

# Aide Mémoire Ocamllex & Ocaml yacc

Armelle Bonenfant

9 novembre 2010

## Résumé

C document est une traduction/résumé des documentations de Ocaml et Ocamlc en ligne :

<http://plus.kaist.ac.kr/~shoh/ocaml/ocaml-ocamlc-ocaml-tutorial.txt>

<http://plus.kaist.ac.kr/~shoh/ocaml/ocaml-ocamlc-ocamlc-tutorial>

(Copyright (C) 2004 SooHyung Oh.).

Il présente les fonctionnalités de Ocaml et d Ocamlc.

# Table des matières

<b>1</b>	<b>Ocamllex</b>	<b>7</b>
1.1	Commandes et structure	7
1.2	Input	7
1.3	Patterns	7
1.4	Match	9
1.5	Action	10
1.6	Output : le scanner généré	11
1.7	Démarrage	11
1.8	Interface avec Ocamlyacc	12
1.9	Options	12
<b>2</b>	<b>Ocamlyacc</b>	<b>14</b>
2.1	Introduction	14
2.2	Concepts	14
2.2.1	Langages et grammaires hors-contexte (type 2)	14
2.2.2	Des règles formelles à l'entrée Ocamlyacc	15
2.2.3	Valeurs sémantiques	16
2.2.4	Actions sémantiques	16
2.2.5	Positions	17
2.2.6	Sorti d'Ocamlyacc : le fichier du parser	17
2.2.7	Les étapes pour utiliser Ocamlyacc	18
2.3	Le fichier de grammaire Ocamlyacc	18
2.3.1	Structure du fichier	18
2.3.2	Symbols, terminaux et non-terminaux	19
2.3.3	Syntax des règles de grammaire	20
2.3.4	Règles récursives	21
2.3.5	Définir la sémantique d'un langage	22
2.3.6	Suivre les positions	23
2.3.7	Déclarations Ocamlyacc	25
2.4	L'interface du parser	27
2.4.1	La fonction Parsr	27
2.4.2	La fonction d'analyse lexical	27
2.4.3	La fonction de rapport d'erreur	28
2.5	L'algorithme du parser Ocamlyacc	28
2.5.1	Les tokens en avant	28
2.5.2	Conflits shift/réduire	29
2.5.3	Précédence des opérateurs	30
2.5.4	6.4. Contexte-Dépendant Précedence	32
2.5.5	Etats du parser	33

2.5.6	Les conflits R duc /R duc . . . . .	33
2.5.7	Mystérieux R duc /R duc Conflits . . . . .	35
2.6	Error récupération . . . . .	37
2.7	Debugger son parser . . . . .	38
2.8	Exécuter Ocamlyacc . . . . .	39
2.8.1	Ocamlyacc Options . . . . .	39
<b>3</b>	<b>Astuces</b>	<b>40</b>

# Liste des tableaux

1.1	Tableau des motifs . . . . .	8
1.2	Quelques fonctions du module L xing . . . . .	11
1.3	Options d'Ocaml x . . . . .	13

# Table des figures

1.1	Ligne de commande . . . . .	7
1.2	Structure du fichier d'entrée ".ml1" . . . . .	8
1.3	Importance de l'ordre (1) . . . . .	9
1.4	Importance de l'ordre (2) . . . . .	9
1.5	Plus long motif reconnu . . . . .	10
1.6	Suppression de "zap m " . . . . .	10
1.7	Compression . . . . .	10
1.8	L(s) point(s) d'entrée . . . . .	11
1.9	Règles conditionnelles . . . . .	12
1.10	Utilisation de token lex - yacc . . . . .	12
2.1	Simplification de fonction C traduite en tokens . . . . .	15
2.2	Règle de somme . . . . .	17
2.3	Structure du fichier de grammaire Ocamlyacc . . . . .	18
2.4	Structure de règle de grammaire Ocamlyacc . . . . .	20
2.5	Exemple de la combinaison par PLUS . . . . .	20
2.6	Plusieurs règles . . . . .	20
2.7	Règle vide . . . . .	21
2.8	Règle récursive - virgules . . . . .	21
2.9	Règle récursive droite . . . . .	21
2.10	Règles mutuellement récursives . . . . .	22
2.11	Action . . . . .	23
2.12	Type de positions . . . . .	23
2.13	Initialiser les positions du lex/pars . . . . .	24
2.14	Position . . . . .	24
2.15	Fonctions de position . . . . .	24
2.16	Fonctions pour le membre gauche . . . . .	24
2.17	Exemple d'utilisation des positions . . . . .	25
2.18	Déclaration de type de token . . . . .	25
2.19	Syntaxe des déclarations de précédence . . . . .	26
2.20	Appel des fonctions . . . . .	27
2.21	Token en avant . . . . .	29
2.22	Conflits shift/réduire . . . . .	29
2.23	Entrées équivalentes si shift . . . . .	30
2.24	Grammaire à conflit . . . . .	30
2.25	Grammaire à conflit soluble par précédence . . . . .	31
2.26	Exemple de grammaire . . . . .	31
2.27	Définition de plusieurs précédences gauches . . . . .	32
2.28	Plusieurs précédences groupées . . . . .	32

2.29	Définition d s précéd nc s d cont xt . . . . .	33
2.30	Err ur d grammair provoquant un conflit r duc /r duc . . . . .	34
2.31	Résolution du conflit r duc /r duc . . . . .	34
2.32	Ex mpl 2 d'un conflit r duc /r duc . . . . .	34
2.33	Résolution 1 d l' x mpl 2 d'un conflit r duc /r duc . . . . .	35
2.34	Résolution 2 d l' x mpl 2 d'un conflit r duc /r duc . . . . .	35
2.35	Conflit r duc /r duc non dét etabl . . . . .	36
2.36	Résolution d'un conflit r duc /r duc non dét etabl . . . . .	36
2.37	Résolution d'un conflit r duc /r duc non dét etabl - 2èm possibilité . . . . .	37
2.38	Rattrapag d' rr ur . . . . .	37
2.39	Stratégi d récupération d' rr ur . . . . .	38
2.40	Err ur causé ntr d ux délimit urs . . . . .	3800(.)T0_g36

# Chapitre 1

## Ocamllex

Ocamllex est un outil d'OCaml pour générer des *scanners* : programmes qui reconnaissent des motifs lexicaux à l'aide de *descriptions* sous forme de paire (*expression régulière, règle*). Ocamllex génère un exécutable qui effectue l'analyse lexicale tout en générant le code OCaml correspondant.

### 1.1 Commandes et structure

La commande pour lancer l'analyse est :

```
ocamllex program.mll
```

FIG. 1.1 – Ligne de commande

### 1.2 Input

Le fichier d'entrée (fig 1.2) d'Ocamllex contient 4 sections : *header*, *definitions*, *rules* et *trailer*. Son extension est ".mll".

Il faut remarquer que :

- *header* et *rules* sont obligatoires
- *header* et *trailer* sont accolés, code OCaml qui sera remis "tel quel" dans le fichier OCaml de sortie.
- *definitions* contiennent les déclarations des expressions régulières et leurs identifiants.
- *rules* contiennent des points d'entrée pour définir l'effet de l'analyse lexicale.

Les patterns sont décrits dans la section suivante 1.3.

### 1.3 Patterns

Les motifs ou *patterns* ont le format d'expressions régulières dans le style de lex à l'aide d'une syntaxe Caml-like. La table 1.1 répertorie les expressions régulières.



```

(* header section *)
{ h a d r }

(* definitions section *)
let id nt = r g xp
let ...

(* rules section *)
rul ntrypoint [arg1... argn] = pars
  | patt rn { action }
  | ...
  | patt rn { action }
and ntrypoint [arg1... argn] = pars
  ...
and ...

(* trailer section *)
{ trail r }

```

FIG. 1.2 – Structur du fichi r d' ntré ".ml1"

'c'	match th charact r 'c'
—	match any charact r
of	match an nd-of-fil
"foo"	match th lit ral string "foo"
['x' 'y' 'z']	match s ith r an 'x', a 'y', or a 'z'
['a' 'b' 'j'—'o' 'Z']	match s an 'a', a 'b', any l tt r from 'j' through 'o', or a 'Z'
[^ 'A'—'Z']	a "n gat d charact r s t"
[^ 'A'—'Z' 'n']	any charact r EXCEPT an upp rcas l tt r or a n wlin
r*	z ro or mor r's, wh r r is any r gular xpr ssion
r+	on or mor r's, wh r r is any r gular xpr ssion
r?	"an optional r"
id nt	th xpansion of th "id nt"
(r)	par nth s s ar us d to ov rrid pr c d nc (s b low)
rs	"concat nation"
r s	ith r an r or an s
r#s	match th diff r nc of th two sp cifi d charact r s ts
r as id nt	bind th string match d by r to id ntifi r id nt

TAB. 1.1 – Tabl au d s motifs

L s xpr ssions régulièr s cité s dans la tabl 1.1 sont trié s s lon l s priorités d précéd nc du plus fort au plus faibl : '\*' t '+' ont la plus fort précéd nc , puis '?', 'concat nation', '|', t nfin 'as'.

Par x mpl , "foo" | "bar"\* = ("foo")|("bar"\*) = "foo" ou z ro-ou-plusi urs "bar", t ("foo|"bar")\* = z ro-ou-plusi urs "foo"-ou-"bar" ( x : "foobarbarbarfooobar", "", "barbar"...)

Remarque : l'expression n'wline peut être reconnue par "[^'A'-'Z']" à moins que "n" soit explicite ment présente dans la négation : "[^'A'-'Z'\n]" (non standard).

## 1.4 Match

Le principe est le suivant : si il y a plusieurs patterns qui correspondent à l'expression saisie, le principe du "plus long match" est appliqué. C'est-à-dire qu'entre les patterns, "ding", "dong" et "dingdong", l'expression saisie sera associée au pattern "dingdong". S'il y a plusieurs patterns de longueur égale, on prendra le premier de la liste.

Lorsqu'un motif correspondant est trouvé, le texte correspondant, appelé *token* est considéré comme disponible sous la forme d'un chaîne de caractères. L'action correspondant au motif filtré est alors exécutée (voir description des actions dans la section 1.5) et l'input est scanné pour un motif suivant.

Si aucun correspondant n'est trouvé, le scanner lève l'exception "lexing : empty token".

Voici quelques exemples qui montrent le fonctionnement du filtrage.

