

quel point commun ?



Langages de programmation et compilation

École Normale Supérieure

Jean-Christophe Filliâtre

Cours 6 / 10 novembre 2011

Analyse syntaxique

l'objectif de l'analyse syntaxique est de reconnaître les phrases appartenant à la syntaxe du langage

son entrée est le flot des lexèmes construits par l'analyse lexicale, sa sortie est un arbre de syntaxe abstraite

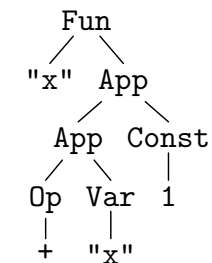
Analyse syntaxique

suite de lexèmes

fun x -> (x + 1)

↓
analyse syntaxique

↓
syntaxe abstraite



en particulier, l'analyse syntaxique doit détecter les erreurs de syntaxe et

- les localiser précisément
- les identifier (le plus souvent seulement « erreur de syntaxe » mais aussi « parenthèse non fermée », etc.)
- voire, reprendre l'analyse pour découvrir de nouvelles erreurs

pour l'analyse syntaxique, on va utiliser

- une **grammaire non contextuelle** pour décrire la syntaxe
- un **automate à pile** pour la reconnaître

c'est l'analogue des expressions régulières / automates finis utilisés dans l'analyse lexicale

on suppose connues les notions de grammaire, dérivation et automate à pile; cf. cours langages formels

dans ce qui suit, on note T l'ensemble des terminaux, N l'ensemble des non terminaux et S le symbole de départ d'une grammaire donnée

Exemple

pour un langage d'expressions arithmétiques construites à partir de constantes entières, d'additions, de multiplications et de parenthèses, on se donne la grammaire suivante :

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

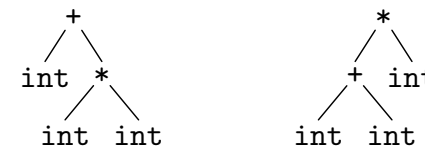
les terminaux de la grammaire sont les lexèmes produits par l'analyse lexicale, ici $T = \{+, *, (,), \text{int}\}$

note : `int` désigne le lexème correspondant à une constante entière (*i.e.* sa nature, pas sa valeur)

Ambiguïté

cette grammaire est **ambiguë** car certains mots admettent plusieurs arbres de dérivation

exemple : le mot `int + int * int` admet les deux arbres de dérivations

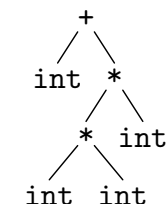


pour ce langage-là, il est néanmoins possible de proposer une autre grammaire, non ambiguë, qui définit le même langage

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

cette nouvelle grammaire traduit la priorité de la multiplication sur l'addition, et le choix d'une associativité à gauche pour ces deux opérations

ainsi, le mot `int + int * int * int` n'a plus qu'un seul arbre de dérivation, à savoir



correspondant à la dérivation gauche

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

Objectif

déterminer si une grammaire est ou non ambiguë n'est **pas décidable**

on va donc se restreindre à des classes de grammaires non ambiguës, et où on sait en outre construire un automate à pile **déterministe** correspondant

ces classes s'appellent LL(1), LR(0), SLR(1), LALR(1), LR(1), etc.

analyse descendante

idée : procéder par expansions successives du non terminal le plus à gauche (on construit donc une dérivation gauche) en partant de S et en se servant d'une **table** indiquant, pour un non terminal X à développer et les k premiers caractères de l'entrée, l'expansion $X \rightarrow \beta$ à effectuer (en anglais on parle de *top-down parsing*)

supposons $k = 1$ par la suite et notons $T(X, c)$ cette table

en pratique on suppose qu'un symbole terminal $\#$ dénote la fin de l'entrée, et la table indique donc également l'expansion de X lorsque l'on atteint la fin de l'entrée

on utilise une pile qui est un mot de $(N \cup T)^*$; initialement la pile est réduite au symbole de départ

à chaque instant, on examine le sommet de la pile et le premier caractère c de l'entrée

- si la pile est vide, on s'arrête; il y a succès si et seulement si c est $\#$
- si le sommet de la pile est un terminal a , alors a doit être égal à c , on dépile a et on consomme c ; sinon on échoue
- si le sommet de la pile est un non terminal X , alors on remplace X par le mot $\beta = T(X, c)$ en sommet de pile, le cas échéant, en empilant les caractères de β en partant du dernier; sinon, on échoue

