcc le fichier lex.yy.c produit par l'étape précédente :
gcc -o lex.yy.c scan -lfl
4. Lancer l'analyseur en utilisant le nom de celui-ci : scan

Flex: Fast Lexer

- Par défaut, le texte à analyser est lu à partir de l'entrée standard (le clavier) et le résultat de l'analyse lexical est affiché sur la sortie standard (l'écran).
- On peut indiquer au compilateur le fichier source à analyser :
scan < input.txt
- On peut indiquer au compilateur d'afficher le résultat de l'analyse dans un fichier :
scan > output.txt
- On peut combiner les deux :
scan < input.txt > output.txt

© Compilation / M.MOURCHID,
SMI/2021-2022

9

Flex: Fast Lexer

Exemple 1:

Ecrire un analyseur lexical qui remplace toute occurrence de lettre "a" par une lettre "z", et inversement.

Solution:

```
%%
"a" printf("z");
"z" printf("a");
```

© Compilation / M.MOURCHID,
SMI/2021-2022

10

Flex: Fast Lexer

Fonctionnement

- Les deux symboles %% marquent le debut de la **section des règles**. Ils doivent être écrits en **première colonne**.
- Le code :
`"a" printf("z")`
 est une **règle**.
- Chaque règle contient deux éléments :
 une expression régulière;
 et du code C.
- Dans le code précédent, l'expression régulière est "a" et le code C est l'instruction `printf("z")`.

© Compilation / M.MOURCHID,
SMI/2021-2022

11

Flex: Fast Lexer

Fonctionnement

- Signification du code : lorsque l'analysateur reconnaît une lettre "a" dans le fichier source, il va la remplacer par une lettre "z".
- On a un résultat similaire avec la seconde règle :
`"z" printf("a")`
 par laquelle l'analyseur remplacera chaque lettre "z" par une lettre "a".

© Compilation / M.MOURCHID,
SMI/2021-2022

12

Flex: Fast Lexer

Remarque:

- Une expression régulière doit toujours commencer en première colonne.
- Il faut laisser au moins un caractère blanc (espace ou tabulation) entre l'expression régulière et le code C associé.
- Le code C peut contenir plusieurs instructions. Les accolades { et } doivent être utilisées si c'est le cas.

© Compilation / M.MOURCHID,
SMI/2021-2022

13

Flex: Fast Lexer

Exemple 2:

Ecrire un analyseur lexical qui calcule le nombre de lignes contenu dans le fichier source.

```
%{
    int nblignes = 0;
}%
%%
\n nb_lignes++;
%%
main() {
    yylex();
    printf("nombre de lignes : %d", nblignes); }
```

© Compilation / M.MOURCHID,
SMI/2021-2022

14

Flex: Fast Lexer

Fonctionnement

- La zone entre les symboles %{ et %} représente la **section des déclarations**.
- Dans cette section, on peut :
 - déclarer des variables globales.
 - déclarer des fonctions globales.
 - inclure des fichiers.
- Après le deuxième %, c'est la section du **code utilisateur** : la fonction **main** et les autres fonctions.
- La règle :
 \n nblignes- - ;
 signifie : à la rencontre d'un caractère de saut de ligne "\n" dans le fichier source, incrémenter la variable nblignes.

© Compilation / M.MOURCHID,
SMI/2021-2022

15

Flex: Fast Lexer

Fonctionnement

- La fonction "main" contient deux instructions :
 - yylex : c'est une routine qui effectue l'appel explicite de l'analyseur lexical.
 - une instruction *printf* par laquelle on affiche le nombre de lignes calculé.

© Compilation / M.MOURCHID,
SMI/2021-2022

16

Flex: Fast Lexer

Format d'un fichier Flex

- Un fichier Flex contient trois sections :
 - La section des **déclarations** et des **définitions**.
 - La section des **règles**.
 - La section du **code utilisateur**.
- Chaque section est séparée de la suivante par une ligne contenant juste les deux symboles `%%`.

Flex: Fast Lexer

Format d'un fichier Flex

Un fichier *Flex* se présente comme suit:

```
%{
/* Déclarations */
}%
/* Définitions */
%%
/* Règles */
%%
/* Code utilisateur */
```

Flex: Fast Lexer

Format d'un fichier Flex

- Une **définition** permet de **nommer** une expression régulière.
- Une définition est un couple formé d'un **identificateur** et d'une **expression régulière**, séparée par au moins un caractère espace (ou tabulation).
- Une expression régulière peut contenir une **référence** à un identificateur déjà défini. Dans ce cas, l'identificateur doit être écrit entre les symboles ” {” et ”}”
- Une définition doit être écrite sur une seule ligne.

Flex: Fast Lexer

Exemple:

- **CHIFFRE** [0 - 9]
- **ENTIER** (+|-)? {**CHIFFRE**}⁺

Flex: Fast Lexer

- Une **règle** Flex est la donnée :
d'une expression régulière **e**.
et d'une action (celle qui sera exécutée lorsqu'une séquence d'entrée est reconnue par **e**).

Exemple:

```
"int" {printf("Mot clé int");}
```

A chaque fois que la chaîne de caractères "int" sera reconnue, le message "Mot clé int" sera affichée sur la sortie standard.

Flex: Fast Lexer

Les règles:

- Une expression régulière spécifie une unité lexicale du langage source.
- *Flex* utilise les **expressions régulières étendues**.
- Dans *Flex*, les expressions régulières sont construites à partir des caractères et d'opérateurs.
- Les opérateurs de *Flex* sont :
" \ [] ~ - ? . * + | () \$ / { } % < >
- Pour utiliser ces opérateurs comme caractères ordinaires, il faut les protéger en les plaçant dans une chaîne entourée de double-quotes (") ou en les plaçant après un caractère \.
- Par exemple, \n, \t et \b correspondent comme en C au saut de ligne, à la tabulation et au retour en arrière.