Les deux premiers cas (1.3, 1.4) illustrent qu'il y a deux motifs peuvent être reconnus, Ocamll choisit le premier de la liste. Ainsi, dans le deuxième cas (1.4), le pattern "ding" est donc inutile.

```
(* rules section *)
rul tok n = pars
| "ding"          { print_ ndlin  "Ding" }                (* "ding"
pattern *)
| ['a'-'z']+ as word                (* "word"
pattern *)
{ print_ ndlin  ("Word:~" ^ word) }
...
```

FIG. 1.3 – Importance de l'ordre (1)

```
(* rules section *)
rul tok n = pars
| ['a'-'z']+ as word                (* "word" pattern *)
{ print_ ndlin  ("Word:~" ^ word) }
| "ding"          { print_ ndlin  "Ding" }                (* "ding"
pattern *)
| ...
```

FIG. 1.4 – Importance de l'ordre (2)

Dans ce troisième exemple (1.5), on choisit le plus long motif reconnu. Il y a trois patterns : ding, dong et dingdong.

```

(* rules section *)
rul tok n = pars
| "ding"      { print_ ndlin "Ding" }          (* "ding"
pattern *)
| "dong"      { print_ ndlin "Dong" }          (* "dong"
pattern *)
| "dingdong"  { print_ ndlin "Ding-Dong" }      (* "dingdong" pattern
*)
...

```

FIG. 1.5 – Plus long motif r connu

Lorsqu "dingdong" st donné n input, deux choix sont possibles : soit l s motifs ding + dong soit l motif dingdong. Par l princip du plus long motif r connu, l motif dingdong st choisi.

Il st possible d'appliqu r l princip du plus court motif r connu n r mplaçant l mot clé pars par `shortst`. L princip du pr mi r d la list st conservé.

## 1.5 Action

Chaqu motif dans un règl a un action correspondant , qui st un xpr ssion Ocaml. Par x mpl , voici un programm qui ffac tout s l s occurr nc s d "zap m " :

```

(* rules section *)
rul tok n = pars
| "zap_m"     { tok n l xbuf }                 (* ignore this token: no
processing and continue *)
| _ as c      { print_char c; tok n l xbuf }

```

FIG. 1.6 – Suppr ssion d "zap m "

Ex mpl d'un programm qui compr ss l s spac s t tabulations multipl s n un s ul spac , t qui supprim un spac n fin d lign :

```

{}
rul tok n = pars
| [ ' ' '\t' ]+ { print_char ' '; tok n l xbuf }
| [ ' ' '\t' ]+ '\n' { tok n l xbuf }          (* ignore this token *)

```

FIG. 1.7 – Compr ssion

L s actions p uv nt inclur du cod Ocaml qui r nvoi un val ur. Chaqu fois qu la fonction d'analys ur l xical st app lé , ll continu d s'appliqu r sur l s *tokens* d puis l squ ls il a analysé n d rni r jusqu'à att indr la fin du fichi r.

R marqu : c qu fait l'analys ur, c' st simpl m nt d vérifi r qu l'input vérifi l s règl s l xical s, il n'y a pas d "résultat" à propr m nt parlé, mais simpl m nt, n cas d succès, l'assuranc qu l'analys ur r connaît l'input. L s actions p rm tt nt d'agir sur l s mots r connus au cours d l'analys .

Les actions sont évaluées après que la *lexbuf* soit associée au buffer lexical courant et à l'identifiant qui suit le mot clé (ou à la chaîne de caractères correspondante). L'utilisation du *lexbuf* est fournie par la librairie standard du module *Lxing* dont la table 1.2 présente quelques extraits.

<code>Lxing.lexm lexbuf</code>	Valeur de la chaîne reconnue
<code>Lxing.lexm _char lexbuf n</code>	nième caractère de la chaîne reconnue (début à 0)
<code>Lxing.lexm _start lexbuf</code>	Position de la chaîne depuis le début de l'input (début à 0)
<code>Lxing.lexm _end lexbuf</code>	Position de la chaîne depuis la fin de l'input
<code>Lxing.lexm _start_pos lexbuf</code>	Position de type position (voir documentation Ocaml 3.08)
<code>ntrypoint [xp1... xpn] lexbuf</code>	Appelle un autre lexer sur le point d'entrée donné. lexbuf doit être le même argument

TAB. 1.2 – Quelques fonctions du module *Lxing*

## 1.6 Output : le scanner généré

La sortie *output* générée est un fichier ".ml" du même nom que celui invoqué par `Ocamllex`. Le fichier généré contient les fonctions de scanner, des tables utilisées par le scanner pour faire correspondre les tokens et des petites fonctions auxiliaires. Les fonctions de scanner sont déclarées de la façon suivante (1.8) :

```
let ntrypoint [arg1... argn] lexbuf =
  ...
and ...
```

FIG. 1.8 – Les point(s) d'entrée

où la fonction a  $n + 1$  arguments. Les  $n$  arguments proviennent de la définition des règles. La fonction scanner résultante nécessite un argument supplémentaire, appelé *lexbuf* de type *Lxing.lexbuf* qui doit être le même argument.

Lorsqu'un point d'entrée est appelé, il scanne les tokens à partir du l'argument *lexbuf*. Quand il trouve un motif correspondant, il exécute les actions correspondantes et retourne le *lexbuf* (amputé du mot reconnu). De façon à continuer l'analyse lexicale après l'évaluation d'une action, il faut appeler la fonction de scanner récursivement.

## 1.7 Démarrage

Il est possible d'activer les règles de façon conditionnelle. Quand on veut activer l'autre règle, il suffit d'appeler le point d'entrée de l'autre fonction. Par exemple, l'input suivant est constitué de deux règles, l'une vérifiant les tokens et l'autre "sautant" les commentaires (1.9).

```

{}
rul tok n = pars
| [ ' ' '\t' '\n' ]+
    (* skip spaces *)
    { tok n l xbuf }
| "("
    (* activate "comment" rule *)
    { comm nt l xbuf }
...
and comm nt = pars
| ")"
    (* go to the "token" rule *)
    { tok n l xbuf }
| -
    (* skip comments *)
    { comm nt l xbuf }
...

```

FIG. 1.9 – Règles conditionnelles

Lorsqu'un scanner génère et rencontre un commentaire ouvrant "(" à la règle `tok n`, il active l'autre règle de commentaire. Lorsqu'il rencontre un commentaire fermant ")" dans la règle `comm nt`, il retourne à la règle `tok n`.

## 1.8 Interface avec Ocaml yacc

L'un des principaux usages de Ocamllex est l'association avec son compagnon Ocaml yacc pour la génération de parsers. Les parsers d'Ocaml yacc appellent l'un des fonctions de scanner pour trouver le token d'entrée suivant. La fonction sélectionne et renvoie le type du token suivant avec sa valeur associée. Pour utiliser Ocamllex avec Ocaml yacc, les fonctions de scanner doivent utiliser un module de parser pour définir les types de tokens, qui sont définis dans les attributs '%token' apparaissant dans l'input de Ocaml yacc. Par exemple, si le fichier d'entrée de Ocaml yacc est `parser.mly` et que l'un des tokens est "NUMBER", un parti du scanner devra ressembler à :

```

{
  open Pars
}

rul tok n = pars
...
| ['0'-'9']+ as num { NUMBER (int_of_string num) }
...

```

FIG. 1.10 – Utilisation tok n l x - yacc

## 1.9 Options

Voici quelques options pour Ocamllex (tableau 1.3) :

<code>-o output-fil</code>	Chang l nom d la sorti
<code>-ml</code>	L'automat st codé av c d s fonctions ocaml plutôt qu'un automat built-in. Pour l débuggag .
<code>-q</code>	Pour supprim r l s m ssag s d'information.

TAB. 1.3 – Options d'Ocamll x

# Chapitre 2

## Ocamlyacc

### 2.1 Introduction

Ocamlyacc est un générateur de parsers qui convertit un grammaire (type LALR(1)) en un programme Ocaml qui parse cette grammaire. Quand on maîtrise Ocamlyacc, on peut l'utiliser pour construire des parsers d'un large gamm de langages, de la simple calculatrice au langage de programmation complet.

Ocamlyacc est très proche de yacc (bison) qui sont répandus dans les environnements de programmation C. La maîtrise d'ocaml est nécessaire pour l'utiliser.

Cet tutoriel est fait de chapitres simples qui expliquent les concepts de base et présentent quelques exemples.

### 2.2 Concepts

#### 2.2.1 Langages et grammaires hors-contexte (type 2)

Pour qu'Ocamlyacc parse un langage, le langage doit être décrit par un grammaire hors-contexte (type 2). Ce qui veut dire qu'on doit spécifier un ou plusieurs groupes syntaxiques et donner les règles de constructions correspondantes. Par exemple, le langage C, un sort de groupe syntaxique est appelé "expression". Une règle pour construire une expression peut être "une expression peut être construite à l'aide du signe - et d'une autre expression". Une autre peut être "une expression peut être un entier". Comme on peut le voir dans ces exemples, les règles sont souvent récursives, mais il doit y avoir au moins une règle qui mène hors de la récursion.

Le système formel le plus commun pour présenter les règles de façon lisible est la norme Backus-Naur Form ou "BNF". Tout grammaire exprimée en BNF est un grammaire hors-contexte. L'entrée d'Ocamlyacc est essentiellement BNF.

Seuls les grammaires LALR(1) peuvent être traitées par Ocamlyacc. Pour simplifier, on doit pouvoir parser n'importe quel mot de l'entrée avec un seul token d'avance, ce qui est une description d'un grammaire LR(1), les grammaires LALR(1) rajoutent des restrictions plus complètes à expliquer, mais il est rare qu'un grammaire LR(1) ne soit pas LALR(1). Pour plus d'information à ce propos, il faut se référer à la sous-section 2.5.7.

Dans une règle de grammaire formelle pour un langage, chaque unité ou groupe syntaxique est désigné par un symbole. Les symboles non-terminaux désignent ceux qui sont construits par des unités ou groupes plus petits, les symboles terminaux (ou tokens) ne peuvent pas être subdivisés. Un morceau/mot de l'entrée qui correspond à un symbole terminal est un token, tandis que s'il correspond à un non-terminal, il se dénomme un groupe /groupement.

En utilisant le langage C comme exemple, voici quelques symboles terminaux et non-terminaux vus précédemment. Les tokens de C sont les identifiants, les constantes (numériques et chaînes de caractères), les différents mots clés, les opérateurs arithmétiques et les ponctuations. La grammaire de C inclut les symboles terminaux suivants : "identifiant", "nombre", "string", un symbole par mot clé, opérateur ou ponctuation ("if", "return", "const", "static", "int", "char", "plus-sign", "op-n-brac", "comma"...).

