

École Normale Supérieure

Langages de programmation et compilation

Jean-Christophe Filliâtre

Cours 7 / 16 novembre 2015

- notion de grammaire
 - dérivation
 - grammaire ambiguë
- analyse descendante
 - une table nous indique quelle expansion choisir
 - calcul de NULL, FIRST, FOLLOW
 - calcul de point fixe
 - grammaire LL(1)

$$\begin{array}{l} E \rightarrow E + E \\ \quad | E * E \\ \quad | (E) \\ \quad | \text{int} \end{array}$$

analyse ascendante

l'idée est toujours de lire l'entrée de gauche à droite, mais on cherche maintenant à reconnaître des membres droits de productions pour construire l'arbre de dérivation de bas en haut (*bottom-up parsing*)

l'analyse manipule toujours une pile qui est un mot de $(T \cup N)^*$

à chaque instant, deux actions sont possibles

- opération de **lecture** (*shift* en anglais) : on lit un terminal de l'entrée et on l'empile
- opération de **réduction** (*reduce* en anglais) : on reconnaît en sommet de pile le membre droit β d'une production $X \rightarrow \beta$, et on remplace β par X en sommet de pile

dans l'état initial, la pile est vide

lorsqu'il n'y a plus d'action possible, l'entrée est reconnue si elle a été entièrement lue et si la pile est réduite à S

	pile	entrée	action
	ϵ	int+int*int	lecture
	int	+int*int	réduction $F \rightarrow \text{int}$
	F	+int*int	réduction $T \rightarrow F$
	T	+int*int	réduction $E \rightarrow T$
$E \rightarrow E + T$	E	+int*int	lecture
T	$E+$	int*int	lecture
$T \rightarrow T * F$	$E+\text{int}$	*int	réduction $F \rightarrow \text{int}$
F	$E+F$	*int	réduction $T \rightarrow F$
$F \rightarrow (E)$	$E+T$	*int	lecture
int	$E+T*$	int	lecture
	$E+T*\text{int}$		réduction $F \rightarrow \text{int}$
	$E+T*F$		réduction $T \rightarrow T*F$
	$E+T$		réduction $E \rightarrow E+T$
	E		succès

comment prendre la décision lecture / réduction ?

en se servant d'un automate fini et en examinant les k premiers caractères de l'entrée ; c'est l'analyse LR(k)

(LR signifie « **L**eft to right scanning, **R**ightmost derivation »)

en pratique $k = 1$

i.e. on examine uniquement le premier caractère de l'entrée

la pile est de la forme

$$s_0 x_1 s_1 \dots x_n s_n$$

où s_i est un état de l'automate et $x_i \in T \cup N$ comme auparavant

soit a le premier caractère de l'entrée ; on regarde la transition de l'automate pour l'état s_n et l'entrée a

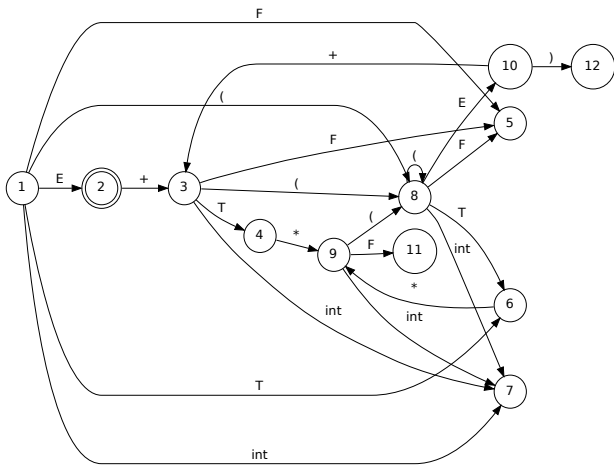
- si c'est un succès ou un échec, on s'arrête
- si c'est une lecture, alors on empile a et l'état résultat de la transition
- si c'est une réduction $X \rightarrow \alpha$, avec α de longueur p , alors on doit trouver α en sommet de pile