Exemple

transformons encore la grammaire des expressions arithmétiques et considérons la table d'expansion suivante

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \\ &\quad | \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \\ &\quad | \epsilon \\ F &\rightarrow (E) \\ &\quad | \text{int} \end{aligned}$$

	+	*	()	int	#
E			TE'		TE'	
E'	+TE'			ϵ		ϵ
T			FT'		FT'	
T'	ϵ	*FT'		ϵ		ϵ
F			(E)		int	

Exemple

illustrons l'analyse descendante du mot

int + int * int

	+	*	()	int	#
E			TE'		TE'	
E'	+TE'			ϵ		ϵ
T			FT'		FT'	
T'	ϵ	*FT'		ϵ		ϵ
F			(E)		int	

pile	entrée
E	int+int*int#
E'T	int+int*int#
E'T'F	int+int*int#
E'T'int	int+int*int#
E'T'	+int*int#
E'	+int*int#
E'T+	+int*int#
E'T	int*int#
E'T'F	int*int#
E'T'int	int*int#
E'T'	*int#
E'T'F*	*int#
E'T'F	int#
E'T'int	int#
E'T'	#
E'	#
ϵ	#

un analyseur descendant se programme très facilement en introduisant une fonction pour chaque non terminal de la grammaire

chaque fonction examine l'entrée et, selon le cas, la consomme ou appelle récursivement les fonctions correspondant à d'autres non terminaux, selon la table d'expansion

faisons le choix d'une programmation purement applicative, où l'entrée est une liste de lexèmes du type

```
type token = Tplus | Tmult | Tleft | Tright | Tint | Teof
```

on va donc construire cinq fonctions qui « consomment » la liste d'entrée

```
val e : token list → token list
val e' : token list → token list
val t : token list → token list
val t' : token list → token list
val f : token list → token list
```

et la reconnaissance d'une entrée pourra alors se faire ainsi

```
let recognize l = e l = [Teof];;
```

les fonctions procèdent par filtrage sur l'entrée, et suivent la table

	+	*	()	int	#
<i>E</i>			<i>TE'</i>		<i>TE'</i>	

```
let rec e = function
| (Tleft | Tint) :: _ as m → e' (t m)
| _ → error ()
```

	+	*	()	int	#
<i>E'</i>	+ <i>TE'</i>			ε		ε

```
and e' = function
| Tplus :: m → e' (t m)
| (Tright | Teof) :: _ as m → m
| _ → error ()
```

	+	*	()	int	#
<i>T</i>			<i>FT'</i>		<i>FT'</i>	

```
and t = function
| (Tleft | Tint) :: _ as m → t' (f m)
| _ → error ()
```

	+	*	()	int	#
<i>T'</i>	ε	* <i>FT'</i>		ε		ε

```
and t' = function
| (Tplus | Tright | Teof) :: _ as m → m
| Tmult :: m → t' (f m)
| _ → error ()
```

	+	*	()	int	#
F			(E)		int	

```

and f = function
  | Tint :: m → m
  | Tleft :: m → begin match e m with
    | Tright :: m → m
    | _ → error ()
  end
  | _ → error ()
    
```

remarques

- la table d'expansion n'est pas explicite : elle est dans le code de chaque fonction
- la pile non plus n'est pas explicite : elle est réalisée par la pile d'appels
- on aurait pu les rendre explicites
- on aurait pu aussi faire le choix d'une programmation plus impérative

```
val next_token : unit → token
```

Construire la table d'expansion

reste une question d'importance : comment construire la table ?

l'idée est simple : pour décider si on réalise l'expansion $X \rightarrow \beta$ lorsque le premier caractère de l'entrée est c , on va chercher à déterminer si c fait partie des **premiers** caractères des mots reconnus par β

une difficulté se pose pour une production telle que $X \rightarrow \epsilon$, et il faut alors considérer aussi l'ensemble des caractères qui peuvent **suivre** X

déterminer les premiers et les suivants nécessite également de déterminer si un mot peut se dériver en ϵ

Définitions

Définition (NULL)

Soit $\alpha \in (T \cup N)^*$. $\text{NULL}(\alpha)$ est vrai si et seulement si on peut dériver ϵ à partir de α i.e. $\alpha \xrightarrow{*} \epsilon$.

Définition (FIRST)

Soit $\alpha \in (T \cup N)^*$. $\text{FIRST}(\alpha)$ est l'ensemble de tous les premiers terminaux des mots dérivés de α , i.e. $\{a \in T \mid \exists w. \alpha \xrightarrow{*} aw\}$.