La figure 2.1 présente une simple fonction C subdivisée en tokens :

```
int          /* mot clé "int" */
sqr (int x)  /* identifiant, op-n-brac, identifiant, identifiant,
  clos-brac */
{           /* op-n-brac */
  return x * x; /* mot clé "return", identifiant, ast-risk, identifiant,
    , s-micolon */
}          /* clos-brac */
```

FIG. 2.1 – Simple fonction C traduite en tokens

Les groupes syntaxiques de C incluent l'expression, la déclaration, la définition de fonction. Ils peuvent être représentés dans la grammaire de C par les symboles non-terminaux "expression", "statm", "déclaration" et "fonctiondéfinition". La grammaire complète utilise différents constructeurs supplémentaires, chacun ayant son symbole non-terminal de façon à exprimer ces quatre groupes syntaxiques. L'exemple précédent (figure 2.1) est une "fonctiondéfinition", composée d'une "déclaration" et d'un "statm". Dans ce "statm", chaque  $x$  est une "expression" ainsi que  $x * x$ .

Chaque symbole non-terminal doit avoir une règle grammaticale montrant comment il est construit à partir de plusieurs règles plus simples. Par exemple, un statm en C est `return`; il serait décrit avec une règle de grammaire qui pourrait s'écrire "un 'statm' peut être construit d'un 'return', un 'expression' et un point virgule ('s-micolon')". Et ainsi pour tous les statms possibles en C.

Il faut distinguer un symbole non-terminal : celui qui énonce le langage complet. Il s'appelle le symbole de départ (start symbol). Dans un compilateur, cela signifie l'entrée complète du programme. En C, le symbole non-terminal "séquence de définitions et déclarations" joue ce rôle.

Par exemple, `1 + 2` est une expression C valide (une partie valide d'un programme C), mais n'est pas valide en tant que programme C entier. Dans la grammaire hors-contexte de C, `1 + 2` qui est une expression n'est pas le symbole de départ.

Le parser Ocamlyacc, lit une séquence de tokens comme les entrées et regroupe les tokens en utilisant les règles de grammaire. Si l'entrée est valide, la totalité des tokens est réduite à un seul groupe dont le symbole est le symbole de départ de la grammaire. Si l'entrée est invalide, le parser renvoie un erreur. En C, l'entrée doit être une "séquence de définitions et déclarations".

## 2.2.2 Des règles formelles à l'entrée Ocamlyacc

Une grammaire formelle est une construction mathématique. Pour définir le langage pour Ocamlyacc, il faut créer un fichier exprimant la grammaire dans la syntaxe Ocamlyacc : un fichier de grammaire Ocamlyacc (voir la section sur le format des fichiers Ocamlyacc 2.3)

Un symbole non-terminal d'une grammaire formelle est représenté dans l'input Ocamlyacc par un identifiant, comme en Ocaml. C'est un symbole Caml, excepté qu'il ne peut pas finir par



’ t qu’il doit commencer par un minuscule (  $x$  : `xpr, stmt, déclaration` ).

La représentation Ocamlyacc d’un symbole terminal est aussi appelée type de token. Les types de token doivent être déclarés dans la section déclaration d’Ocamlyacc et ils doivent être ajoutés comme constructeurs pour des types concernés de token. En tant que constructeurs, ils doivent commencer par un majuscule (  $x$  : `Int gr` , `Identifieur` , `IF` ou `RETURN` ). Le symbole terminal d’erreur est réservé pour la récupération d’erreurs (voir la section sur les symboles 2.3.2).

Les règles de grammaire ont également un expression syntaxe Ocamlyacc. Dans l’exemple suivant, une règle Ocamlyacc possible pour le retour de `C` :

```
stmt: RETURN xpr SEMICOLON ;
```

Voir la section sur la syntaxe des règles de grammaire 2.3.3.

### 2.2.3 Valeurs sémantiques

La grammaire forme la sélection des tokens seulement d’après leur classification. Par exemple, sur une règle mentionnant le symbole terminal ‘int gr constant’, cela signifie qu’il n’importe que l’entier constant soit valide pour cette position. La valeur précise de la constante est ignorée lors du parsing de l’entrée : si  $x + 4$  est grammaticalement correct, alors  $x + 1$  ou  $x + 3989$  est également grammaticalement correct.

Par contre, la valeur précise est très importante pour ce que l’entrée signifie une fois parsée. Un compilateur est inutile s’il échoue à distinguer entre les constantes 4, 1 et 3989 dans le programme. C’est pourquoi chaque token de la grammaire Ocamlyacc possède un type de token et une valeur sémantique. Voir la section sur la sémantique 2.3.5.

Un type de token est un symbole terminal défini dans la grammaire, tel que `INTEGER`, `IDENTIFIER` ou `SEMICOLON`. Il contient les informations nécessaires pour décider où la validité du token peut apparaître et comment le grouper avec les autres tokens. Les règles de grammaire ne considèrent que les types de tokens, pas les tokens eux-mêmes.

La valeur sémantique contient le reste de l’information à propos des tokens, par exemple la valeur d’un entier, le nom d’un identifieur. Attention, un token comme `SEMICOLON` n’a pas de valeur sémantique.

Par exemple, un token d’entrée peut être classé comme un token de type `INTEGER` et avoir une valeur sémantique 4. Un autre token d’entrée peut avoir le même type de token `INTEGER` mais la valeur 3989. Lorsqu’une règle de grammaire indique qu’un `INTEGER` est autorisé, ces deux tokens sont acceptables parce qu’ils sont tous deux des `INTEGER`. Lorsqu’il pars une acceptation de token, il garde une trace de la valeur sémantique.

Chaque groupe peut aussi avoir une valeur sémantique comme son symbole non-terminal. Par exemple, dans un calculatrice, une expression a typiquement une valeur sémantique qui est un nombre. Dans un compilateur, pour un langage de programmation, une expression a typiquement une valeur sémantique qui est une structure d’arbre décrivant les nœuds de l’expression.

### 2.2.4 Actions sémantiques

Pour qu’un programme soit utile, il doit être plus que correct à parser. Il doit aussi produire des sorties basées sur les entrées... Dans une grammaire Ocamlyacc, une règle de grammaire peut avoir une action associée écrite en Ocaml. Chaque fois qu’il pars une reconnaissance pour une règle, une action est exécutée. Voir la sous-section sur les actions 2.3.5.

La plupart du temps, l’objectif de l’action est de calculer la valeur sémantique globale de la construction globale à partir des sous-parties. Par exemple, si on a une règle qui dit qu’une expression peut être la somme de deux expressions, lorsqu’il pars une reconnaissance, il

chaque sous-*xpr* ssion possède un val ur sémantiqu e qui décrit comment il a été construit . L'action d'un t ll règl e vrait cré r un val ur du mêm e g nr pour l' *xpr* ssion r connu .

Par x mpl , voici figur 2.2 un règl e qui dit qu'un *xpr* ssion p ut êtr e la somm e d d ux *xpr* ssions :

```
xpr :  xpr PLUS xpr    { $1 + $3 }
      ;
```

FIG. 2.2 – Règl e d somm

L'action indiqu e comment produire la val ur sémantiqu e d la somm e d s *xpr* ssions à partir d s val urs sémantiqu s d s d ux sous- *xpr* ssions.

### 2.2.5 Positions

D nombr us s applications, t ls int rprét urs ou compilat urs, doiv nt produire d s m s- sag s d' rr urs utilis t détaillés. Afin d fair e la, on doit êtr e capabl e d'id ntifi r la position dans un t xt d chaqu e construct ur syntaxiqu e. Ocaml yacc p rm t d fair e la.

Chaqu e tok n a un val ur sémantiqu e t un position associé . Mais l typ e d position st l mêm e pour tous l s tok ns t group s. D plus, la sorti du pars r st fait d structur d donné s p rm ttant d stock r l s positions. Voir la sous-s ction position 2.3.6.

Many applications, lik int rpr t rs or compil rs, hav to produc v rbo s and us ful rror m ssag s. To achi v this, on must b abl to k p track of th t xtual position, or location, of ach syntactic construct. Ocaml yacc provid s a m chanism for handling th s locations. On p ut utilis r c s positions pour att indr l s actions utilisant l s fonctions du modul Parsing.

### 2.2.6 Sortie d'Ocaml yacc : le fichier du parser

Lorsqu'on xécut Ocaml yacc , on donn e un fichi r d grammair Ocaml yacc comm ntré . La sorti st un fichi r Ocaml qui pars e l langag e décrit par la grammair . C fichi r st app lé un pars ur Ocaml yacc . Att ntion Ocaml yacc st l'outil dont la sorti st l pars ur Ocaml yacc

La fonction du pars ur Ocaml yacc st d r group r l s tok ns dans l s group s s lon l s règl s d grammair . Par x mpl , d construire id ntifi r t opérat urs n *xpr* ssion. Lorsqu'il fait e r group m nt, il ff ctu e l s actions d s règl s d grammair s qu'il utilis e.

L s tok ns vi nn nt d'un fonction app lé l'analys ur l xical qu l'on doit fournir d'un façon ou d'un autr . L pars ur Ocaml yacc app ll l'analys ur l xical chaqu fois qu'il a b so in d'un nouv au tok n. Il n sait pas c qu conti nt l s tok ns, (bi n qu l ur val ur sémantiqu e p ut êtr e récupéré ). Typiqu m nt, l'analys ur l xical fabriqu nt l s tok ns n parsant l s caractèr s d'un t xt , tandis qu'Ocaml yacc n dép nd pas d ça. Voir la s ction 2.4.2 sur la fonction d'analys l xical .

L fichi r pars ur d Ocaml yacc st n Ocaml. Il définit l s fonctions qui implém nt nt la grammair . Chaqu fonction d' ntré du cod Ocaml généré st nommé d'après l s symbol s d départ du fichi r d grammair . L s fonctions n constitu nt pas un programm Ocaml compl t : il faut fournir qu lqu s fonctions supplém ntair s. La pr mièr st l'analys ur l xical qui doit êtr e donné comm argum nt d la fonction d' ntré du pars ur. Un autr fonction st la fonction d rapport d' rr ur app lé par l pars ur afin d rapport r un rr ur. Et nfin, un programm compl t Ocaml doit pouvoir app ll r un (ou plusi urs) fonction d' ntré généré ou alors l pars ur n pourra pas d' xécut r. Voir la s ction 2.4 sur l'int rfac e du pars ur.