$$s_0 x_1 s_1 \dots x_{n-p} s_{n-p} \mid \alpha_1 s_{n-p+1} \dots \alpha_p s_n$$

on dépile alors α et on empile Xs , où s est l'état résultat de la transition $s_{n-p} \xrightarrow{X} s$, *i.e.*

$$s_0 x_1 s_1 \dots x_{n-p} s_{n-p} X s$$

dans l'exemple plus haut, on s'est servi de cet automate

$$\begin{array}{l}
 E \rightarrow E + T \\
 \quad | \quad T \\
 T \rightarrow T * F \\
 \quad | \quad F \\
 F \rightarrow (E) \\
 \quad | \quad \text{int}
 \end{array}$$


construction de l'automate

fixons pour l'instant $k = 0$

on commence par construire un automate **asynchrone**

c'est-à-dire contenant des transitions spontanées
appelées ϵ -**transitions** et notées $s_1 \xrightarrow{\epsilon} s_2$

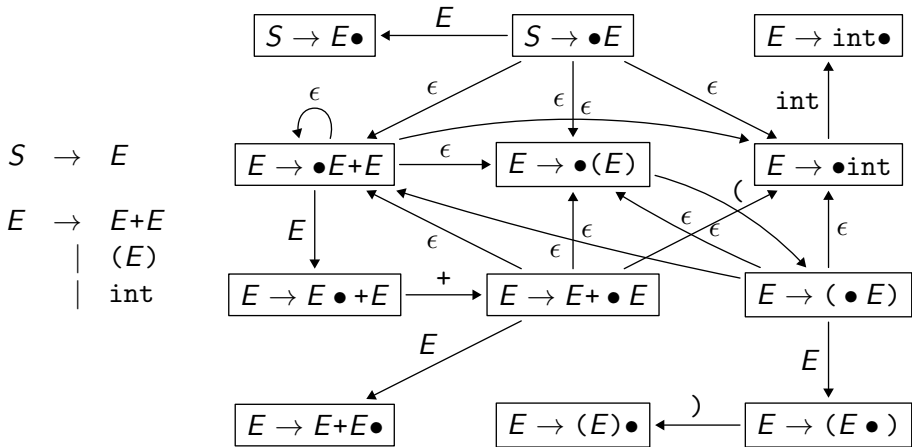
les **états** sont des *items* de la forme

$$[X \rightarrow \alpha \bullet \beta]$$

où $X \rightarrow \alpha\beta$ est une production de la grammaire ; l'intuition est « je cherche à reconnaître X , j'ai déjà lu α et je dois encore lire β »

les **transitions** sont étiquetées par $T \cup N$ et sont les suivantes

$$\begin{aligned} [Y \rightarrow \alpha \bullet a\beta] &\xrightarrow{a} [Y \rightarrow \alpha a \bullet \beta] \\ [Y \rightarrow \alpha \bullet X\beta] &\xrightarrow{X} [Y \rightarrow \alpha X \bullet \beta] \\ [Y \rightarrow \alpha \bullet X\beta] &\xrightarrow{\epsilon} [X \rightarrow \bullet \gamma] \quad \text{pour toute production } X \rightarrow \gamma \end{aligned}$$



déterminisons l'automate LR(0)

pour cela, on regroupe les états reliés par des ϵ -transitions

les états de l'automate déterministe sont donc des ensembles d'*items*, tel que

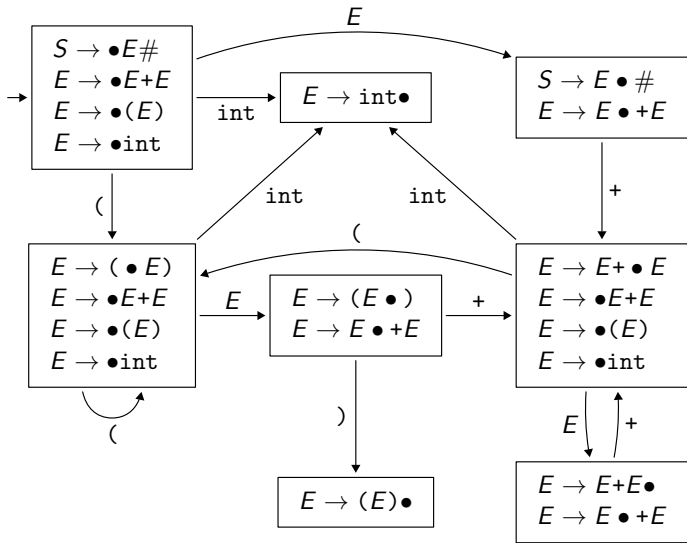
$E \rightarrow E+ \bullet E$ $E \rightarrow \bullet E+E$ $E \rightarrow \bullet (E)$ $E \rightarrow \bullet \text{int}$