Définition (FOLLOW)

Soit $X \in N$. $\text{FOLLOW}(X)$ est l'ensemble de tous les terminaux qui peuvent apparaître après X dans une dérivation, i.e. $\{a \in T \mid \exists u, w. S \xrightarrow{*} uXaw\}$.

pour calculer $NULL(\alpha)$ il suffit de déterminer $NULL(X)$ pour $X \in N$

$NULL(X)$ est vrai si et seulement si

- il existe une production $X \rightarrow \epsilon$,
- ou il existe une production $X \rightarrow Y_1 \dots Y_m$ où $NULL(Y_i)$ pour tout i

problème : il s'agit donc d'un ensemble d'équations mutuellement récursives

dit autrement, si $N = \{X_1, \dots, X_n\}$ et si $\vec{V} = (NULL(X_1), \dots, NULL(X_n))$, on cherche la **plus petite solution** d'une équation de la forme

$$\vec{V} = F(\vec{V})$$

Théorème (existence d'un plus petit point fixe (Tarski))

Soit A un ensemble fini muni d'une relation d'ordre \leq et d'un plus petit élément ϵ . Toute fonction $f : A \rightarrow A$ monotone, i.e. telle que $\forall x, y. x \leq y \Rightarrow f(x) \leq f(y)$, admet un plus petit point fixe.

preuve : comme ϵ est le plus petit élément, on a $\epsilon \leq f(\epsilon)$
 f étant monotone, on a donc $f^k(\epsilon) \leq f^{k+1}(\epsilon)$ pour tout k
 A étant fini, il existe donc un plus petit k_0 tel que $f^{k_0}(\epsilon) = f^{k_0+1}(\epsilon)$
 $a_0 = f^{k_0}(\epsilon)$ est donc un point fixe de f

soit b un autre point fixe de f
 on a $\epsilon \leq b$ et donc $f^k(\epsilon) \leq f^k(b)$ pour tout k
 en particulier $a_0 = f^{k_0}(\epsilon) \leq f^{k_0}(b) = b$
 a_0 est donc le plus petit point fixe de f □

Calcul de NULL

Exemple

dans le cas du calcul de NULL, on a $A = \text{BOOL} \times \dots \times \text{BOOL}$ avec $\text{BOOL} = \{\text{false}, \text{true}\}$

on peut munir BOOL de l'ordre $\text{false} \leq \text{true}$ et A de l'ordre point à point

$$(x_1, \dots, x_n) \leq (y_1, \dots, y_n) \text{ si et seulement si } \forall i. x_i \leq y_i$$

le théorème s'applique en prenant $\epsilon = (\text{false}, \dots, \text{false})$

pour calculer les $NULL(X_i)$, on part donc de $NULL(X_1) = \text{false}, \dots, NULL(X_n) = \text{false}$ et on applique les équations jusqu'à obtention du point fixe

```

E  →  T E'
E' →  + T E'
    |  ε
T  →  F T'
T' →  * F T'
    |  ε
F  →  ( E )
    |  int
    
```

E	E'	T	T'	F
false	false	false	false	false
false	true	false	true	false
false	true	false	true	false

de même, les équations définissant FIRST sont mutuellement récursives

$$\text{FIRST}(X) = \bigcup_{X \rightarrow \beta} \text{FIRST}(\beta)$$

et

$$\begin{aligned} \text{FIRST}(a\beta) &= \{a\} \\ \text{FIRST}(X\beta) &= \text{FIRST}(X), \quad \text{si } \neg \text{NULL}(X) \\ \text{FIRST}(X\beta) &= \text{FIRST}(X) \cup \text{FIRST}(\beta), \quad \text{si } \text{NULL}(X) \end{aligned}$$

de même, on procède par calcul de point fixe sur le produit cartésien $A = \mathcal{P}(T) \times \dots \times \mathcal{P}(T)$ muni, point à point, de l'ordre \subseteq et avec $\varepsilon = (\emptyset, \dots, \emptyset)$

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \\ &\quad | \quad \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \\ &\quad | \quad \epsilon \\ F &\rightarrow (E) \\ &\quad | \quad \text{int} \end{aligned}$$

