

Projet d'Analyse syntaxique

Licence d'informatique

—2018-2019—

Le but du projet est d'écrire un analyseur syntaxique en utilisant les outils `flex` et `bison`.

Le langage source est un petit langage de programmation appelé TPC, qui ressemble à un sous-ensemble du langage C.

Le projet est à faire en binôme ou seul. Si vous n'avez pas de binôme, contactez Eric Laporte. Les dates limite de rendu sont :

- le samedi 10 novembre 2018 à 23h55 au plus tard pour une première version de l'analyseur ⁽¹⁾ avec `makefile`, 2 jeux de tests (au moins 1 fichier dans chaque jeu), script de test, affichage des numéros de ligne en cas d'erreur (1/4 de la note de projet) ;
- le samedi 12 janvier 2019 à 23h55 au plus tard pour l'analyseur complet (3/4 de la note de projet).

1 Définition informelle du langage source

Un programme TPC est une suite de fonctions. Chaque fonction est constituée de déclarations de constantes et variables (locales à la fonction), et d'une suite d'instructions. Les fonctions peuvent être récursives. Il peut y avoir des constantes et variables de portée globale. Elles sont alors déclarées avant les fonctions.

Les types de base du langage sont `int` (entier signé) et `char`. Le mot clé `void` est utilisé pour indiquer qu'une fonction ne fournit pas de résultat ou n'a pas d'arguments.

Le langage TPC utilise `print` pour afficher un entier ou un caractère. Le mot-clé `readc` permet d'obtenir un caractère lu au clavier ; `reade` permet de lire un entier.

2 Définition des éléments lexicaux

Les identificateurs sont constitués d'une lettre, suivie éventuellement de lettres, chiffres, symbole souligné ("`_`"). Vous pouvez fixer une longueur maximale pour un identificateur. Il y a distinction entre majuscule et minuscule. Les mots-clés comme `if`, `else`, `return`, etc., doivent être écrits en minuscules. Ils sont reconnus par l'analyseur lexical et ne peuvent pas être utilisés comme identificateurs.

Les éléments lexicaux pour les constantes numériques sont des suites de chiffres.

Les caractères littéraux dans le programme sont délimités par le symbole `'`, comme en C. Dans les caractères littéraux, la barre oblique inverse ("`\`") est utilisée pour déspecialiser `'` et pour specialiser `n` et `t` : `\n` et `\t` sont le caractère fin de ligne et la tabulation.

Les commentaires sont délimités par `/*` et `*/` et ne peuvent pas être imbriqués.

Les différents opérateurs et autres éléments lexicaux sont :

<code>=</code>	: opérateur d'affectation
<code>+</code>	: addition ou plus unaire
<code>-</code>	: soustraction ou moins unaire
<code>*</code>	: multiplication
<code>/</code> et <code>%</code>	: division et reste de la division entière
<code>!</code>	: négation booléenne
<code>==</code> , <code>!=</code> , <code><</code> , <code>></code> , <code><=</code> , <code>>=</code>	: les opérateurs de comparaison
<code>&&</code> , <code> </code>	: les opérateurs booléens
<code>;</code> et <code>,</code>	: le point-virgule et la virgule
<code>(</code> , <code>)</code> , <code>{</code> , <code>}</code> , <code>[</code> et <code>]</code>	: les parenthèses, les accolades et les crochets

Chacun de ces éléments sera identifié par l'analyse lexicale qui devra produire une erreur pour tout élément ne faisant pas partie du lexique du langage.

(1). Déposez votre projet sur la plateforme elearning dans la zone prévue à cet effet, sous la forme d'une archive tar compressée de nom "ProjetASL3_NOM1_NOM2.tar.gz", qui, au désarchivage, crée un répertoire "ProjetASL3_NOM1_NOM2" contenant le projet.

3 Notations et sémantique du langage

Dans ce qui suit,

- **CARACTERE** et **NUM** désignent respectivement un caractère littéral et une constante numérique ;
- **IDENT** désigne un identificateur ;
- **TYPE** désigne un nom de type qui peut être **int** ou **char** ;
- **EQ** désigne les opérateurs d'égalité ('==') et d'inégalité ('!=') ;
- **ORDER** désigne les opérateurs de comparaison '<', '<=', '>' et '>=' ;
- **ADDSUB** désigne les opérateurs '+' et '-' (binaire ou unaire) ;
- **DIVSTAR** désigne les opérateurs '*', '/' et '%' ;
- **OR** et **AND** désignent les deux opérateurs booléens '||' et '&&' .
- Les mots-clés sont notés par des unités lexicales qui leur sont identiques à la casse près.