chaque état s est **saturé** par la propriété

si $Y \rightarrow \alpha \bullet X \beta \in s$
 et si $X \rightarrow \gamma$ est une production
 alors $X \rightarrow \bullet \gamma \in s$

l'état initial est celui contenant $S \rightarrow \bullet E$

$S \rightarrow E$
 $E \rightarrow E+E$
 $E \rightarrow (E)$
 $E \rightarrow int$



en pratique, on ne travaille pas directement sur l'automate mais sur deux tables

- une table d'**actions** ayant pour lignes les états et pour colonnes les terminaux; la case $\text{action}(s, a)$ indique
 - shift s' pour une lecture et un nouvel état s'
 - reduce $X \rightarrow \alpha$ pour une réduction
 - un succès
 - un échec
- une table de **déplacements** ayant pour lignes les états et pour colonnes les non terminaux; la case $\text{goto}(s, X)$ indique l'état résultat d'une réduction de X

on construit ainsi la table action

- $\text{action}(s, \#) = \text{succès}$ si $[S \rightarrow E \bullet \#] \in s$
- $\text{action}(s, a) = \text{shift } s'$ si on a une transition $s \xrightarrow{a} s'$
- $\text{action}(s, a) = \text{reduce } X \rightarrow \beta$ si $[X \rightarrow \beta \bullet] \in s$, pour tout a
- échec dans tous les autres cas

on construit ainsi la table goto

- $\text{goto}(s, X) = s'$ si et seulement si on a une transition $s \xrightarrow{X} s'$

sur notre exemple, la table est la suivante :

	<i>action</i>					<i>goto</i>
état	()	+	int	#	<i>E</i>
1	shift 4			shift 2		3
2	reduce $E \rightarrow \text{int}$					
3			shift 6		succès	
4	shift 4			shift 2		5
5		shift 7	shift 6			
6	shift 4			shift 2		8
7	reduce $E \rightarrow (E)$					
8			shift 6			
	reduce $E \rightarrow E+E$					

la table LR(0) peut contenir deux sortes de conflits

- un conflit **lecture/réduction** (*shift/reduce*), si dans un état s on peut effectuer une lecture mais aussi une réduction
- un conflit **réduction/réduction** (*reduce/reduce*), si dans un état s deux réductions différentes sont possibles

Définition (classe LR(0))

Une grammaire est dite LR(0) si la table ainsi construite ne contient pas de conflit.

on a un conflit lecture/réduction dans l'état 8

$$\begin{array}{l} E \rightarrow E+E\bullet \\ E \rightarrow E\bullet +E \end{array}$$

il illustre précisément l'ambiguïté de la grammaire sur un mot tel que `int+int+int`

on peut résoudre le conflit de deux façons

- si on favorise la **lecture**, on traduira une associativité à droite
- si on favorise la **réduction**, on traduira une associativité à gauche

privilégions la réduction et illustrons sur un exemple

	()	+	int	#	E
1	s4			s2		3
2	reduce $E \rightarrow \text{int}$					
3			s6		ok	
4	s4			s2		5
5		s7	s6			
6	s4			s2		8
7	reduce $E \rightarrow (E)$					
8	reduce $E \rightarrow E+E$					

pile	entrée	action
1	int+int+int	s2
1 int 2	+int+int	$E \rightarrow \text{int}, g3$
1 E 3	+int+int	s6
1 E 3 + 6	int+int	s2
1 E 3 + 6 int 2	+int	$E \rightarrow \text{int}, g8$
1 E 3 + 6 E 8	+int	$E \rightarrow E+E, g3$
1 E 3	+int	s6
1 E 3 + 6	int	s2
1 E 3 + 6 int 2	#	$E \rightarrow \text{int}, g8$
1 E 3 + 6 E 8	#	$E \rightarrow E+E, g3$
1 E 3	#	succès

la construction LR(0) engendre très facilement des conflits
on va donc chercher à limiter les réductions

une idée très simple consiste à poser $\text{action}(s, a) = \text{reduce } X \rightarrow \beta$ si et seulement si

$$[X \rightarrow \beta \bullet] \in s \quad \text{et} \quad a \in \text{FOLLOW}(X)$$