NULL

E	E'	T	T'	F
false	true	false	true	false

FIRST

E	E'	T	T'	F
\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
\emptyset	{+}	\emptyset	{*}	{(, int}
\emptyset	{+}	{(, int}	{*}	{(, int}
{(, int}	{+}	{(, int}	{*}	{(, int}
{(, int}	{+}	{(, int}	{*}	{(, int}

là encore, les équations définissant FOLLOW sont mutuellement récursives

$$\text{FOLLOW}(X) = \bigcup_{Y \rightarrow \alpha X \beta} \text{FIRST}(\beta) \cup \bigcup_{Y \rightarrow \alpha X \beta, \text{NULL}(\beta)} \text{FOLLOW}(Y)$$

on procède par calcul de point fixe, sur le même domaine que pour FIRST

note : il faut introduite # dans les suivants du symbole de départ (ce que l'on peut faire directement, ou en ajoutant une règle $S' \rightarrow S\#$)

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \\ &\quad | \quad \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \\ &\quad | \quad \epsilon \\ F &\rightarrow (E) \\ &\quad | \quad \text{int} \end{aligned}$$

NULL

E	E'	T	T'	F
false	true	false	true	false

FIRST

E	E'	T	T'	F
{(, int}	{+}	{(, int}	{*}	{(, int}

FOLLOW

E	E'	T	T'	F
{#}	\emptyset	\emptyset	\emptyset	\emptyset
{#,)}	{#}	{+, #}	\emptyset	{*}
{#,)}	{#,)}	{+, #,)}	{+, #}	{*, +, #}
{#,)}	{#,)}	{+, #,)}	{+, #,)}	{*, +, #,)}
{#,)}	{#,)}	{+, #,)}	{+, #,)}	{*, +, #,)}

à l'aide des premiers et des suivants, on construit la table d'expansion $T(X, a)$ de la manière suivante

pour chaque production $X \rightarrow \beta$,

- on pose $T(X, a) = \beta$ pour tout $a \in \text{FIRST}(\beta)$
- si $\text{NULL}(\beta)$, on pose aussi $T(X, a) = \beta$ pour tout $a \in \text{FOLLOW}(X)$

E	\rightarrow	TE'
E'	\rightarrow	$+TE'$
		$\mid \epsilon$
T	\rightarrow	FT'
T'	\rightarrow	$*FT'$
		$\mid \epsilon$
F	\rightarrow	(E)
		$\mid \text{int}$

FIRST				
E	E'	T	T'	F
$\{(, \text{int})\}$	$\{+\}$	$\{(, \text{int})\}$	$\{*\}$	$\{(, \text{int})\}$
FOLLOW				
E	E'	T	T'	F
$\{\#, \,)\}$	$\{\#, \,)\}$	$\{+, \#, \,)\}$	$\{+, \#, \,)\}$	$\{*, +, \#, \,)\}$

	+	*	()	int	#
E			TE'		TE'	
E'	$+TE'$			ϵ		ϵ
T			FT'		FT'	
T'	ϵ	$*FT'$		ϵ		ϵ
F			(E)		int	

Grammaire LL(1)

Construire des grammaires LL(1)

Définition (grammaire LL(1))

Une grammaire est dite LL(1) si, dans la table précédente, il y a au plus une production dans chaque case.

LL signifie « Left to right scanning, Leftmost derivation »

il faut souvent transformer les grammaires pour les rendre LL(1)

en particulier, une grammaire **récursive gauche**, i.e. contenant une production de la forme

$$X \rightarrow X\alpha$$

ne sera jamais LL(1)

il faut alors supprimer la récursion gauche (directe ou indirecte)

de même il faut factoriser les productions qui commencent par le même terminal (factorisation gauche)

les analyseurs LL(1) sont relativement simples à écrire
mais ils nécessitent d'écrire des grammaires peu naturelles

on va se tourner vers une autre solution...

analyse ascendante

Idée générale

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*)

Fonctionnement de l'analyse

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

Analyse LR

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 $X s$, où s est l'état résultat de la transition $s_{n-p} \xrightarrow{X} s$, i.e.

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

Analyse LR

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 **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 **goto**(s, X) indique l'état résultat d'une réduction de X

