

# Projet du cours « Compilation »

## Jalon 1 : Analyse lexicale et syntaxique de HOPIX

version Tue Oct 9 15:27:50 CEST 2018

## 1 Spécification de la grammaire

### 1.1 Notations extra-lexicales

Les commentaires, espaces, les tabulations et les sauts de ligne jouent le rôle de séparateurs. Leur nombre entre les différents symboles terminaux peut donc être arbitraire. Ils sont ignorés par l'analyse lexicale : ce ne sont pas des lexèmes.

Les commentaires sont entourés des deux symboles « `(` » et « `)` ». Par ailleurs, ils peuvent être imbriqués.

### 1.2 Symboles

**Symboles terminaux** Les terminaux sont répartis en trois catégories : les mots-clés, les identificateurs et la ponctuation.

Les mots-clés sont les noms réservés aux constructions du langage. Ils seront écrits avec des caractères de machine à écrire (comme par exemple les mots-clés `if` et `while`).

Les identificateurs sont constitués des identificateurs de variables, d'étiquettes, de constructeurs de données et de types ainsi que des littéraux, comprenant les constantes entières, les caractères et les chaînes de caractères. Ils seront écrits dans une police **sans-serif** (comme par exemple `type_con` ou `int`). La classification des identificateurs est définie par les expressions rationnelles suivantes :

<code>alien_infix_id</code>	$\equiv$ <code>' [a-z 0-9 + - * / = _ ! ?]^+ '</code>	<i>Identificateur d'opérateurs infixes</i>
<code>alien_prefix_id</code>	$\equiv$ <code>' [a-z 0-9 + - * / = _ ! ?]^+ '</code>	<i>Identificateur d'opérateurs préfixes</i>
<code>var_id</code>	$\equiv$ <code>[a-z] [A-Z a-z 0-9 _]*</code>   <code>alien_prefix_id</code>	<i>Identificateur de variables préfixe</i>
<code>all_var_id</code>	$\equiv$ <code>var_id</code>   <code>alien_infix_id</code>	<i>Identificateur de variables de toutes sortes</i>
<code>constr_id</code>	$\equiv$ <code>' A-Z [A-Z a-z 0-9 _]*</code>	<i>Identificateur de constructeurs de données</i>
<code>label_id</code>	$\equiv$ <code>[a-z] [A-Z a-z 0-9 _]*</code>	<i>Identificateur d'étiquettes d'enregistrement</i>
<code>type_con</code>	$\equiv$ <code>' A-Z [A-Z a-z 0-9 _]*</code>	<i>Identificateur de constructeurs de type</i>
<code>type_variable</code>	$\equiv$ <code>[a-z] [A-Z a-z 0-9 _]*</code>	<i>Identificateur de variables de type</i>
<code>int</code>	$\equiv$ <code>-?[0-9]^+   0x[0-9 a-f A-F]^+   0b[0-1]^+   0o[0-7]^+</code>	<i>Littéraux entiers</i>
<code>char</code>	$\equiv$ <code>'atom'</code>	<i>Littéraux caractères</i>
<code>atom</code>	$\equiv$ <code>\000   ...   \255   \0[x][0-9 a-f A-F]^2   [printable]   \\   \'   \n   \t   \b   \r</code>	
<code>string</code>	$\equiv$ <code>" (atom   [']) -{ " , \ ' }   \" } * "</code>	<i>Littéraux chaîne de caractères</i>

Autrement dit, les identificateurs de valeur, de variable de type et de champs d'enregistrement commencent par une lettre minuscule et peuvent comporter ensuite des majuscules, des minuscules, des chiffres et le caractère souligné `_`. Les identificateurs de constructeurs de données et de constructeurs de type peuvent comporter les mêmes caractères, mais doivent commencer par une majuscule ou un guillemet arrière. Par ailleurs, un identificateur d'opérateur est formé de symboles et de caractères alphanumériques. S'il est entouré de deux `'`, il est dit *infixe*. Enfin, notez que le caractère `_` représente un espace.

Les constantes entières sont constituées de chiffres en notation décimale, en notation hexadécimale, en notation binaire ou en notation octale. Les entiers utilisent une représentation binaire sur 32 bits en complément à deux. Les constantes entières sont donc prises dans  $[-2^{31}; 2^{31} - 1]$ .

Les constantes de caractères sont écrites entre guillemets simples (ce qui signifie en particulier que les guillemets simples doivent être échappés dans les constantes de caractères). On y trouve tous les symboles ASCII affichables (voir la spécification de ASCII pour plus de détails). Par ailleurs, sont des caractères valides : les séquences d'échappement usuelles, ainsi que les séquences d'échappement de trois chiffres décrivant le code ASCII du caractère en notation décimale ou encore les séquences d'échappement de deux chiffres décrivant le code ASCII en notation hexadécimale.