Définition (classe SLR(1))

Une grammaire est dite SLR(1) si la table ainsi construite ne contient pas de conflit.

(SLR signifie *Simple LR*)

la grammaire

$$\begin{aligned} S &\rightarrow E\# \\ E &\rightarrow E + T \\ &\quad | T \\ T &\rightarrow T * F \\ &\quad | F \\ F &\rightarrow (E) \\ &\quad | \text{int} \end{aligned}$$

est SLR(1)

exercice : le vérifieur (l'automate contient 12 états)

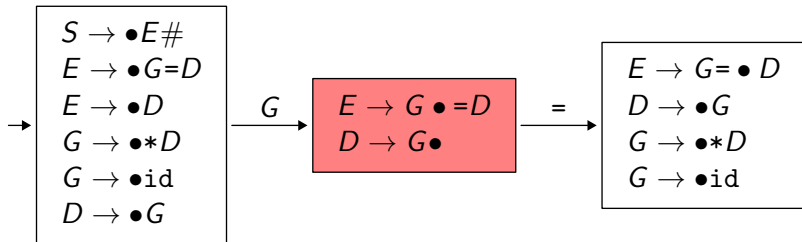
limites de l'analyse SLR(1)

en pratique, la classe SLR(1) reste trop restrictive

exemple :

$$\begin{aligned} S &\rightarrow E\# \\ E &\rightarrow G = D \\ &\quad | D \\ G &\rightarrow *D \\ &\quad | \text{id} \\ D &\rightarrow G \end{aligned}$$

	=	
1
2	shift 3 reduce $D \rightarrow G$...
3	⋮	⋮



on introduit une classe de grammaires encore plus large, **LR(1)**, au prix de tables encore plus grandes

dans l'analyse LR(1), les *items* ont maintenant la forme

$$[X \rightarrow \alpha \bullet \beta, a]$$

dont la signification est : « je cherche à reconnaître X , j'ai déjà lu α et je dois encore lire β puis vérifier que le caractère suivant est a »

les transitions de l'automate LR(1) non déterministe sont

$$\begin{aligned}
 [Y \rightarrow \alpha \bullet a\beta, b] &\xrightarrow{a} [Y \rightarrow \alpha a \bullet \beta, b] \\
 [Y \rightarrow \alpha \bullet X\beta, b] &\xrightarrow{X} [Y \rightarrow \alpha X \bullet \beta, b] \\
 [Y \rightarrow \alpha \bullet X\beta, b] &\xrightarrow{\epsilon} [X \rightarrow \bullet \gamma, c] \quad \text{pour tout } c \in \text{FIRST}(\beta b)
 \end{aligned}$$

l'état initial est celui qui contient $[S \rightarrow \bullet \alpha, \#]$

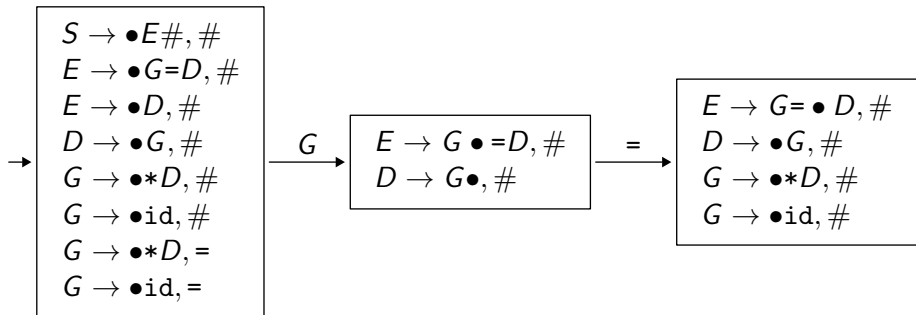
comme précédemment, on peut déterminer l'automate et construire la table correspondante ; on introduit une action de réduction pour (s, a) seulement lorsque s contient un item de la forme $[X \rightarrow \alpha \bullet, a]$

Définition (classe LR(1))

Une grammaire est dite LR(1) si la table ainsi construite ne contient pas de conflit.