commençons par $k = 0$

on construit un automate LR(0) non déterministe

dont les états sont des items de la forme

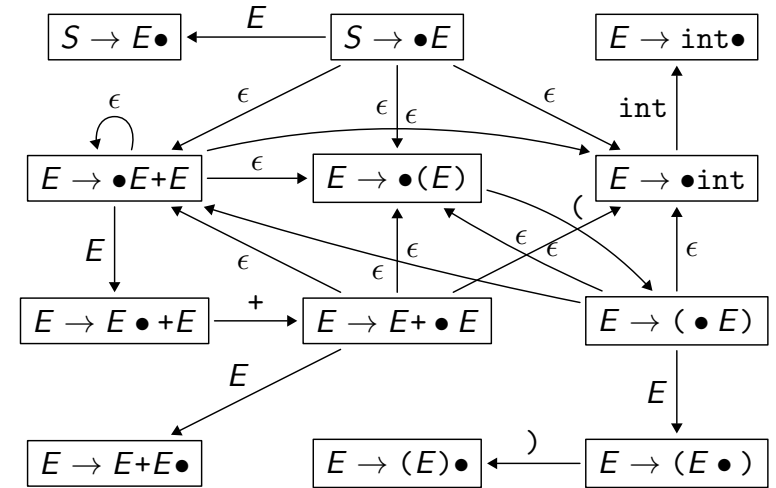
$$[X \rightarrow \alpha \bullet \beta]$$

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 »

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

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

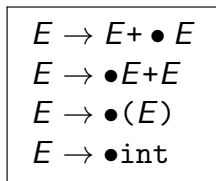
$$\begin{aligned}
 S &\rightarrow E \\
 E &\rightarrow E+E \\
 &\quad (E) \\
 &\quad \text{int}
 \end{aligned}$$



déterminisons l'automate LR(0)

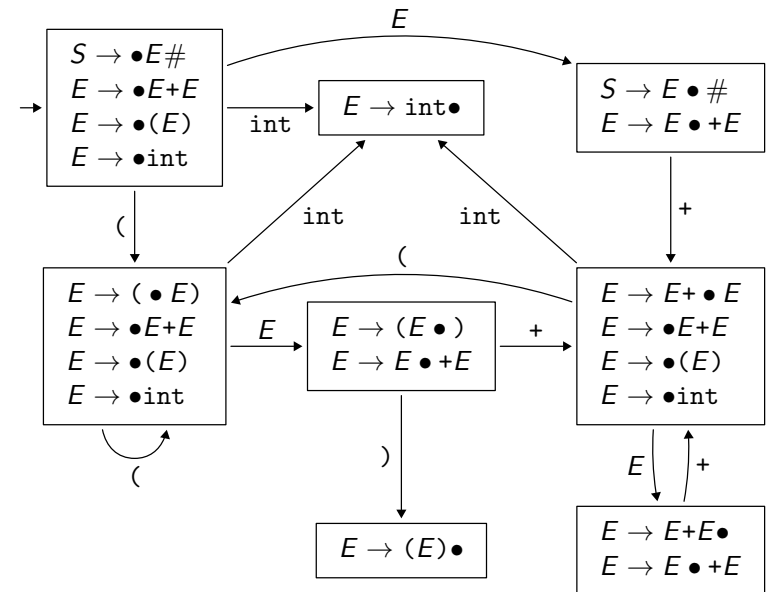
pour cela, on regroupe les états reliés par des ϵ -transitions

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



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

$$\begin{aligned}
 S &\rightarrow E \\
 E &\rightarrow E+E \\
 &\quad (E) \\
 &\quad \text{int}
 \end{aligned}$$



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 $s \xrightarrow{X} s'$

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.

Exemple

sur notre exemple, la table est la suivante :

état	action					goto
	()	+	int	#	E
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 $E \rightarrow E+E$					

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 → int					
3			s6		ok	
4	s4			s2		5
5		s7	s6			
6	s4			s2		8
7	reduce E → (E)					
8	reduce E → E+E					

pile	entrée	action
1	int+int+int	s2
1 int 2	+int+int	E → int, g3
1 E 3	+int+int	s6
1 E 3 + 6	int+int	s2
1 E 3 + 6 int 2	+int	E → int, g8
1 E 3 + 6 E 8	+int	E → E+E, g3
1 E 3	+int	s6
1 E 3 + 6	int	s2
1 E 3 + 6 int 2	#	E → int, g8
1 E 3 + 6 E 8	#	E → 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 $action(s, a) = reduce X \rightarrow \beta$ si et seulement si

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

Définition (classe SLR(1))

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

(SLR signifie Simple LR)

Exemple

la grammaire

$S \rightarrow E\#$
 $E \rightarrow E + T$
 $\quad | T$
 $T \rightarrow T * F$
 $\quad | F$
 $F \rightarrow (E)$
 $\quad | int$

est SLR(1)

