

Compilateur Pascal avec modules

première partie : analyse syntaxique

à rendre avant le lundi 12 octobre à 12h00

Le langage considéré dans ce projet est un sous-ensemble de Pascal (pas de pointeurs, pas de fonctions locales, etc.), augmenté d'un système de modules à la Caml.

Dans la suite, nous utilisons les notations suivantes dans les grammaires :

| | |
|--------------------------------------|---|
| $\langle \text{r\`egle} \rangle^*$ | répétition de la règle $\langle \text{r\`egle} \rangle$ un nombre quelconque de fois (y compris aucune) |
| $\langle \text{r\`egle} \rangle_t^*$ | répétition de la règle $\langle \text{r\`egle} \rangle$ un nombre quelconque de fois (y compris aucune), les occurrences étant séparées par le terminal t |
| $\langle \text{r\`egle} \rangle^+$ | répétition de la règle $\langle \text{r\`egle} \rangle$ au moins une fois |
| $\langle \text{r\`egle} \rangle_t^+$ | répétition de la règle $\langle \text{r\`egle} \rangle$ au moins une fois, les occurrences étant séparées par le terminal t |
| $\langle \text{r\`egle} \rangle ?$ | utilisation optionnelle de la règle $\langle \text{r\`egle} \rangle$ (i.e. 0 ou 1 fois) |
| $(\langle \text{r\`egle} \rangle)$ | parenthésage ; attention à ne pas confondre ces parenthèses avec les terminaux $($ et $)$ |

1 Analyse lexicale

Espaces, tabulations et retour-chariots constituent les blancs. Les commentaires peuvent prendre deux formes :

- débutant par $(*$ et s'étendant jusqu'à $*)$, et ne pouvant être imbriqués ;
- débutant par $\{$ et s'étendant jusqu'à $\}$, et ne pouvant être imbriqués.

Les identificateurs obéissent à l'expression régulière $\langle \text{ident} \rangle$ suivante :

$$\begin{aligned} \langle \text{chiffre} \rangle &::= 0-9 \\ \langle \text{alpha} \rangle &::= a-z \mid A-Z \\ \langle \text{ident} \rangle &::= \langle \text{alpha} \rangle (\langle \text{alpha} \rangle \mid _ \mid \langle \text{chiffre} \rangle)^* \end{aligned}$$

Les identificateurs suivants sont des mots clés :

| | | | | | |
|------|--------|-------|---------|-----------|---------|
| AND | ARRAY | BEGIN | BOOLEAN | DIV | DO |
| ELSE | END | FALSE | FOR | IF | INTEGER |
| MOD | NOT | OF | OR | PROCEDURE | PROGRAM |
| THEN | TO | TRUE | VAR | WHILE | STRUCT |
| SIG | MODULE | OPEN | | | |

Pour les identificateurs et les mots-clés, la casse des caractères est indifférente : ainsi, xy , Xy , xY et XY représentent tous les même identificateur, et de même $begin$, $BEGIN$ ou encore $BeGin$ représentent tous le même mot-clé.

Les constantes entières obéissent à l'expression régulière $\langle \text{entier} \rangle$ suivante :

$$\langle \text{entier} \rangle ::= \langle \text{chiffre} \rangle^+$$

2 Analyse syntaxique

La grammaire des fichiers sources considérée est donnée figure 1. Le point d'entrée est le non-terminal *fichier*.

Les associativités et précédences des divers opérateurs sont données par la table suivante, de la plus faible à la plus forte précedence :

| opérateur | associativité |
|-----------------------|---------------|
| < <= > >= = <> | à gauche |
| + - or | à gauche |
| * / div mod and | à gauche |
| + (unaire) - (unaire) | — |
| not | — |

3 Localisation des erreurs

Lorsqu'une erreur lexicale ou syntaxique est détectée par votre compilateur, elle doit être signalée le plus précisément possible, par sa nature et sa localisation dans le fichier source. On adoptera le format suivant pour cette signalisation :

```
File "test.pas", line 4, characters 5-6:  
syntax error
```

L'anglicisme de la première ligne est nécessaire pour que la fonction `next-error` d'Emacs puisse interpréter la localisation et placer ainsi automatiquement le curseur sur l'emplacement de l'erreur. (En revanche, le message d'erreur proprement dit pourra être écrit en français si vous le souhaitez.)

Les localisations peuvent être obtenues grâce aux fonctions `Lexing.lexeme_start_p`, `Lexing.lexeme_end_p`, `Parsing.rhs_start_pos`, et `Parsing.rhs_end_pos`.

Anticipation. Dans la deuxième partie, les erreurs sémantiques doivent également être localisées avec précision, c'est-à-dire indiquer le sous-terme en cause. Pour cela il est donc nécessaire que l'arbre de syntaxe abstraite contienne cette information de localisation. Il vous est donc conseillé d'anticiper ce besoin et de construire dès à présent un arbre de syntaxe abstraite adéquat.

4 Travail demandé

Écrire un compilateur, appelons-le `compilo`, acceptant sur sa ligne de commande exactement un fichier Pascal et éventuellement l'option `-parse-only`. Cette option est sans effet dans cette première partie et son utilisation sera justifiée par la suite.

Si le fichier est conforme à la syntaxe décrite dans ce document, le programme doit terminer avec le code de sortie 0 (`exit 0` explicite ou terminaison normale du programme), sans rien afficher. En cas d'erreur lexicale ou syntaxique, celle-ci doit être signalée (de la manière indiquée ci-dessus) et le programme doit terminer avec le code de sortie 1 (`exit 1`). En cas d'autre erreur (une erreur du compilateur lui-même), le programme doit terminer avec le code de sortie 2 (`exit 2`).