$S \rightarrow E\#$
 $E \rightarrow G = D$
 $\quad | D$
 $G \rightarrow * D$
 $\quad | id$
 $D \rightarrow G$

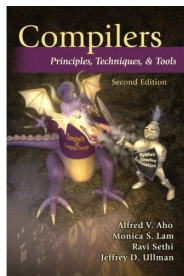
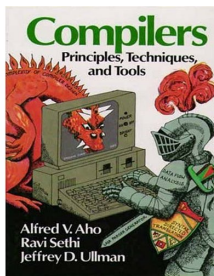
	#	=	
1
2	reduce $D \rightarrow G$	shift 3	...
3	:	:	..



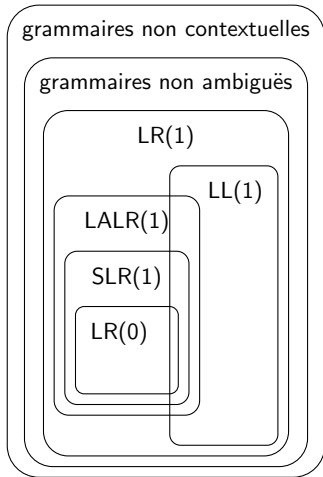
la construction LR(1) pouvant être coûteuse, il existe des approximations

la classe LALR(1) (*lookahead LR*) est une telle approximation, utilisée notamment dans les outils de la famille yacc

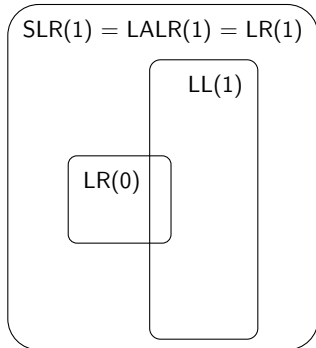
plus d'info : voir par exemple *Compilateurs : principes techniques et outils* (dit « le dragon ») de A. Aho, R. Sethi, J. Ullman, section 4.7



grammaires



langages



l'analyse ascendante est puissante mais le calcul des tables est complexe

le travail est automatisé par de nombreux outils

c'est la grande famille de `yacc`, `bison`, `ocamlyacc`, `cups`, `menhir`, ...
(YACC signifie *Yet Another Compiler Compiler*)

l'outil Menhir

Menhir est un outil qui transforme une grammaire en un analyseur OCaml ; Menhir est basé sur une analyse LR(1)

chaque production de la grammaire est accompagnée d'une **action sémantique** *i.e.* du code OCaml construisant une valeur sémantique (typiquement un arbre de syntaxe abstraite)

Menhir s'utilise conjointement avec un analyseur lexical (typiquement `ocamllex`)

un fichier Menhir porte le suffixe `.mly` et a la structure suivante

```
%{  
  ... code OCaml arbitraire ...  
%}  
...déclaration des tokens...  
...déclaration des précédences et associativités...  
...déclaration des points d'entrée...  
%%  
non-terminal-1:  
| production { action }  
| production { action }  
;  
  
non-terminal-2:  
| production { action }  
...  
%%  
  ... code OCaml arbitraire ...
```

```
%token PLUS LPAR RPAR EOF
```

```
%token <int> INT
```

```
%start <int> phrase
```

```
%%
```

```
phrase:
```

```
    e = expression; EOF  { e }
```

```
;
```

```
expression:
```

```
| e1 = expression; PLUS; e2 = expression { e1 + e2 }
```

```
| LPAR; e = expression; RPAR           { e }
```

```
| i = INT                               { i }
```

```
;
```

on compile le fichier `arith.mly` de la manière suivante

```
% menhir -v arith.mly
```

on obtient du code OCaml pur dans `arith.ml(i)`, qui contient notamment

- la déclaration d'un type `token`

```
type token = RPAR | PLUS | LPAR | INT of int | EOF
```