## 2.2.7 Les étapes pour utiliser Ocaml yacc

Voici les étapes de construction d'un langage utilisant Ocaml yacc, partant de la spécification de la grammaire allant jusqu'au compilateur ou interpréteur :

- Former le langage spécifique de la grammaire sous une forme connue par Ocaml yacc (voir la section sur les fichiers de grammaire 2.3). Pour chaque règle dans le langage, décrire l'action qui doit être effectuée lorsqu'une instance de cette règle sera connue. L'action est décrite par un séquence de statements Ocaml.
- Écrire un analyseur lexical pour analyser l'entrée et transmettre les tokens au parser. L'analyseur lexical peut être écrit à la main en Ocaml (voir la section 2.4.2). Il peut aussi être produit par Ocamllex, voir la partie 1.
- Écrire une fonction de contrôle qui applique le parser produit par Ocaml yacc
- Écrire les fonctions de rapport d'erreurs.

Pour transformer le code exécutable, il faut suivre les étapes suivantes :

- Exécuter Ocaml yacc sur la grammaire pour produire le parser.
- Compiler le code sorti d'Ocaml yacc ainsi que les autres fichiers sources.
- Lier les fichiers objets pour produire le produit fini.

## 2.3 Le fichier de grammaire Ocaml yacc

Ocaml yacc prend en entrée une spécification de grammaire hors-contexte et produit une fonction Ocaml qui reconnaît les instances correctes de cette grammaire. Le fichier d'entrée de grammaire Ocaml yacc a par convention une extension en ".mly". Voir la section Invoquer Ocaml yacc 2.8.

### 2.3.1 Structure du fichier

Le fichier de grammaire Ocaml yacc est constitué de quatre sections principales, présentées ici (figure 2.3) avec leurs délimiteurs :

```
%{  
    Header — Ocaml declarations (Ocaml code)  
}%  
    Ocaml yacc declarations  
%%  
    Grammar rules  
%%  
    Trailer — Additional Ocaml code (Ocaml code)
```

FIG. 2.3 – Structure du fichier de grammaire Ocaml yacc

- Les `%%`, `%{` et `%}` sont des ponctuations qui permettent de séparer les sections.
- Dans la section *header* on définit les types, variables et fonctions utilisés dans les actions.
- Dans la section *Ocaml yacc declarations*, on déclare les noms des symboles terminaux et non-terminaux. On peut aussi décrire les précédents des opérateurs et les types des données des valeurs sémantiques des différents symboles.
- La section *grammar rules* définit comment construire chaque non-terminal à partir des symboles terminaux.
- La section *Trailer* peut contenir du code Ocaml.

Par défaut, les commentaires sont encadrés entre "/" et "/" (comme en C) excepté dans

La valeur renvoyée par la fonction `lexer` est toujours l'un des symboles terminaux. Chaque type de token contient une valeur de type variant dans le fichier `parser`, afin que la fonction `lexer` puisse retourner un.

Parce que la fonction `lexer` est définie dans un fichier séparé, il faut faire en sorte que les définitions de types de tokens soient disponibles à ce niveau. Après avoir invoqué `ocamlyacc filename.mly`, le fichier `"filename.mli"` qui est généré contiendra les définitions de types de tokens. Il est utilisé dans la fonction `lexer`.

Le symbole d'erreur est un symbole terminal réservé pour récupérer les erreurs (voir la section 2.6 sur le rattrapage d'erreur), il doit rester réservé à cet usage.

### 2.3.3 Syntaxe des règles de grammaire

Une règle de grammaire Ocamlyacc a la forme générale suivante (figure 2.4) :

```

r_suit :
    symbol ... symbol { semantic-action }
    | ...
    | symbol ... symbol { semantic-action }
    ;

```

FIG. 2.4 – Structure d'une règle de grammaire Ocamlyacc

où `r_suit` est le symbole non-terminal que cette règle décrit, `symbol` sont des symboles terminaux ou non-terminals qui sont assemblés par cette règle (voir la section 2.3.2 sur les symboles).

Par exemple la règle suivante (figure 2.5) qui définit deux groupes de type `xp` avec un token `PLUS` au milieu peut être combinée en un groupe plus large de type `xp`.

```

xp :
    xp PLUS xp {}
    ;

```

FIG. 2.5 – Exemple de la combinaison par PLUS

Les espaces dans les règles sont nécessaires seulement pour séparer les symboles, des espaces supplémentaires peuvent être ajoutés.

À la suite des composants de la règle, il doit y avoir l'action qui détermine la sémantique de la règle. Une action doit être définie de cette façon `{ Ocaml code }`, voir la section 2.3.5 sur les actions.

Il est possible d'écrire des règles séparées pour le même résultat, mais on peut aussi les joindre à l'aide d'un barre verticale `"|"` comme dans la figure 2.6. Les règles sont considérées

```

r_suit :
    rule1-symbol ... rule1-symbol { rule1-semantic-action }
    | rule2-symbol ... rule2-symbol { rule2-semantic-action }
    | ...
    ;

```

FIG. 2.6 – Plusieurs règles

comme distincts, même quand ils sont écrits de cette manière.

Si la composition d’une règle est vide, cela veut dire que le résultat peut correspondre à un chaîne de caractères vide. Par exemple, voici comment définir un séquence de zéro séparé de virgules ou de plus d’un group d’xp (figure 2.7) :

```
xps q : /* mpty */ {}
      | xps q1 {}
      ;

xps q1 : xp {}
        | xps q1 COMMA xp {}
        ;
```

FIG. 2.7 – Règle vide

L’usage veut qu’on écrive le commentaire `/* mpty */` pour chaque règle sans composant.

### 2.3.4 Règles récursives

Une règle est appelée récursive si son résultat non-terminal apparaît aussi dans le membre droit de la règle. Pratiquement tout les grammairiens Ocamlacc utilisent la récursion, parce que c’est le seul moyen de définir un séquence d’un nombre indéterminé d’un chaîne de caractères. Voici l’exemple (figure 2.8) de la définition récursive d’un séquence séparé de virgule d’un ou plusieurs xpressions :

```
xps q1 : xp {}
        | xps q1 COMMA xp {}
        ;
```

FIG. 2.8 – Règle récursive - virgules

Cette récursion est appelée récursion gauche car l’utilisation du symbole `xps q1` est le symbole le plus à gauche dans le membre droit de la règle. Voici la même construction avec une récursion droite (figure 2.9) :

```
xps q1 : xp {}
        | xp COMMA xps q1 {}
        ;
```

FIG. 2.9 – Règle récursive droite

Tout sorte de séquence peut être défini en utilisant la récursion gauche ou droite, mais il est préférable d’utiliser toujours la récursion gauche car cela utilise un nombre limité d’espace dans la pile. La récursion droite utilise un espace dans la pile Ocamlacc proportionnel au nombre d’éléments de la séquence, parce que tous les éléments doivent être empilés avant que la règle puisse être appliquée au moins une fois. Voir la section 2.5 sur l’algorithme du parcours.

On peut aussi trouver de la récursion mutuelle ou indirecte lorsque le résultat de la règle n’apparaît pas directement dans le membre droit, mais qu’il apparaît dans des règles d’autres non-terminals qui apparaissent eux dans le membre droit de cette règle.

Par exemple dans la figure 2.10 suivante, on définit deux non-terminaux mutuellement récursifs, chacun faisant référence à l'autre.

```
xpr:      primary {}
        | primary PLUS primary {}
        ;

primary:   constant {}
        | LPAREN xpr RPAREN {}
        ;
```

FIG. 2.10 – Règles mutuellement récursives

### 2.3.5 Définir la sémantique d'un langage

Les règles de grammaire d'un langage ne déterminent que la syntaxe. La sémantique est déterminée par les valeurs sémantiques associées aux différents tokens et par les actions à effectuer lorsque les différents groupes sont reconnus.

Par exemple, un calculateur calculerait ment parce que la valeur associée à chaque expression est le nombre correct. Il ajouterait ment parce que l'action du groupe  $x + y$  est d'ajouter les nombres associés à  $x$  et  $y$ .

#### Type de données des valeurs sémantiques

Dans un simple programme, il peut être suffisant d'utiliser le même type de données pour les valeurs sémantiques de tous les constructeurs du langage. Dans la plupart des programmes, il est cependant nécessaire d'avoir des types de données différents pour des types de tokens ou des groupes différents. Par exemple, un constant numérique peut nécessiter le type `int` ou `float` alors qu'un constant chaîne de caractères ou un identifiant peut nécessiter le type `string`.

Afin d'utiliser plus d'un type de données pour les valeurs sémantiques dans un parser, Ocaml yacc impose de choisir un type de données pour chaque symbole (terminal ou non-terminal) pour lequel la valeur sémantique est utilisée. Pour les tokens, cela se fait à l'aide de `%token` Ocaml yacc d'claration (voir la sous-section 2.3.7 sur les noms de type de tokens) et pour les groupes avec `%nonterminal` (voir la sous-section 2.3.7 sur les symboles non-terminaux).

#### Actions

Une action est associée à une règle syntaxique et contient du code Ocaml qui est exécuté chaque fois qu'une instance de la règle est reconnue. La tâche de la plupart des actions est de calculer une valeur sémantique pour le groupe construit par la règle à partir des valeurs sémantiques des tokens et des sous-groupes qui constituent la règle.

Une action est construite de manière statique en Ocaml encadrée par des accolades. Les règles n'ont qu'une action à la fin de la règle suivie par tous les composants.

```

xp :      ...
      |   xp PLUS  xp { $1 +. $3 }

```

FIG. 2.11 – Action

Cette règle construit un `xp` à partir de deux groupes `xp` plus petits connectés avec le token du signe "+". Dans l'action `$1` et `$3` font référence aux valeurs sémantiques des deux composants des groupes `xp`, qui sont le premier et le troisième symboles du membre droit de la règle. La somme est renvoyée afin que la valeur sémantique de l'expression d'addition soit reconnue par la règle. S'il y avait un valeur sémantique utilisable avec le token PLUS, on pourrait y référer avec `$2`.

### 2.3.6 Suivre les positions

Bien que les règles de grammaire et les actions sémantiques soient suffisantes pour écrire un parser entièrement fonctionnel, il peut être utile de pouvoir rajouter des informations complémentaires, en particulier les symboles de positions.

La façon dont les positions sont traitées est définie en fournissant le type de données et les actions à effectuer quand les règles sont appliquées.

#### Type de données des positions

Le contenu de cette section est valable depuis Ocaml 3.08.

Le type de données pour les positions est construit comme ceci (figure 2.12) :

```

type position = {
  pos_fnam  : string;           (* file name *)
  pos_lnum  : int;              (* line number *)
  pos_bol   : int;              (* the offset of the beginning of the
                                line *)
  pos_cnum  : int;              (* the offset of the position *)
}

```

FIG. 2.12 – Type de positions

La valeur du champ `pos_bol` est le nombre de caractères entre le début du fichier et le début de la ligne, la valeur du champ `pos_cnum` est le nombre de caractères entre le début du fichier et la position.