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

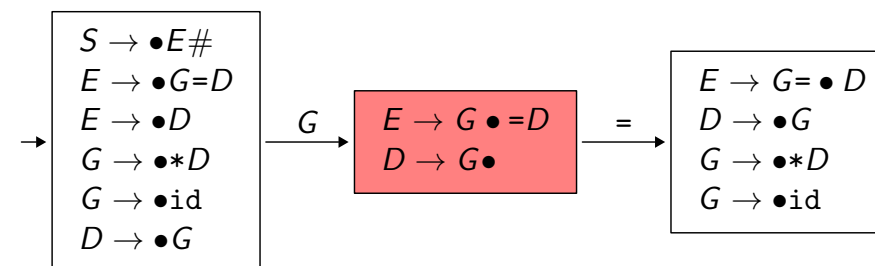
Limites de l'analyse SLR(1)

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

exemple :

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

	=	
1
2	shift 3 reduce D → 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 $[X \rightarrow \alpha \bullet, a]$

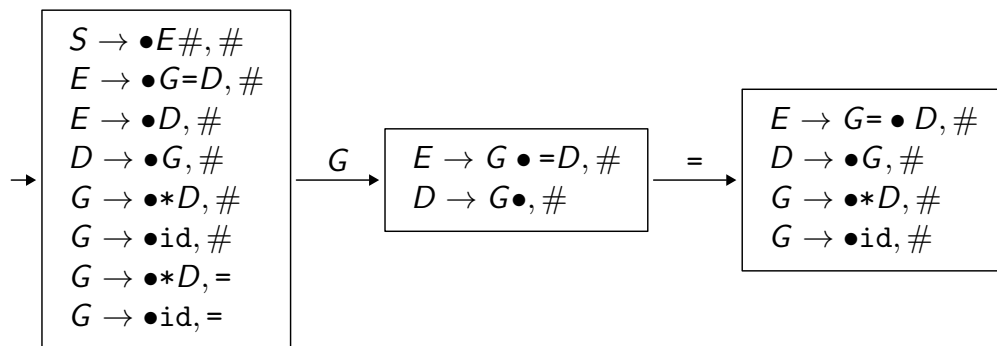
Définition (classe LR(1))

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

Exemple

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

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

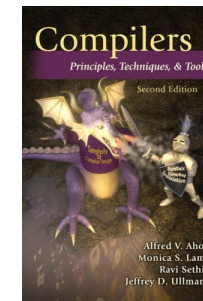
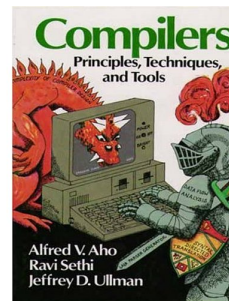


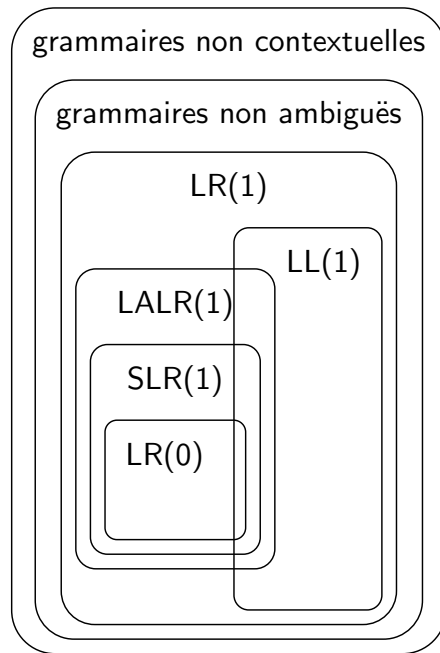
LALR(1)

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





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, ocaml yacc, cups, menhir, ...
(YACC signifie *Yet Another Compiler Compiler*)

Menhir

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 (tel 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 }
;
```

Exemple

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

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

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

la déclaration d'un type `token`

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

et, 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 précédent)

Conflits

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 && ocamlc -c -i arith.ml
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
```

Résolution des conflits

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

pour cela, on donne 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)

les priorités des lexèmes sont données explicitement par l'utilisateur

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
```

Résolution des conflits

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
```

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 !

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.ml1` fourni dans le TD 3)

une façon de conserver l'information de localisation dans l'arbre de syntaxe abstraite est la suivante

```
type expression =
  { expr_desc: expr_desc;
    expr_loc  : Lexing.position × Lexing.position }

and expr_desc =
  | Econst of int
  | Eplus  of expression × expression
  | Eneg   of expression
  | ...
```

Plus d'information

pour plus d'informations, consulter le manuel de Menhir, accessible depuis la page du cours

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)

La suite

- TD de mercredi
 - analyseur LL(1)

- Cours de jeudi
 - production de code