- pour chaque non terminal déclaré avec `%start`, une fonction du type

```
val phrase: (Lexing.lexbuf -> token) -> Lexing.lexbuf -> int
```

comme on le voit, cette fonction prend en argument un analyseur lexical, du type de celui produit par `ocamllex` (cf cours 5)

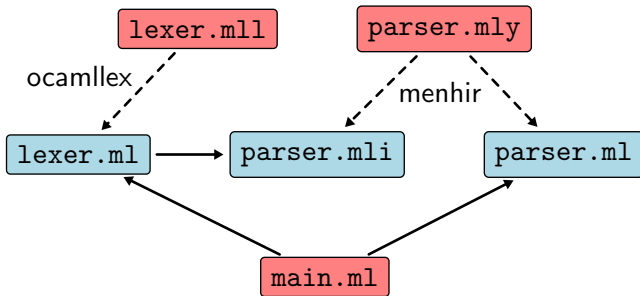
quand on combine ocamllex et menhir

- lexer.mll fait référence aux lexèmes définis dans parser.mly

```
{  
  open Parser  
}  
...
```

- l'analyseur lexical et l'analyseur syntaxique sont combinés ainsi :

```
let c = open_in file in  
let lb = Lexing.from_file c in  
let e = Parser.phrase Lexer.token lb in  
...
```



= source utilisateur

= construit automatiquement

→ = dépendance

lorsque la grammaire n'est pas LR(1), Menhir présente les **conflits** à l'utilisateur

- le fichier `.automaton` contient une description de l'automate LR(1); les conflits y sont mentionnés
- le fichier `.conflicts` contient, le cas échéant, une explication de chaque conflit, sous la forme d'une séquence de lexèmes conduisant à deux arbres de dérivation

sur la grammaire ci-dessus, Menhir signale un conflit

```
% menhir -v arith.mly
Warning: one state has shift/reduce conflicts.
Warning: one shift/reduce conflict was arbitrarily resolved.
```

le fichier arith.automaton contient notamment

```
State 6:
expression -> expression . PLUS expression [ RPAR PLUS EOF ]
expression -> expression PLUS expression . [ RPAR PLUS EOF ]
-- On PLUS shift to state 5
-- On RPAR reduce production expression -> expression PLUS expression
-- On PLUS reduce production expression -> expression PLUS expression
-- On EOF reduce production expression -> expression PLUS expression
** Conflict on PLUS
```


le fichier `arith.conflicts` contient une explication limpide

```
** Conflict (shift/reduce) in state 6.  
** Token involved: PLUS  
** This state is reached from phrase after reading:  
  
expression PLUS expression  
  
** In state 6, looking ahead at PLUS, shifting is permitted  
** because of the following sub-derivation:  
  
expression PLUS expression  
           expression . PLUS expression  
  
** In state 6, looking ahead at PLUS, reducing production  
** expression -> expression PLUS expression  
** is permitted because of the following sub-derivation:  
  
expression PLUS expression // lookahead token appears
```

une manière de résoudre les conflits est d'indiquer à Menhir comment choisir entre lecture et réduction

pour cela, on peut donner des **priorités** aux lexèmes et aux productions, et des règles d'**associativité**

par défaut, la priorité d'une production est celle de son lexème le plus à droite (mais elle peut être spécifiée explicitement)

si la priorité de la production est supérieure à celle du lexème à lire,
alors la réduction est favorisée

inversement, si la priorité du lexème est supérieure,
alors la lecture est favorisée

en cas d'égalité, l'associativité est consultée : un lexème associatif à
gauche favorise la réduction et un lexème associatif à droite la lecture

dans notre exemple, il suffit d'indiquer par exemple que PLUS est associatif à gauche

```
%token PLUS LPAR RPAR EOF
%token <int> INT
%left PLUS
%start <int> phrase
%%
phrase:
    e = expression; EOF { e }
;
expression:
| e1 = expression; PLUS; e2 = expression { e1 + e2 }
| LPAR; e = expression; RPAR { e }
| i = INT { i }
;
```

pour associer des priorités aux lexèmes, on utilise la convention suivante :