Les constantes de chaîne de caractères sont formées d'une séquence de caractères. Cette séquence est entourée de guillemets (ce qui signifie en particulier que les guillemets doivent être échappés dans les chaînes).

Les symboles seront notés avec la police "machine à écrire" (comme par exemple « `(` » ou « `=` »).

**Symboles non-terminaux** Les symboles non-terminaux seront notés à l'aide d'une police légèrement inclinée (comme par exemple *expr*).

Une séquence entre crochets est optionnelle (comme par exemple « *[ ref ]* »). Attention à ne pas confondre ces crochets avec les symboles terminaux de ponctuation notés *[* et *]*. Une séquence entre accolades se répète zéro fois ou plus, (comme par exemple « *( arg { , arg } )* »).

## 2 Grammaire en format BNF

La grammaire du langage est spécifiée à l'aide du format BNF.

**Programme** Un programme est constitué d'une séquence de définitions de types et de valeurs.

<i>p</i> ::= { <i>definition</i> }	<i>Programme</i>
<i>definition</i> ::= <b>type</b> <i>type_con</i> [ < <i>type_variable</i> { , <i>type_variable</i> } > ] [= <i>tdefinition</i> ]   <b>extern</b> <i>all_var_id</i> : <i>type_scheme</i>   <i>vdefinition</i>	<i>Définition de type</i> <i>Valeurs externes</i> <i>Définition de valeur(s)</i>
<i>tdefinition</i> ::= [ <i>l</i> ] <i>constr_id</i> [ ( <i>type</i> { , <i>type</i> } ) ] { <i>l</i> <i>constr_id</i> [ ( <i>type</i> { , <i>type</i> } ) ] }   { <i>label_id</i> : <i>type</i> { ; <i>label_id</i> : <i>type</i> } }	<i>Type somme</i> <i>Type produit étiqueté</i>
<i>vdefinition</i> ::= <b>val</b> <i>var_id</i> [ : <i>type_scheme</i> ] = <i>expr</i>   <b>def</b> <i>fundef</i> { <b>and</b> <i>fundef</i> }	<i>Valeur simple</i> <i>Fonction(s)</i>
<i>fundef</i> ::= <i>all_var_id</i> [ : <i>type_scheme</i> ] ( [ <i>var_id</i> { , <i>var_id</i> } ] ) = <i>expr</i>	

**Types de données** La syntaxe des types est donnée par la grammaire suivante :

<i>type</i> ::= <i>type_con</i> [ < <i>type</i> { , <i>type</i> } > ]   <i>type</i> { * <i>type</i> } -> <i>type</i>   <i>type_variable</i>   ( <i>type</i> )
<i>type_scheme</i> ::= [ <b>forall</b> <i>type_variable</i> { , <i>type_variable</i> } . ] <i>type</i>

**Expression** La syntaxe des expressions du langage est donnée par la grammaire suivante.

<i>expr</i> ::= <b>int</b>	<i>Entier positif</i>
<b>char</b>	<i>Caractère</i>
<b>string</b>	<i>Chaîne de caractères</i>
<i>all_var_id</i> [ < [ <i>type</i> { , <i>type</i> } ] > ]	<i>Variable</i>
<i>constr_id</i> [ < [ <i>type</i> { , <i>type</i> } ] > ] [ ( <i>expr</i> { , <i>expr</i> } ) ]	<i>Construction d'une donnée</i>
{ <i>label_id</i> = <i>expr</i> { ; <i>label_id</i> = <i>expr</i> } } [ < [ <i>type</i> { , <i>type</i> } ] > ]	<i>Construction d'un enregistrement</i>
<i>expr</i> . <i>label_id</i>	<i>Accès à un champ</i>
<i>expr</i> ; <i>expr</i>	<i>Séquencement</i>
<i>vdefinition</i> ; <i>expr</i>	<i>Définition locale</i>
<b>fun</b> ( [ <i>var_id</i> { , <i>var_id</i> } ] ) => <i>expr</i>	<i>Fonction anonyme</i>
<i>expr</i> ( [ <i>expr</i> { , <i>expr</i> } ] )	<i>Application</i>
<i>expr</i> <i>binop</i> <i>expr</i>	<i>Application infixe</i>
<b>case</b> <i>expr</i> { <i>branches</i> }	<i>Analyse de motifs</i>
<b>if</b> <i>expr</i> <b>then</b> <i>expr</i> [ <b>else</b> <i>expr</i> ]	<i>Conditionnelle</i>
<b>ref</b> <i>expr</i>	<i>Allocation</i>
<i>expr</i> := <i>expr</i>	<i>Affectation</i>
! <i>expr</i>	<i>Lecture</i>
<b>while</b> <i>expr</i> { <i>expr</i> }	<i>Boucle non bornée</i>
<b>for</b> <i>var_id</i> = <i>expr</i> <b>to</b> <i>expr</i> [ <b>by</b> <i>expr</i> ] { <i>expr</i> }	<i>Boucle bornée</i>
( <i>expr</i> )	<i>Parenthésage</i>
( <i>expr</i> : <i>type</i> )	<i>Annotation de type</i>