Le mot lexical parvient seulement au champ `pos_cnum` du `lxbuf.lx_curr_p` avec le nombre de caractères lus à partir du début du lxbuf. Le programmeur est donc responsable de l'exactitude des autres champs. Avant d'utiliser la position dans le parser, il faut initialiser `Lxing.lxbuf.lx_curr_p` correctement dans le lexer, en utilisant une fonction comme celle-ci (figure 2.13) :

#### Actions et positions

Les actions sont utilisées non seulement pour définir la sémantique du langage mais aussi (à l'aide des positions) pour décrire le comportement du parser. Le moyen le plus simple de construire des positions de groupes syntaxiques est très similaire à la façon dont les valeurs



```

let incr_lin num l xbuf =
  let pos = l xbuf.L xing.l x_curr_p in
  l xbuf.L xing.l x_curr_p <- { pos with
    L xing.pos_lnum = pos.L xing.pos_lnum + 1;
    L xing.pos_bol = pos.L xing.pos_cnum;
  }
;;

```

FIG. 2.13 – Initialis r l s positions du l x r/pars r

sémantiqu s sont calculé s. Pour un règl donné , plusi urs construct urs p uv nt êtr utilisés pour accéd r aux positions d s élém nts corr spondants. La position du nièm composant du m mbr droit p ut êtr obt nu à l'aid d (figur 2.14) :

```

val Parsing.rhs_start : int -> int
val Parsing.rhs_nd : int -> int

```

FIG. 2.14 – Position

`Parsing.rhs_start n` r nvoi l'écart av c l pr mi r caractère du nièm it m du m mbr droit d la règl . `Parsing.rhs_nd n` r nvoi l'écart av c l d rni r caractère d l'it m. C s fonctions doit êtr app lée s par l s actions. n vaut 1 pour l'it m l plus à gauch t l s pr mi r caractère dans un fichi r a un écart d 0.

On p ut aussi utilis r l s fonctions suivant s (figur 2.15) :

```

val Parsing.rhs_start_pos : int -> L xing.position
val Parsing.rhs_nd_pos : int -> L xing.position

```

FIG. 2.15 – Fonctions d position

(D puis Ocaml 3.08) L s fonctions suivant s r tourn nt un position plutôt qu'un écart (voir sous-sous-ction 2.3.6).

La position du group du m mbr d gauch p ut êtr obt nu par (figur 2.16) :

```

val Parsing.symbol_start : unit -> int
val Parsing.symbol_nd : unit -> int

```

FIG. 2.16 – Fonctions pour m mbr d gauch

`symbol_start ()` r nvoi l'écart av c l pr mi r caractère du m mbr gauch t `symbol_nd ()` l'écart av c l d rni r caractère .

(D puis Ocaml 3.08) L s fonctions suivant s sont comm `symbol_start` and `symbol_nd` , xc pté qu' ll s r nvoi nt un position plutôt qu'un écart (voir sous-sous-ction 2.3.6).

```

val Parsing.symbol_start_pos : unit -> L xing.position
val Parsing.symbol_nd_pos : unit -> L xing.position

```

Voici un exemple basique qui utilise le type donné par défaut des positions (figure 2.17) :

```

xp :      ...
        | xp DIVIDE xp
          { if $3 <> 0.0 then $1 /. $3
            else (
              let start_pos = Parsing.rhs_start_pos 3 in
              let nd_pos = Parsing.rhs_nd_pos 3 in
              printf "%d.%d-%d.%d: _division_by_zero"
                start_pos.pos_lnum (start_pos.pos_cnum -
                  start_pos.pos_bol)
                nd_pos.pos_lnum ( nd_pos.pos_cnum - nd_pos.
                  pos_bol );
              1.0
            )
          }$

```

FIG. 2.17 – Exemple d'utilisation des positions

### 2.3.7 Déclarations Ocamlyacc

La section "déclaration Ocamlyacc" de la grammaire Ocamlyacc définit les symboles utilisés pour formuler la grammaire et les types de données des valeurs sémantiques. Voir la sous-section 2.3.2 sur les symboles.

Tous les types de tokens doivent être déclarés. Les symboles non-terminaux doivent être déclarés s'il est nécessaire de spécifier quel type de données utiliser pour la valeur sémantique (voir la sous-sous-section 2.3.5 sur les types de données des valeurs sémantiques).

La première règle dans le fichier spécifie aussi le symbole de départ, par défaut. Si on veut un autre symbole comme départ, on doit le déclarer explicitement (voir la sous-section 2.2.1 sur les grammaires hors-contexte).

#### Noms des types de token

Le moyen classique pour déclarer un nom de type de token (symbole terminal) est le suivant (figure 2.18) :

```

%token nam ... nam
%token <type> nam ... nam

```

FIG. 2.18 – Déclaration de type de token

Ocamlyacc convertit la sous-forme d'un type dans le parser, de façon à ce que la fonction `lexer` puisse utiliser le nom correspondant au code du type de token.

Si le token a une valeur, l'argument de la déclaration `%token` doit inclure le type de données encadré par les "<>" (voir la sous-sous-section 2.3.5 sur les types de données des valeurs sémantiques).

Par exemple : `%token <float> NUM /* d fin token NUM and its type */`

La parti "typ " doit être une expression de type Caml. La parti **<type>** est copié dans le fichier sorti en ".mli", tous les noms de constructeur de type doivent être valides (par exemple `Modul _nom .typ _nom`).

## Opérateurs de précédences

On utilise les déclarations **%left**, **%right** ou **%nonassoc** pour spécifier les précédences et associativités. Voir la section 2.5.3 sur les conflits résolus par les opérateurs de précédence.

La syntaxe de déclaration de précédence est la suivante (figure 2.19) :

```
%right symbols ... symbols
%nonassoc symbols ... symbols
```

FIG. 2.19 – Syntaxe des déclarations de précédence

Les déclarations de précédence précisent l'associativité et la précédence relative pour tous les symboles :

- L'associativité d'un opérateur "op" déterminant comment interpréter la répétition de l'opérateur : si "x op y op z" est considéré comme "(x op y) op z" ou "x op (y op z)". **%left** spécifie l'associativité gauche "(x op y) op z" et **%right** spécifie l'associativité droite "x op (y op z)". **%nonassoc** ne spécifie pas d'association : "x op y op z" sera considéré comme un erratum de syntaxe.
- La précédence d'un opérateur détermine son comportement lorsqu'il est encapsulé avec d'autres opérateurs. Tous les tokens déclarés dans la même déclaration de précédence ont une précédence équivalente, ce qui signifie qu'ils sont encapsulés selon leur associativité. Lorsque deux tokens déclarés dans deux déclarations de précédence différentes sont associés, celui déclaré plus tard possède la précédence la plus forte, il est donc groupé (au sens syntaxique) en premier.

## Symboles non-terminaux

On peut déclarer le type de chaque symbole non-terminal pour lesquels les valeurs sont utilisées. Cela se fait à l'aide de la déclaration **%type** comme ça :

**%type <type> nonterminal ... nonterminal**. "nonterminal" est le nom d'un symbole non-terminal et "typ " est le nom du type voulu. On peut donner autant de symboles non-terminaux de départ qu'on souhaite dans la même déclaration **%type** s'ils ont le même type. On utilise les espaces pour séparer les noms de symboles.

Cette déclaration est nécessaire pour les symboles de départ. Voir la sous-section 2.3.7 sur les noms de types de tokens.

## Le symbole de départ

On doit déclarer un symbole de départ/démarrage en utilisant la déclaration **%start** de la façon suivante :

```
%start symbol ... symbol
```

Chaque symbole de départ a une fonction de parsing du même nom dans le fichier de sortie. On s'en sert comme point d'entrée pour la grammaire. Rappel : chaque symbole de départ doit avoir un type associé, on utilise alors la déclaration **%type**, voir sous-section précédente 2.3.7 sur les symboles non-terminaux.

## Récapitulatif des déclarations Ocaml yacc

Voici le récapitulatif des déclarations nécessaires pour définir une grammaire :

Here is a summary of the declarations used to define a grammar :

- **%token** Déclaration d'un symbole terminal (nom du type du token) sans précédence ni associativité (voir 2.3.7).
- **%right** Déclaration d'un symbole terminal (nom du type du token) qui a une associativité droite (voir 2.3.7).
- **%left** Déclaration d'un symbole terminal (nom du type du token) qui a une associativité gauche (voir 2.3.7).
- **%nonassoc** Déclaration d'un symbole terminal (nom du type du token) qui n'est pas associatif (pour la syntaxe si c'était le cas, voir 2.3.7).
- **%type** Déclaration du type de la valeur sémantique d'un symbole non-terminal (voir 2.3.7).
- **%start** Précise le symbole de départ de la grammaire (voir 2.3.7).

## 2.4 L'interface du parser

Le parser Ocaml yacc se fait un simple module de fonction Ocaml nommé `s` à partir des symboles de départ de la grammaire (voir sous-section 2.3.7 sur les symboles de départ). Cette section décrit les conventions des fonctions du parser et les autres fonctions nécessaires.

### 2.4.1 La fonction Parser

Afin de faire un parsing, il faut appeler la fonction de parser avec deux paramètres. Le premier paramètre est la fonction d'analyse lexicale de type `L xing.l xbuf -> token` et le second est une valeur de type `L xing.l xbuf`.

Soit `prse` le symbole de départ dans le fichier `prser.mly` et `token` la fonction `lxr` dans le fichier `lexer.mll`, on procède de la manière suivante (figure 2.20) :

```
let l xbuf = L xing.from_channel stdin in
...
let result = Parser.parse L x r.token l xbuf in
...
```

FIG. 2.20 – Appel des fonctions

La fonction `parse` lit les tokens, exécute les actions et s'arrête soit quand elle rencontre un fin d'entrée, soit un erreur de syntaxe non récupérée.

### 2.4.2 La fonction d'analyse lexicale

La fonction d'analyse lexicale est nommée d'après les déclarations de règles. Elle connaît les tokens et puis le flux d'entrée et les renvoie au parser. Ocaml yacc ne crée pas cette fonction automatiquement : il faut l'écrire afin que la fonction de parser puisse l'appeler. On peut assimiler cette fonction à un scanner lexical. Elle est généralement générée par Ocamllex (voir partie 1 ou le chapitre 12 du manuel Ocaml <http://cml.inria.fr/pub/docs/manual-ocaml/manual1026.html>).