- l'ordre de déclaration des associativités fixe les priorités (les premiers lexèmes ont les priorités les plus faibles)
- plusieurs lexèmes peuvent apparaître sur la même ligne, ayant ainsi la même priorité

exemple :

```
%left PLUS MINUS  
%left TIMES DIV
```

la grammaire suivante contient un conflit

```
expression:  
| IF e1 = expression; THEN; e2 = expression  
  { ... }  
| IF e1 = expression; THEN; e2 = expression;  
  ELSE; e3 = expression  
  { ... }  
| i = INT  
  { ... }  
| ...
```

il correspond à la situation

```
IF a THEN IF b THEN c ELSE d
```

pour associer le ELSE au THEN le plus proche, il faut privilégier la lecture

```
%nonassoc THEN  
%nonassoc ELSE
```

(connu en anglais sous le nom de *dangling else*)

Menhir offre de nombreux avantages par rapport aux outils traditionnels tels que `ocaml yacc` :

- non-terminaux paramétrés par des (non-)terminaux
 - en particulier, facilités pour écrire des expressions régulières ($E?$, E^* , E^+), des listes avec séparateur
- explication des conflits
- mode interactif
- analyse LR(1) plutôt que LALR(1)

lire le manuel de Menhir! (accessible depuis la page du cours)

pour que les phases suivantes de l'analyse (typiquement le typage) puissent **localiser** les messages d'erreur, il convient de conserver une information de localisation dans l'arbre de syntaxe abstraite

Menhir fournit cette information dans `$startpos` et `$endpos`, deux valeurs du type `Lexing.position`; cette information lui a été transmise par l'analyseur lexical

attention : `ocamllex` ne maintient automatiquement que la position absolue dans le fichier; pour avoir les numéros de ligne et de colonnes à jour, il faut un traitement spécial du retour chariot dans l'analyseur lexical (voir par exemple `lexer.mll` fourni dans le TD 3)

une façon de conserver l'information de localisation dans l'arbre de syntaxe abstraite est la suivante (cf cours 4)

```
type expression =  
  { desc: desc;  
    loc : Lexing.position * Lexing.position }  
  
and desc =  
  | Econst of int  
  | Eplus  of expression * expression  
  | Eneg   of expression  
  | ...
```

la grammaire peut donc ressembler à ceci

```
expression:
```

```
| d = desc { { desc = d; loc = $startpos, $endpos } }  
;
```

```
desc:
```

```
| i = INT                                { Econst i }  
| e1 = expression; PLUS; e2 = expression { Eplus (e1, e2) }  
| ...
```

comme dans le cas d'ocamllex, il faut s'assurer de l'application de menhir avant le calcul des dépendances

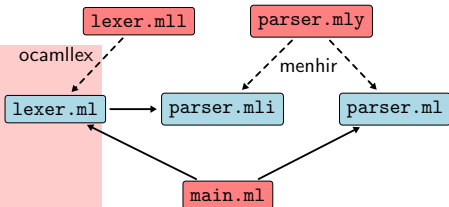
le Makefile ressemble donc à ceci :

```
lexer.ml: lexer.mll
    ocamllex lexer.mll

parser.mli parser.ml: parser.mly
    menhir -v parser.mly

.depend: lexer.ml parser.mli parser.ml
    ocamldep *.ml *.mli > .depend

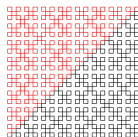
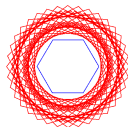
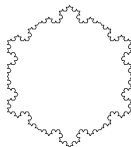
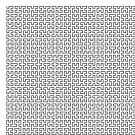
include .depend
```



(cf. Makefile fourni avec le TD 2 ou 3 / sinon utiliser ocamlbuild)

- TD cette semaine

utilisation d'ocamllex + menhir
sur un petit langage Logo
(tortue graphique)



- cours semaine prochaine
 - production de code