$\langle \text{fichier} \rangle ::= \text{program } \langle \text{ident} \rangle ; \langle \text{définition} \rangle^* \langle \text{bloc} \rangle . \text{ EOF}$

$\langle \text{définition} \rangle ::= \langle \text{vars} \rangle$
 $\quad | \langle \text{decl_proc} \rangle \langle \text{vars} \rangle ? \langle \text{bloc} \rangle ;$
 $\quad | \text{module } \langle \text{ident} \rangle \langle \text{params_mod} \rangle ? (: \langle \text{signature} \rangle) ? := \langle \text{expr_mod} \rangle ;$
 $\quad | \text{open } \langle \text{lident} \rangle ;$

$\langle \text{vars} \rangle ::= \text{var } \langle \text{decl_var} \rangle^+$
 $\langle \text{decl_var} \rangle ::= \langle \text{ident} \rangle^+ : \langle \text{type} \rangle ;$
 $\langle \text{decl_proc} \rangle ::= \text{procedure } \langle \text{ident} \rangle \langle \text{params} \rangle ? ;$
 $\langle \text{signature} \rangle ::= \text{sig } \langle \text{déclaration} \rangle^* \text{end}$
 $\langle \text{déclaration} \rangle ::= \langle \text{vars} \rangle | \langle \text{decl_proc} \rangle$
 $\quad | \text{module } \langle \text{ident} \rangle \langle \text{params_mod} \rangle ? : \langle \text{signature} \rangle ;$

$\langle \text{params} \rangle ::= (\langle \text{param} \rangle^+ ;)$
 $\langle \text{param} \rangle ::= \text{var } ? \langle \text{ident} \rangle^+ : \langle \text{type} \rangle$
 $\langle \text{params_mod} \rangle ::= (\langle \text{param_mod} \rangle^+ ;)$
 $\langle \text{param_mod} \rangle ::= \langle \text{ident} \rangle : \langle \text{signature} \rangle$

$\langle \text{expr_mod} \rangle ::= \text{struct } \langle \text{définition} \rangle^* (\langle \text{bloc} \rangle ;) ? \text{end}$
 $\quad | \langle \text{lident} \rangle (\langle \text{lident} \rangle^+ ;)$

$\langle \text{type} \rangle ::= \text{boolean} | \text{integer}$
 $\quad | \text{array } [\langle \text{entier} \rangle .. \langle \text{entier} \rangle] \text{ of } \langle \text{type} \rangle$

$\langle \text{constant} \rangle ::= \text{true} | \text{false} | \langle \text{entier} \rangle$
 $\langle \text{expr} \rangle ::= \langle \text{constant} \rangle | \langle \text{lvalue} \rangle$
 $\quad | \text{not } \langle \text{expr} \rangle | - \langle \text{expr} \rangle | + \langle \text{expr} \rangle$
 $\quad | \langle \text{expr} \rangle \langle \text{opérateur} \rangle \langle \text{expr} \rangle$
 $\quad | (\langle \text{expr} \rangle)$

$\langle \text{lident} \rangle ::= \langle \text{ident} \rangle | \langle \text{lident} \rangle . \langle \text{ident} \rangle$
 $\langle \text{lvalue} \rangle ::= \langle \text{lident} \rangle | \langle \text{lvalue} \rangle [\langle \text{expr} \rangle]$
 $\langle \text{opérateur} \rangle ::= = | <> | < | <= | > | >=$
 $\quad | + | - | * | / | \text{div} | \text{mod} | \text{and} | \text{or}$

$\langle \text{instruction} \rangle ::= \langle \text{lvalue} \rangle := \langle \text{expr} \rangle$
 $\quad | \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{instruction} \rangle$
 $\quad | \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{instruction} \rangle \text{ else } \langle \text{instruction} \rangle$
 $\quad | \text{while } \langle \text{expr} \rangle \text{ do } \langle \text{instruction} \rangle$
 $\quad | \text{for } \langle \text{lident} \rangle := \langle \text{expr} \rangle \text{ to } \langle \text{expr} \rangle \text{ do } \langle \text{instruction} \rangle$
 $\quad | \langle \text{bloc} \rangle$
 $\quad | \langle \text{lident} \rangle (\langle \text{expr} \rangle^+ ;)$
 $\quad | \langle \text{lident} \rangle$

$\langle \text{bloc} \rangle ::= \text{begin } \langle \text{instruction} \rangle^* ; (; ?) \text{end}$

FIG. 1 – Grammaire des fichiers Pascal

Votre code devra comprendre (au moins) les fichiers suivants :

- `lexer.mll` : le lexer ocamllex
- `ast.mli` : le(s) type(s) des arbres de syntaxe abstraite
- `parser.mly` : le parseur ocaml yacc
- `main.ml` : le fichier principal

Il est possible d'utiliser d'autres fichiers, si cela est justifié dans le rapport.

Modalités de remise de votre projet. Votre projet doit se présenter sous forme d'une archive tar compressée (option "z" de tar), appelée *vos_noms.tgz* qui doit contenir un répertoire appelé *vos_noms* (exemple : *dupont-durand.tgz*). Dans ce répertoire doivent se trouver les *sources* de votre programme (ne donnez pas les fichiers compilés). Quand on se place dans ce répertoire, la commande `make` doit créer votre compilateur, qui sera appelé `compilo`. La commande `make clean` doit effacer tous les fichiers que `make` a engendré et ne laisser dans le répertoire que les fichiers sources. Le répertoire doit également contenir un rapport d'une page ou deux au format PDF.

Le projet est à faire impérativement en binôme. Il doit être remis à votre chargé de TP par e-mail¹ au plus tard le lundi 12 octobre à 12h00.

¹francois.bobot@lri.fr, dephine.longuet@lri.fr, louis.mandel@lri.fr