### 2.4.3 La fonction de rapport d'erreur

Le parser Ocamlyacc détecte les erreurs de parsing ou de syntaxe lorsqu'il lit un token qui ne satisfait aucun des règles de syntaxe. Une action de la grammaire peut également provoquer une erreur, utilisant la libé `Parsing.Pars _ rror`.

Le parser Ocamlyacc va rapporter l'erreur en appelant une fonction de rapport d'erreur appelée `p_rse_error` qui est optionnelle. Par défaut, la fonction `p_rse_error` ne fait aucune action. Elle est appelée par la fonction de parser lorsqu'un erreur de syntaxe est trouvée et reçoit un argument. Pour une erreur de parsing, l'erreur est généralement `synt xe error`.

La définition suivante est généralement suffisante pour les programmes simples :

```
let pars _ rror s = print _ ndlin s
```

Si la fonction `p_rse_error` renvoie quelque chose à la fonction de parsing, elle essaiera de rattraper l'erreur s'il existe des règles (valides) de grammaire qui récupèrent les erreurs (voir la section 2.6 sur le rattrapage d'erreurs). Si la récupération est impossible, la fonction de parsing lèvera l'exception `Parsing.Pars _ rror`.

## 2.5 L'algorithme du parseur Ocamlyacc

Au fur et à mesure que Ocamlyacc lit les tokens, ils sont empilés avec leur valeur sémantique. La pile est appelée *pile du parser*. On appelle ça du *shifting*.

Par exemple, supposons que la calculatrice infix a lu  $1 + 5 * 3$  et qu'il reste à lire. La pile aura 4 éléments, un par token shifté.

La pile n'a pas toujours un élément pour chaque token. Lorsque les différents tokens sont groupés et shiftés correspondent aux composants d'une règle de grammaire, ils peuvent être combinés suivant cette règle. C'est ce qu'on appelle la *réduction*. Ces tokens sont remplacés dans la pile par un simple groupe dont le symbole est le résultat (membre de gauche) de la règle. Exécuter l'action de la règle fait partie du processus de réduction, c'est-à-dire le calcul de la valeur sémantique du groupe résultant.

Par exemple, si la pile de la calculatrice infix contient  $1 + 5 * 3$  et que le prochain token est un saut de ligne, alors les trois (différents) éléments peuvent être réduits à 15 en appliquant la règle `xpr: xpr MULTIPLY xpr;`. La pile contiendra alors les trois éléments suivants :  $1 + 15$ . Une réduction supplémentaire peut être effectuée, qui produit la valeur 16. Et dès lors le saut de ligne peut être shifté.

Le parser essaie, en shiftant ou réduisant, de réduire l'input complet ment à un groupe simple dont le symbole est le symbole de départ de la grammaire (voir la sous-section 2.2.1 sur les grammaires hors-contexte).

Ce genre de parser est connu dans la littérature comme un parser ascendant.

### 2.5.1 Les tokens en avant

Le parser Ocamlyacc ne réduit pas toujours immédiatement, aussitôt que les différents tokens sont groupés et correspondent à une règle. Une stratégie aussi simple ne couvre pas la plupart des langages. En fait, lorsqu'une réduction est possible, le parser garde "en avant" le token suivant de façon à décider quoi faire.

Quand un token est lu, il n'est pas immédiatement shifté. Il doit d'abord être une *token en avant*, qui n'est pas sur la pile. Le parser peut alors effectuer une ou plusieurs réductions de tokens et groupes sur la pile, tandis que le token en avant reste de côté. Quand plus aucune réduction ne peut être effectuée, le token en avant est shifté sur la pile. Il peut donc rester des réductions possibles, mais selon le type du token en avant, certains des règles peuvent choisir de retarder leur application.

Voici figur 2.21 un x mpl d cas simpl d'utilisation du tok n n avant. C s d ux règl s définiss nt d s xpr ssions cont nant l'opérat ur d'addition t l'opérat ur postfix d factori ll (**FACTORIAL** pour "!") t l'ncadr m nt par par nthès s.

```
xpr :      t rm PLUS  xpr
        |  t rm
        ;

t rm :      LPAREN  xpr RPAREN
        |  t rm FACTORIAL
        |  NUMBER
        ;
```

FIG. 2.21 – Tok n n avant

Supposons qu l s tok ns **1 + 2** sont lus t shifté, qu fair nsuit ?

- Si l tok n suivant st **RPAREN** alors l s trois pr mi rs tok ns doiv nt ètr réduit comm **xpr**. C s ra alors la s ul possibilité, parc qu shift r **RPAREN** produira un séqu nc **t rm RPAREN** t aucun règl n l p rm t.
- Si l tok n suivant st **FACTORIAL**, il doit ètr shifté dir ct m nt afin d réduire **2 FACTORIAL** n **t rm**. Si l pars r voulait réduire avant d shift r **1 + 2** d vi ndrait un **xpr**. On n pourrait donc plus shift r **FACTORIAL** parc qu'on aurait sur la pil un séqu nc d symbol s **xpr FACTORIAL**, qui n' st pas autorisé dans l s règl s.

## 2.5.2 Conflits shift/reduce

Voici un x mpl figur 2.22 d règl s pour un langag ayant d s if-th n t if-th n- ls qu l'on souhait pars r :

```
if_stmt :
        IF  xpr THEN stmt
        |  IF  xpr THEN stmt ELSE stmt
        ;
```

FIG. 2.22 – Conflits shift/r duc

**IF, THEN, t ELSE** sont d s symbol s t rminaux pour l s mots clés corr spondants. Quand l tok n **ELSE** st lu t d vi nt un tok n n avant, l cont nu d la pil (supposant l' ntré corr ct ) st valid pour un réduction par la pr mièr règl . Mais il st aussi légitim d shift r l **ELSE** parc qu c la pourrait conduire à un év ntu ll réduction par la d uxièm règl .

C tt situation, lorsqu shift r ou réduire p uv nt ètr tout d ux valid s, st app lé un conflit shift/r duc . Ocamlyacc st conçu pour résoudre c problèm n choisissant d shift r, sauf s'il y a d'autr s précisions dans l s déclarations d précéd nc sur l s opérat urs.

La raison pour laqu ll on choisit plutôt d shift r qu d réduire st xposé ci-après. Parc qu l pars r préfèr shift r l **ELSE** l résultat st d'associ r la claus " ls " au stat m nt "if" l plus imbriqué, n r ndant l s d ux ntré s suivant s équival nt s (figur 2.23) :

Il y a un conflit car la grammair t ll qu' ll st écrit st ll -mêm ambiguë : on p ut pars r un stat m nt "if" simpl imbriqué. L s conv ntions établi s sont qu l s ambiguïtès sont

```

if x then if y then win (); else los ;

if x then do; if y then win (); else los ; end;

```

FIG. 2.23 – Entrées équivalentes si shift

résolu s'en associant les clauses "ls" au statut "if" le plus imbriqué, ce qui fait Ocamlyacc choisir plutôt qu de réduire.

Remarque : il serait idéal ment plus propre d'écrire un grammaire non ambiguë, mais ce n'est pas aisé dans ce cas. Cette ambiguïté particulière a été rencontrée dans les spécifications d'*Algol 60* et s'appelle l'ambiguïté "dangling ls".

La définition précédente (figur 2.22) du `if_stmt` est la seule à blâmer pour le conflit, mais le conflit n'apparaît en fait qu'avec l'ajout d'autres règles. Voici un fichier figur 2.24 complet d'un grammaire Ocamlyacc qui provoque un conflit.

```

%token IF THEN ELSE variabl
%%
stmt:
    | xpr
    | if_stmt
    ;

if_stmt:
    IF xpr THEN stmt
    | IF xpr THEN stmt ELSE stmt
    ;

xpr:
    variabl
    ;

```

FIG. 2.24 – Grammaire à conflit

### 2.5.3 Précédence des opérateurs

D'autres conflits shift/réduire apparaissent dans les expressions arithmétiques. Shift ne s'applique pas la même solution. Les déclarations de précédence Ocamlyacc permettent de spécifier quand shift ou quand réduire.

#### Cas de précédence nécessaire

Voici figur 2.25 un grammaire ambiguë (car l'entrée `1 - 2 * 3` peut être parsée de deux façons) :

Si le parseur voit les tokens `1 - 2` devra-t-il réduire avec la règle de soustraction ? Cela dépend du token suivant.

- `)` on doit réduire, shift ne s'applique pas valide car aucune règle ne peut réduire la séquence `MINUS xpr` (c'est-à-dire `-2`) ou n'importe quel groupe commençant par ça.
- `* ou **` on doit choisir : shift ou réduire peuvent être effectués mais avec des résultats différents.

```

xpr :      xpr MINUS  xpr
        |      xpr MULTIPLY  xpr
        |      xpr LT   xpr
        |      LPAREN  xpr RPAREN
        ...
        ;

```

FIG. 2.25 – Grammaire à conflit soluble par précedence

Pour décider ce qu'Ocamlyacc doit faire, il faut regarder les résultats. Si l'opérateur suivant "op" est shifté, alors il doit être réduit d'abord pour permettre une autre possibilité de réduction pour l'-. Le résultat serait  $1 - (2 \text{ op } 3)$ . Sinon, si la soustraction est réduite avant de shift l'opérateur, le résultat est  $(1 - 2) \text{ op } 3$ . Ainsi, le choix entre shift et reduce doit dépendre du relatif précedence entre les opérateurs (ici "-" et "op").

Dans le cas où l'entrée est  $1 - 2 - 5$ , st-ce  $(1 - 2) - 5$  ou  $1 - (2 - 5)$  qui doit être connu? Pour la plupart des opérateurs, on préfère la première possibilité, qu'on appelle associativité gauche. La deuxième, association droite, est préférée pour les opérateurs d'attributions. Le choix d'association gauche ou droite dépend si le parser choisit shift ou reduce lorsque la pile contient  $1 - 2$  et que le token suivant est  $-$ .

### Spécification des précédences d'opérateurs

Ocamlyacc permet de spécifier des choix à l'aide de déclarations de précedence **%left** et **%right** (voir la sous-section 2.3.7 sur la précedence des opérateurs). Chaque déclaration contient un liste de tokens, qui sont des opérateurs dont on déclare la précedence et l'associativité. La déclaration **%left** rend les opérateurs associatifs gauche et **%right** associatifs droite. La troisième alternative **%noassoc** rend incorrecte l'association, c'est-à-dire provoque un erreur de syntaxe si le même opérateur est trouvé deux fois "d suite".