L'instruction nulle est notée ';' .

4 Grammaire du langage TPC

```
Prog      : DeclConsts DeclVars DeclFoncts  ;
DeclConsts : DeclConsts CONST ListConst ';' ;
ListConst  : ListConst ',' IDENT '=' Litteral
            | IDENT '=' Litteral  ;
Litteral   : NombreSigne
            | CARACTERE  ;
NombreSigne : NUM
            | ADDSUB NUM  ;
DeclVars   : DeclVars TYPE Declarateurs ';' ;
Declarateurs : Declarateurs ',' Declarateur
            | Declarateur  ;
Declarateur : IDENT
            | IDENT '[' NUM ']'  ;
DeclFoncts  : DeclFoncts DeclFonct
            | DeclFonct  ;
DeclFonct   : EnTeteFonct Corps  ;
EnTeteFonct : TYPE IDENT '(' Parametres ')'
            | VOID IDENT '(' Parametres ')'  ;
Parametres  : VOID
            | ListTypVar  ;
ListTypVar  : ListTypVar ',' TYPE IDENT
            | TYPE IDENT  ;
Corps       : '{' DeclConsts DeclVars SuiteInstr '}'  ;
SuiteInstr  : SuiteInstr Instr
            | ;
Instr       : Exp ';'
            | ';'
            | RETURN Exp ';'
            | RETURN ';'
            | READE '(' IDENT ')' ';'
            | READC '(' IDENT ')' ';'
            | PRINT '(' Exp ')' ';'
            | IF '(' Exp ')' Instr
            | IF '(' Exp ')' Instr ELSE Instr
            | WHILE '(' Exp ')' Instr
            | '{' SuiteInstr '}'  ;
Exp         : LValue '=' Exp
            | EB  ;
```

```

EB      : EB OR TB
        | TB ;
TB      : TB AND FB
        | FB ;
FB      : FB EQ M
        | M ;
M       : M ORDER E
        | E ;
E       : E ADDSUB T
        | T ;
T       : T DIVSTAR F
        | F ;
F       : ADDSUB F
        | '!' F
        | '(' Exp ')'
        | LValue
        | NUM
        | CARACTERE
        | IDENT '(' Arguments ')' ;
LValue  : IDENT
        | IDENT '[' Exp ']' ;
Arguments : ListExp
        | ;
ListExp  : ListExp ',' Exp
        | Exp ;

```

5 Travail demandé

Analyseur syntaxique Écrire un analyseur syntaxique de ce langage en utilisant **flex** pour l'analyse lexicale et **bison** pour l'analyse syntaxique. Les messages d'erreur doivent donner le numéro de ligne. Vous pouvez modifier la grammaire, mais vos modifications ne peuvent affecter le langage engendré que si cela enrichit le langage TPC.

Le répertoire que vous déposerez doit être organisé correctement : un répertoire pour les sources, un autre pour la documentation, un autre pour les tests... Le répertoire pour les sources doit contenir un Makefile nommé **Makefile** ou **makefile**. L'analyseur créé avec le Makefile doit être nommé **as**. La commande suivante doit exécuter votre analyseur :

```
./as < prog.tpc
```

Le programme devra renvoyer 0 si et seulement si *prog.tpc* est syntaxiquement correct.

Tests Écrire deux jeux de tests, un pour les programmes TPC corrects et un autre pour les programmes incorrects, et un script de déploiement des tests, qui produit un rapport unique donnant les résultats de tous les tests.

Tableaux à plusieurs dimensions Ajouter des règles permettant des tableaux à plusieurs dimensions, avec la syntaxe du C.

Redémarrage après erreur Ajouter des règles permettant un redémarrage après erreur, pour pouvoir détecter plusieurs erreurs pendant une même exécution de l'analyseur.

Documentation Vous décrierez dans votre documentation vos choix et les difficultés que vous avez rencontrées.

Déposez votre projet sur la plateforme elearning dans la zone prévue à cet effet, sous la forme d'une archive tar compressée de nom "ProjetASL3_NOM1_NOM2.tar.gz", qui, au désarchivage, crée un répertoire "ProjetASL3_NOM1_NOM2" contenant le projet.