Voici la grammaire des définitions auxiliaires utilisées par la grammaire des expressions :

$binop ::= + \mid - \mid * \mid / \mid \&\& \mid    \mid =? \mid <=? \mid >=? \mid <? \mid >? \mid \text{alien\_infix\_id}$	<i>Opérateurs binaires</i>
$branches ::= [ \mid ] \text{ branch } \{ \mid \text{ branch } \}$	<i>Liste de cas</i>
$branch ::= pattern \Rightarrow expr$	<i>Cas d'analyse</i>

**Motifs** Les motifs (*patterns* en anglais), utilisés par l'analyse de motifs, ont la syntaxe suivante :

$pattern ::= \text{var\_id}$	<i>Motif universel liant</i>
$\mid -$	<i>Motif universel non liant</i>
$\mid ( pattern )$	<i>Parenthésage</i>
$\mid pattern : type$	<i>Annotation de type</i>
$\mid \text{int}$	<i>Entier</i>
$\mid \text{char}$	<i>Caractère</i>
$\mid \text{string}$	<i>Chaîne de caractères</i>
$\mid \text{constr\_id } [ < [ type \{ , type \} ] > ] [ ( pattern \{ , pattern \} ) ]$	<i>Valeurs étiquetées</i>
$\mid \{ \text{label\_id} = pattern \{ ; \text{label\_id} = pattern \} \} [ < [ type \{ , type \} ] > ]$	<i>Enregistrement</i>
$\mid pattern \mid pattern$	<i>Disjonction</i>
$\mid pattern \& pattern$	<i>Conjonction</i>

**Remarques** Notez bien que la grammaire spécifiée plus haut est ambiguë ! Vous devez fixer des priorités entre les différentes constructions ainsi que des associativités aux différents opérateurs. *In fine*, c'est la batterie de tests en ligne qui vous permettra de valider vos choix. Cependant, il est fortement conseillé de poser des questions sur la liste de diffusion du cours pour obtenir des informations supplémentaires sur les règles de disambiguation associées à cette grammaire.

### 3 Code fourni

Un squelette de code vous est fourni, il est disponible sur le dépôt GIT du cours.

`git@moule.informatique.univ-paris-diderot.fr:Yann/compilation-m1-2018.git`

Vous devez vous connecter sur le Gitlab disponible ici :

`http://moule.informatique.univ-paris-diderot.fr:8080`

et vous créer un dépôt par branchement (*fork*) du projet `compilation-m1-2018`.

L'arbre de sources contient des **Makefiles** ainsi que des modules O'CAML à compléter.

La commande **make** produit un exécutable appelé **flap**. On doit pouvoir l'appeler avec un nom de fichier en argument. En cas de réussite (de l'analyse syntaxique), le code de retour de ce programme doit être 0. Dans le cas d'un échec, le code de retour doit être 1.

### 4 Travail à effectuer

La première partie du projet est l'écriture de l'analyseur lexical et de l'analyseur syntaxique spécifiés par la grammaire précédente.

Le projet est à rendre **avant le** :

**21 octobre 2018 à 23h59**

Pour finir, vous devez vous assurer des points suivants :

- Le projet contenu dans cette archive **doit compiler**.
- Vous devez **être les auteurs** de ce projet.
- Il doit être rendu **à temps**.

Si l'un de ces points n'est pas respecté, la note de 0 vous sera affectée.

## 5 Log

2018-10-09 yrg <yrg@irif.fr>

- \* Corrige 'alien\_infix\_id' et 'alien\_prefix\_id'
- \* Corrige les variantes hexa, binaire et octale de 'int'
- \* Corrige la définition de la règle des déclarations externes
- \* Corrige la syntaxe du while

2018-09-14 yrg <yrg@irif.fr>

- \* Version initiale