La précedence relatif entre les différents opérateurs est défini pas l'ordre dans lequel ils sont déclarés. La première déclaration **%left** ou **%right** précise les opérateurs pour lesquels la précedence est la plus faible et on peut ordonner comme ça les précedences entre les opérateurs par déclaration de précedence successive.

### Exemples de précedence

Pour la règle de grammaire de la figure 2.26 suivante :

```

xpr :      xpr MINUS  xpr
        |      xpr MULTIPLY  xpr
        |      xpr LT   xpr
        |      LPAREN  xpr RPAREN
        ...
        ;

```

FIG. 2.26 – Exemple de grammaire

on va définir les règles de précedence suivantes :

Dans un exemple plus complet, qui élargirait à plusieurs autres opérateurs, on déclarerait les opérateurs par groupe de précedence. Par exemple figure 2.28 "+" est déclaré avec "-".



```
%left LT
%left MINUS
%left MULTIPLY
```

FIG. 2.27 – Définition d'opérateurs précédence gauche

```
%left LT GT EQ NE LE GE
%left PLUS MINUS
%left MULTIPLY DIVIDE
```

FIG. 2.28 – Opérateurs précédence groupés

On note ici **NE** pour l'opérateur "not equal" et ainsi de suite pour les autres opérateurs.

### Fonctionnement de la précedence

Les premières listes de déclarations de précédence sont d'effet de niveau aux précédences aux symboles terminaux déclarés. La deuxième liste d'effet de niveau aux précédences à certains règles : chaque règle obtient sa précédence du dernier symbole terminal mentionné dans ses composants.

On peut aussi spécifier explicitement la précédence d'une règle, voir la sous-section 2.5.4 sur les précédences dépendantes du contexte.

En définitive, la résolution des conflits fonctionne en comparant les précédences des règles considérées avec le token à venir. Si la précédence du token est plus élevée, le choix est de shift. Si la précédence de la règle est plus élevée, le choix est de réduire. S'ils ont une précédence égale, le choix se fait à partir de l'associativité du niveau de précédence. Le fichier de sortie "explicit" obtenu à l'aide de l'option **-v** (voir la section Invoker Ocaml yacc 2.8) comment se résout chaque conflit.

Ni toutes les règles, ni tous les tokens n'ont de règles de précédence. Si aucun des règles ou tokens à venir n'ont de précédence, on shift par défaut.

### 2.5.4 6.4. Context-Dependent Precedence

Souvent, la précédence d'un opérateur dépend du contexte. Par exemple, un signe moins a typiquement un très fort précédence en tant qu'opérateur unaire, mais on dirait un précédence un peu moindre (moindre que la multiplication) en tant qu'opérateur binaire.

Les déclarations de précédence, **%left**, **%right** et **%nonassoc** ne peuvent être utilisées qu'une fois pour un token donné. Ainsi un token n'a qu'une seule déclaration de précédence. Afin de déclarer des précédences dépendantes du contexte, il faut utiliser un autre mécanisme, le modificateur **%prec** pour les règles.

Le modificateur **%prec** déclare la précédence d'une règle en particulier en spécifiant le symbole terminal pour lequel la précédence doit être utilisée dans cette règle. Il n'est pas nécessaire que ce symbole apparaisse ailleurs dans la règle. La syntaxe du modificateur est : **%prec terminal-symbol** et est écrit à la suite des composants de la règle. Son effet est d'affecter à la règle la précédence du symbole terminal, écrasant la précédence qui aurait pu être déduit classiquement. La précédence de la règle altérée est utilisée pour résoudre les conflits (voir la sous-sous-section 2.3.7 sur les opérateurs de précédence).

Voici figure 2.29 comment **%prec** résoudre le problème du moins unaire. On déclare d'abord

un précédenc pour le symbol terminal fictif UMINUS. Il n'y a pas de token de ce type mais le symbol sert pour marquer sa précédence :

```
%...
%left PLUS MINUS
%left MULTIPLY
%left UMINUS
....
%  xp :      ...
%          |   xp MINUS  xp
%          |   ...
%          | MINUS  xp %prec UMINUS
```

FIG. 2.29 – Définition des précédences de context

### 2.5.5 Etats du parser

La fonction *yyparse* est implémentée en utilisant un machine à états finis. Les valeurs "pushed" sur la pile du parser ne sont pas simplement des codes de type de token, ils sont présents dans la séquence symbol terminaux et non-terminaux sur ou près du sommet de la pile. L'état courant accumule toutes les informations concernant les entrées précédentes qui sont significatives pour décider de ce qui doit être fait ensuite.

Chaque fois qu'un token à valoir est lu, l'état courant du parser, avec le type du token à valoir sont associés dans un tableau. Le tableau des entrées peut indiquer "shift le token à valoir". Dans ce cas, il peut aussi spécifier le nouvel état du parser, qui est "pushed" au sommet de la pile du parser. Le tableau des entrées peut aussi indiquer "réduire" en utilisant la règle *r*. Ceci signifie qu'un certain nombre de tokens ou groupes sont dépilés et remplacés par un groupe. En d'autres termes, un nombre d'états est dépilé de la pile et un nouvel état est "pushed".

Il existe une autre alternative : le tableau peut indiquer qu'un token à valoir est erroné dans l'état courant. Cela provoque un erreur de parsing (voir la section sur la récupération d'erreur : 2.6).

### 2.5.6 Les conflits Reduce/Reduce

Un conflit reduce/reduce arrive s'il y a deux ou plusieurs règles qui s'appliquent à la même séquence d'entrée. Généralement, cela indique une sérieuse erreur dans la grammaire.

Par exemple figure 2.30, voici un exemple pour définir une séquence de zéro ou plus groupes de mots.

L'erreur est une ambiguïté : il y a plus d'un moyen de parser un simple mot dans une séquence. Il pourrait être réduit comme **maybe word** puis comme **s qu ne** à l'aide de la deuxième règle. Mais également, "rien du tout" pourrait être réduit en **s qu ne** à l'aide de la première règle et être combiné avec le mot en utilisant la 3ème règle de **s qu ne**.

Il y a également plusieurs moyens de réduire "rien du tout" en un **s qu ne**. Cela peut être fait directement à l'aide de la première règle ou indirectement à l'aide de la règle **maybe word** puis la 2ème règle.

On pourrait penser que cette distinction ne fait pas de différence, puisque ça ne change pas le fait qu'une entrée soit valide ou non. Mais cela affecte qu'il s'agit d'actions pouvant être appliquées. Un ordre de parsing applique les actions de la seconde règle, l'ordre d'un autre parsing applique

```

s qu nc : /* mpty */ { printf " mpty_s qu nc \n" }
          | mayb word  {}
          | s qu nc word { printf "add d_word_%s\n" $2 }
          ;

mayb word: /* mpty */ { printf " mpty_mayb word\n" }
          | word      { printf "singl _word_%s\n" $1 }
          ;

```

FIG. 2.30 – Err ur d grammair provoquant un conflit r duc /r duc

l s actions d la pr mièr puis d la troisièm règl . Dans c t x mpl , la sorti du programm chang .

Ocamlyacc résout l s conflits r duc /r duc n choisissant d'utilis r l s règl s dans l'ordr d'apparition dans la grammair , mais c la p ut êtr risqué d s r pos r s ul m nt là d ssus. Chaque conflit r duc /r duc doit êtr étudié t éliminé. Voici figur 2.31 un moy n corr ct d définir `s qu nc` :

```

s qu nc : /* mpty */ { printf " mpty_s qu nc \n" }
          | s qu nc word { printf "add d_word_%s\n" $2 }
          ;

```

FIG. 2.31 – Résolution du conflit r duc /r duc

Voici figur 2.32 un autr x mpl d' rr ur commun qui conduit à un conflit r duc /r duc :

```

s qu nc : /* mpty */
          | s qu nc words
          | s qu nc r dir cts
          ;

words:    /* mpty */
          | words word
          ;

r dir cts:/* mpty */
          | r dir cts r dir ct
          ;

```

FIG. 2.32 – Ex mpl 2 d'un conflit r duc /r duc

L'int ntion ici st d définir un séqu nc qui p ut cont nir soit d s group s `word` soit d s group s `r dir ct` . L s définitions individu ll s d `s qu nc` , `words t r dir cts` n'ont pas d' rr urs, mais l'ass mblag d tout s apport un ambiguïté subtil : un ntré vid p ut êtr parsé par un infinité d moy n!

Ici, "ri n du tout" p ut st un `words` , ou 2 `words` à la suit ou 3 ou n'import qu l nombr . "ri n du tout" p ut tout aussi êtr `r dir cts` un , d ux ou n'import qu l nombr d fois. Ou `words` suivi d trois `r dir cts` t d'autr s `words`...

Voici donc figur 2.33 un pr mi r moy n d corrig r c s règl s, n n'ayant qu'un niv au d s qu nc :

```
s qu nc : /* mpty */
        | s qu nc word
        | s qu nc r dir ct
        ;
```

FIG. 2.33 – Résolution 1 d l' x mpl 2 d'un conflit r duc /r duc

Figur 2.34, on corrig n évitant soit **word** soit **r dir cts** d'êtr vid :

```
s qu nc : /* mpty */
        | s qu nc words
        | s qu nc r dir cts
        ;

words :   word
        | words word
        ;

r dir cts : r dir ct
          | r dir cts r dir ct
          ;
```

FIG. 2.34 – Résolution 2 d l' x mpl 2 d'un conflit r duc /r duc

## 2.5.7 Mystérieux Reduce/Reduce Conflits

Il arriv qu l s conflits r duc /r duc puiss nt s produire sans dét ction garanti . Figur 2.35 montr un x mpl :

Il pourrait s mbl r qu c tt grammair puiss êtr parsé av c s ul m nt un simpl tok n à v nir : lorsqu **param\$\_\$sp c** st lu, un **ID** st un **nam** si "," ou ":" ou un **type** si un autr **ID** st à la suit . En d'autr s t rm s, c tt grammair st LR(1).

C p ndant, Ocamlyacc , comm la plupart d s générat urs d pars urs, n p ut pas support r tout s l s grammair s LR(1). Dans la grammair d c t x mpl 2.35, d ux cont xt s : après un **ID** au début d'un **param\$\_\$sp c** t **ID** au début d'un **r turn\$\_\$sp c**, sont suffisamm nt similar s pour qu Ocamlyacc p ns nt qu'ils soi nt id ntiqu s. Ils apparaiss nt id ntiqu s à Ocamlyacc car l mêm ns mbl d règl s pourrait êtr activé (c ll pour réduire n un **nam** t c ll pour réduire n un **type** . Ocamlyacc n p ut pas dét rmin r à c mom nt d la génération du pars ur, qu l s règl s néc ssit ront un tok n à v nir différ nt dans l s d ux cont xt s, il construit donc un s ul état d pars r pour ux. Or, combin r l s d ux cont xt s provoqu un conflit ultéri ur m nt. En t rm d pars r, c cas illustr qu'il n s'agit pas d'un grammair LALR(1).