Flex: Fast Lexer

Specification d'expressions régulières avec Flex:

- " : la chaîne elle-même.
"abc" : spécifie la chaîne abc
- [] : un des éléments de l'ensemble.
[abc] : a, b ou c
[a - z] : toutes les lettres minuscules
- . : tout caractère sauf \n.
- | : l'alternance (l'union).
(ab|bc) : la chaîne ab ou la chaîne bc.
- * : zéro ou plusieurs fois (l'étoile).
(x|y)* : 0 ou plusieurs caractères x ou y.
- + : un ou plusieurs fois (l'étoile positive).
[a z]+ : 1 ou plusieurs lettres minuscules.

© Compilation / M.MOURCHID,
SMI/2021-2022

23

Flex: Fast Lexer

Specification d'expressions régulières avec Flex:

- ? : opérateur d'occurrence 0 ou 1 fois. ab?c : la chaîne abc ou ac.
- / : condition de reconnaissance.
ab/cd : la chaîne ab seulement si elle est suivie de la chaîne cd.
- \$: reconnaissance en fin de ligne.
ab\$: la chaîne ab seulement si elle est en fin de ligne.
- ^ : reconnaissance en début de ligne.
^ ab : la chaîne ab seulement si elle est en début de ligne.
- { } : répétition bornée.
a{1,5} : les chaînes a, aa, aaa, aaaa ou aaaaa.
a{2,} : les chaînes aa, aaa, aaaa, ...etc.

© Compilation / M.MOURCHID,
SMI/2021-2022

24

Flex: Fast Lexer

Résolution des conflits

- En cas de conflit, Flex choisit toujours la règle qui produit le plus long lexème.

Exemple:

- "prog" action1
 - "program" action2
- La deuxième règle sera choisie en cas de conflit.

Flex: Fast Lexer

Résolution des conflits

- Si plusieurs règles donnent des lexèmes de mêmes longueurs, Flex choisit la première.

Exemple:

- "prog" action1
 - [a — z]+ action2

La première règle sera choisie en cas de conflit de lexèmes ayant mêmes longueurs.

Flex: Fast Lexer

Résolution des conflits

- Si aucune règle ne correspond au flot d'entrée, Flex choisit sa règle par défaut implicite :

.| \n {ECHO }

- Cette règle par défaut, **recopie** le flot d'entrée sur le flot de sortie.

© Compilation / M.MOURCHID,
SMI/2021-2022

27

Flex: Fast Lexer

Les actions:

- Une **action** est un bloc d'instructions qui est exécutée lorsque la chaîne de caractères lue correspond à la chaîne spécifiée avant l'action.
- Flex dispose de quelques **variables** et **actions prédéfinies**.

© Compilation / M.MOURCHID,
SMI/2021-2022

28

Flex: Fast Lexer

Quelques variables de Flex:

- **yytext** : la chaîne reconnue par l'expression régulière (le lexème courant).
- **yytext** : la longueur de yytext.
- **yyin** : le fichier d'entrée (c'est un **FILE***).
- **yyout** : le fichier de sortie (c'est un **FILE***).

Flex: Fast Lexer

Quelques macros de Flex:

- **ECHO** : recopie le lexème courant.
`[a-z]+ ECHO;`
 ce qui équivaut à `{printf("%s", yytext); }`
- **REJECT** : permet d'envisager la "deuxième meilleure" règle pour laquelle le flot d'entrée (ou un de ses préfixes) est reconnu.
`foo {f (); REJECT; }`
`[^\\t\\n]+ { --word_count; }`
 ce qui permet de compter le nombre total des caractères et exécuter la fonction `f()` chaque fois que la chaîne `foo` est rencontrée.

Flex: Fast Lexer

Quelques directives de Flex:

- **yymore()** : indique que lors de la prochaine application d'une règle, la chaîne reconnue doit être concaténée à *yytext* qui ne doit donc pas écraser *yytext*.

Exemple:

mega- { **ECHO**; *yymore()*; }

octets **ECHO**;

*Pour l'entrée mega-octets, l'analyseur fournira
mega-mega-octets.*

Flex: Fast Lexer

Quelques directives de Flex:

- **yylless(n)** : indique que les *n* premiers caractères de *yytext* doivent être pris en compte lors de la prochaine application d'une règle.

Exemple:

foobar { **ECHO**; *yylless(3)*; }

[a-z]+ **ECHO**;

Pour l'entrée foobar, l'analyseur fournira foobarbar.

Flex: Fast Lexer

Quelques directives de Flex

- **yymore()** : indique que lors de la prochaine application d'une règle, la chaîne reconnue doit être concaténée à *yytext* qui ne doit donc pas écraser *yytext*.

Exemple:

*mega- { **ECHO**; yymore(); } octets **ECHO**;*

- *Pour l'entrée mega-octets, l'analyseur fournira mega-mega-octets.*

Flex: Fast Lexer

Quelques directives de Flex

- **yymore()** : indique que lors de la prochaine application d'une règle, la chaîne reconnue doit être concaténée à *yytext* qui ne doit donc pas écraser *yytext*.

Exemple:

*mega- { **ECHO**; yymore(); } octets **ECHO**;*

Pour l'entrée mega-octets, l'analyseur fournira mega-mega-octets.

Flex: Fast Lexer

Quelques directives de Flex

- **yymore()** : indique que lors de la prochaine application d'une règle, la chaîne reconnue doit être concaténée à *yytext* qui ne doit donc pas écraser *yytext*.

Exemple:

*mega- { **ECHO**; yymore(); } octets **ECHO**;*

*Pour l'entrée mega-octets, l'analyseur fournira
mega-mega-octets.*