

# Projet du cours « Compilation »

## Jalon 1 : Analyse lexicale et syntaxique de HOPIX

version numéro jeudi 12 novembre 2015, 08 :21 :59 (UTC+0100)

## 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.

On peut aussi introduire un commentaire à l'aide du symbole « `**` » : tout ce qui le suit jusqu'à la fin de la ligne est interprété comme un commentaire.

### 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>	$\equiv$ ' [A-Z a-z 0-9 + - * / < = > _] <sup>+</sup> '	<i>Identificateur d'opérateurs infixes</i>
<code>alien_prefix_id</code>	$\equiv$ ' [A-Z a-z 0-9 + - * / < = > _] <sup>+</sup>	<i>Identificateur d'opérateurs préfixes</i>
<code>var_id</code>	$\equiv$ [a-z] [A-Z a-z 0-9 _] <sup>*</sup>   <code>alien_prefix_id</code>	<i>Identificateur de variables</i>
<code>label_id</code>	$\equiv$ [a-z] [A-Z a-z 0-9 _] <sup>*</sup>	<i>Identificateur de champs</i>
<code>constr_id</code>	$\equiv$ [A-Z _] [A-Z a-z 0-9 _] <sup>*</sup>	<i>Identificateur de constructeurs de données</i>
<code>type_con</code>	$\equiv$ [a-z] [A-Z a-z 0-9 _] <sup>*</sup>	<i>Identificateur de constructeurs de type</i>
<code>type_variable</code>	$\equiv$ ' [a-z] [A-Z a-z 0-9 _] <sup>*</sup>	<i>Identificateur de variables de type</i>
<code>int</code>	$\equiv$ [0-9] <sup>+</sup>   0[xX][0-9 a-f A-F] <sup>+</sup>   0[Bb][0-1] <sup>+</sup>	<i>Littéraux entiers</i>
<code>char</code>	$\equiv$ 'atom'	<i>Littéraux caractères</i>
<code>atom</code>	$\equiv$ \000   ...   \255   \0[xX][0-9 a-f A-F] <sup>2</sup>   \0[bB][0-1] <sup>+</sup>   [printable]   \   \   \   \n   \t   \b   \r	
<code>string</code>	$\equiv$ " ( atom -{ " }   \ " ) <sup>*</sup> "	<i>Littéraux chaîne de caractères</i>

Autrement dit, les identificateurs de variables, de champs et de types 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 peuvent comporter les mêmes caractères, mais doivent commencer par une majuscule ou par un caractère `_`. Par ailleurs, un identificateur d'opérateur est formé de symboles et de caractères alphanumériques. S'il est entouré de deux `'`, il est dit *infixe*. S'il est préfixé par un `'`, il est dit *préfixe*.

Les constantes entières sont constituées de chiffres en notation décimale, en notation hexadécimale ou en notation binaire. Les constantes entières sont prises dans  $[-2^{31}; 2^{31} - 1]$ .

Les constantes de caractères sont décrites entre guillemets simples (ce qui signifie en particulier que les guillemets simples doivent être échappés dans les constantes de caractères). On y trouve en particulier les symboles ASCII affichables. 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> ::= <i>type</i> <i>type_con</i> [ [ <i>type_variable</i> [ { , <i>type_variable</i> } ] ] ] [ := <i>tdefinition</i> ] .	<i>Définition de type</i>
<i>extern</i> <i>var_id</i> : <i>type</i> .	<i>Valeurs externes</i>
<i>vdefinition</i> .	<i>Définition de valeur(s)</i>
<i>tdefinition</i> ::= { <i>label_id</i> : <i>type</i> { ; <i>label_id</i> : <i>type</i> } }	<i>Type d'enregistrement</i>
{ [   ] <i>constr_id</i> [ : <i>type</i> { * <i>type</i> } ] }	<i>Type somme</i>
{   <i>constr_id</i> [ : <i>type</i> { * <i>type</i> } ] } }	
<i>vdefinition</i> ::= <i>val</i> <i>var_id</i> { <i>pattern</i> } [ : <i>type</i> ] := <i>expr</i>	<i>Valeur simple</i>
<i>rec</i> <i>var_id</i> { <i>pattern</i> } [ : <i>type</i> ] := <i>expr</i>	<i>Valeurs mutuellement récursives</i>
<i>and</i> <i>var_id</i> { <i>pattern</i> } [ : <i>type</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_variable</i>
( <i>type</i> )

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

<i>expr</i> ::= <i>int</i>	<i>Entier</i>
<i>char</i>	<i>Caractère</i>
<i>string</i>	<i>Chaîne de caractères</i>
<i>var_id</i>	<i>Variable</i>
<i>constr_id</i> [ ( <i>expr</i> { , <i>expr</i> } ) ]	<i>Construction d'une donnée étiquetée</i>
{ <i>label_id</i> := <i>expr</i> { ; <i>label_id</i> := <i>expr</i> } }	<i>Construction d'un enregistrement</i>
( <i>expr</i> : <i>type</i> )	<i>Annotation de type</i>
<i>do</i> <i>expr</i> { ; <i>expr</i> } [ ; ] <i>done</i>	<i>Séquencement</i>
<i>vdefinition</i> ; <i>expr</i>	<i>Définition locale</i>
<i>expr</i> <i>expr</i>	<i>Application</i>
\ <i>pattern</i> { <i>pattern</i> } => <i>expr</i>	<i>Fonction anonyme</i>
<i>expr</i> <i>binop</i> <i>expr</i>	<i>Opérations binaires</i>
<i>expr</i> # <i>label_id</i>	<i>Accès à un champ</i>
<i>expr</i> # <i>label_id</i> <- <i>expr</i>	<i>Modification d'un champ</i>
<i>expr</i> ? <i>branches</i>	<i>Analyse de motifs</i>
<i>if</i> <i>expr</i> <i>then</i> <i>expr</i> <i>else</i> <i>expr</i> <i>fi</i>	<i>Conditionnelle</i>
( <i>expr</i> )	<i>Parenthésage</i>
<i>binop</i> ::= +   -   *   /   &&        =   <=   >=   <   >	<i>Opérateurs binaires</i>
<i>branches</i> ::= [   ] <i>branch</i> {   <i>branch</i> }	<i>Liste de cas</i>
{ [   ] <i>branch</i> {   <i>branch</i> } }	<i>... avec des accolades</i>
<i>branch</i> ::= <i>pattern</i> -> <i>expr</i>	<i>Cas d'analyse</i>

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

<code>pattern ::= constr_id [ ( pattern { , pattern } ) ]</code>	<i>Valeurs étiquetées</i>
<code>  { label_id = pattern { ; label_id = pattern } }</code>	<i>Enregistrement</i>
<code>  pattern   pattern</code>	<i>Disjonction</i>
<code>  pattern &amp; pattern</code>	<i>Conjonction</i>
<code>  var_id</code>	<i>Motif universel liant</i>
<code>  -</code>	<i>Motif universel non liant</i>
<code>  ( pattern )</code>	<i>Parenthésage</i>
<code>  ( pattern : type )</code>	<i>Annotation de type</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-m1-2015.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-m1-2015`.

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.

La compilation s'effectue par la commande « **make** ».

Le projet est à rendre **avant le** :

**22 novembre 2015 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