Il st plutôt compl x d corrig r Ocamlyacc pour lui p rm ttr d pars r d s grammair s LR(1) qui n sont pas LALR(1). L s générat urs d pars ur qui l font ont t ndanc à génér r d s pars urs très très gros. En pratiqu , on préfér laiss r Ocamlyacc t l qu'il st.

```

%token ID COMMA COLON

%%
d f:    param_sp c r turn_sp c COMMA
        ;
param_sp c:
        type
        | nam _list COLON type
        ;
r turn_sp c:
        type
        | nam COLON type
        ;
type:    ID
        ;
nam:     ID
        ;
nam _list:
        nam
        | nam COMMA nam _list
        ;

```

FIG. 2.35 – Conflit r duc /r duc non dét ctabl

Lorsqu' l problèm LR(1) s'présnt, on p'ut souv'nt l' résoudre n' id ntifiant l' s' d' ux états d' pars ur qui prèt nt à confusions. t' ajout' r qu' lq' chos qui l' s' f'ront s' distingu' r. Dans l' x' mpl' précéd' nt, ajout' r un' règl' à `r turn$_$sp c` comm' dans l' figur' 2.36 p' rm' t' d' résoudre l' problèm :

```

%token BOGUS
...
%%
...
r turn_sp c:
        type
        | nam COMMA type
        /* This rul' is n' v' r' us' d. */
        | ID BOGUS
        ;

```

FIG. 2.36 – Résolution d'un conflit r duc /r duc non dét ctabl

C'la corr'g' l' problèm n' introduisant la possibilité d'activ' r un' règl' supplém' ntair' dans l' cont'xt' après `ID` au début du `r turn$_$sp c`. C' tt' règl' n' st pas activ' dans l' cont'xt' corr' spondant d' `param_sp c`, t' d' c' tt' façon l' s' d' ux cont'xt' s' r' çoiv' nt d' s' états d' pars ur différ' nts. Du mom' nt qu' l' tok' n' `BOGUS` n' st pas généré par Ocaml' x, la règl' ajoutée n' modifi' pas la façon d' pars' r l' ntré.

Dans c' t' x' mpl' particul' r, il y a un autr' moy' n' d' résoudre l' problèm : réécri' r la

règl d `r turn$_$sp c` afin d'utilis r `ID` dir et m nt, autr m nt qu via `nam`. C la p rm t qu l s d ux cont xt s ai nt d s ns mbl s d règl s activ s différ nts, puisqu la règl pour `r turn$_$sp c` activ la règl modifié pour `r turn$_$sp c` plutôt qu c ll pour `nam`. C ci st prés nté figur 2.37 :

```
param_sp c :
    type
    |   nam _list COMMA type
    ;
r turn_sp c :
    type
    |   ID COMMA type
    ;
```

FIG. 2.37 – Résolution d'un conflit r duc /r duc non dét ctabl - 2èm possibilité

## 2.6 Error récupération

En général, il n faut pas qu'un programm s'int rrompr sur un rrr ur d parsing. Par x mpl , un compilat ur d vrait résist r aux rrr urs suffisamm nt pour pars r l r st du fichi r d' ntré à la r ch rch d' rrr urs; un calculatric d vrait acc pt r un autr xpr ssion.

Lorsqu l'on fait du parsing av c un simpl command int ractiv où chaque ntré st constitué d'un s ul lign , il p ut êtr suffisant d'avoir l'app lant qui attrap l' xc ption t qui ignor l r st d la lign d' ntré lorsqu'un rrr ur survi nt ( t on app ll la fonction d parsing à nouv au pour continu r sur un autr lign ). Mais c la n' st pas adapté pour un compilat ur, parc qu'il "oubl i " tout l cont xt syntaxiqu qui a am né à l' rrr ur. Un rrr ur d syntax dans un fonction du compilat ur n d vrait pas provoqu r qu la "lign " suivant soit traité d la mêm manière qu si c'était l début du fichi r sourc .

On p ut donc définir comm nt résist r à un rrr ur d syntax n écrivant d s règl s qui r connaiss nt un tok n spécial d' rrr ur. C' st un symbol t rminal qui st rés rvé à la récupéra- tion d' rrr ur. L pars ur Ocamlyacc génér un tok n d' rrr ur lorsqu s produit un rrr ur d syntax ; si on définit un règl qui r connaît c tok n dans l cont xt , l parsing p ut alors continu r.

Par x mpl , figur 2.38 :

```
stmnts: /* mpty string */ {}
        | stmnts NEWLINE {}
        | stmnts xp NEWLINE {}
        | stmnts rror NEWLINE {}
```

FIG. 2.38 – Rattrapag d' rrr ur

La quatrièm règl dans c t x mpl dit qu' `rror` suivi d'un nouv ll lign `NEWLINE` s ra un combinaison valid à n'import qu l `stmnts`.

Qu s pass -t-il si un rrr ur d syntax s produit au mili u d'un `xp` ? La règl d récupération d' rrr ur, int rprété strict m nt, va appliqu r la séqu nc précis d s `stmnts`, un `rror` t un `NEWLINE`. Si un rrr ur s produit au mili u d'un `xp`, il y aura probabl m nt

des tokens supplémentaires et des sous-expressions dans la pile après le dernier `stmtnts`, il y aura des tokens à lire avant la prochaine `NEWLINE`. La règle n'est donc pas applicable normalement.

Mais Ocamlyacc peut forcer la situation de façon à correspondre à la règle, en éliminant un parti du contexte sémantique et un parti de l'entrée. Premièrement, Ocamlyacc éliminera les états et les objets de la pile jusqu'à revenir à l'état dans lequel le token `rror` est accepté. (Cela signifie que les sous-expressions déjà parsées sont éliminées, en retournant au dernier `stmtnts` complet.) À ce niveau, le token `rror` peut être shifté. Alors, si l'ancien token à venir n'est pas acceptable pour être shifté à l'état suivant, le parser lit les tokens et les éliminera jusqu'à trouver un token qui est acceptable. Dans cet exemple 2.38, Ocamlyacc lit et éliminera les entrées jusqu'à la prochaine `NEWLINE` afin que la quatrième règle puisse s'appliquer.

Le choix des règles d'erreur dans un grammair est un choix de stratégie pour la récupération d'erreur. Une stratégie simple et utile est simplement de skipper la ligne d'entrée courante ou d'état courant si une erreur est détectée, comme le montre l'exemple de la figure 2.39 :

```
stmtnt: rror SEMICOLON {} /* en cas d'erreur, skipper jusqu'au ;
*/
```

FIG. 2.39 – Stratégie de récupération d'erreur

Il est aussi utile de récupérer jusqu'au délimitur fermant correspondant à un délimitur ouvrant qui aurait déjà été parsé. Sinon, le délimitur fermant apparaîtra probablement isolé et générera un autre message d'erreur (ronné), voir figure 2.40 :

```
primary: LPAREN xpr RPAREN {}
        | LPAREN rror RPAREN {}
        ...
        ;
```

FIG. 2.40 – Erreur causée entre deux délimiturs

Les stratégies de récupération d'erreur sont nécessaires ment de statutatives. Lorsqu'un testatif échoue, une erreur de syntaxe mène souvent à un autre ... Dans l'exemple précédent, la règle de récupération d'erreur suppose qu'une erreur est due à un mauvais entrée à l'intérieur d'un `stmtnts`. Supposons qu'un `;$` soit inséré au milieu d'un `stmtnts` valide. Après avoir utilisé les règles de récupération d'erreur de la première erreur, un autre erreur de syntaxe est trouvée ensuite : puisque le texte qui suit le `;$` ronné est également un `stmtnts` invalide.

Afin d'éviter les messages d'erreur abondants, le parser ne générera pas de message d'erreur pour une erreur de syntaxe qui se produit juste après la première ; ou même après trois tokens d'entrée consécutifs correctement shiftés, les messages d'erreur sont à nouveau générés dans ce cas échéant.

## 2.7 Debugger son parser

Pour déboguer le parser généré par Ocamlyacc :

- En utilisant la commande `ocamlyacc -v filename.mly`, on génère les informations de parsing dans un fichier appelé `filename.output`. Ces informations consistent en un tableau de parsing et des indications concernant les conflits.

- Mettre l'option `p` à la variable d'environnement `OCAMLRUNPARAM` (`export OCAMLRUNPARAM='p'` dans un shell `bash`).

Le parser affichera des messages à propos des actions et des qu'il shift un token, réduire une règle ...

On peut obtenir les numéros de la règle ou les numéros d'états mentionnés dans les messages dans le fichier `.output`.

## 2.8 Executer Ocamllyacc

La manière habituelle d'invoquer `ocamllyacc` : `ocamllyacc fichier .mly`. Où `fichier .mly` est le nom du fichier de grammaire. Le nom du fichier de parsing est obtenu en remplaçant le `.mly` par `.ml`. `Ocamllyacc` génère un `.ml` à partir du `.mly`

### 2.8.1 Ocamllyacc Options

Voici une liste d'options qui peuvent être utilisées avec `Ocamllyacc` :

- `-v` Par défaut, cette option génère un fichier `fichier .output`. Il contient des informations de parsing et la description des tables de parsing et un rapport des ambiguïtés de la grammaire.
- `-bfichier` Permet de changer le nom du fichier de sortie par `fichier .ml` `fichier .mli` et `fichier .output`



## Chapitre 3

### Astuces

- Quand on fait du parsing " n lign " plutôt qu'à l'aid d'un fichi r, CTRL D pour fin d fichi r
- Utilis r < t > pour ntré t sorti du programm compl t : `./prog < ntr > sorti`
- Pour vérifi r la grammair , lir .output
- Fair un Mak fil quand on travaill sur l proc ssus compl t (analys l xical , syntaxiqu t actions sémantiqu s...). Voir l' x mpl figur 3.1.

```
.PRECIOUS: %.mli

all: prog

prog: l x r.cmo pars r.cmo prog.cmo
      ocamlc -o $@ $^

%.cmo: %.ml
      ocamlc -c $<

%.cmi: %.mli
      ocamlc -c $<

%.ml: %.mll
      ocaml1 x $<

%.ml %.mli: %.mly
      ocaml yacc -v $<

l x r.cmo: pars r.cmi

cl an :
      rm -f *.cmi *.cmo
```

FIG. 3.1 – Ex mpl d Mak fil