

Analyse syntaxique

Grammaires

Langage et syntaxe

- ▶ Les phrases d'un langage ont une structure déterminée par une grammaire

- ▶ **Exemple** : une phrase correcte est constituée par un sujet suivi d'un verbe

`if (expression) instruction else instruction`

alternative

- ▶ Ceci peut être exprimé par une grammaire

phrase → sujet verbe

instr → if (expr) instr else instr

- ▶ Si on complète avec deux nouvelles règles

sujet → pierre | claire

verbe → court | marche

- ▶ 4 phrases possibles :

pierre court | pierre marche | claire court | claire marche

Exemple 1

Nombre -> Chiffre | Chiffre Nombre

Chiffre -> 0|1|2|3|4|5|6|7|8|9

gènère

0

17

547

7350

etc...

Exemple 2 bloc de code en Java

En Pascal, est une suite d'instructions séparées par des points-virgules entre une { et }

bloc \rightarrow { *liste_opt_instr* }

liste_opt_instr \rightarrow *liste_instr* | ϵ

liste_instr \rightarrow *liste_instr* ; *instr* | *instr*

liste des symboles vides

Grammaire non contextuelle

Une grammaire non contextuelle comprend quatre composants $G=(\Sigma, V, P, S)$:

- ▶ un ensemble d'unités lexicales appelées symboles terminaux Σ (lexèmes ou tokens) ;
 - ▶ un ensemble de non-terminaux V (disjoint de Σ) ;
 - ▶ un ensemble de règles de production P où chaque règle est constituée d'un non-terminal, appelé *partie gauche*, d'une flèche appelée *partie droite* de la forme
 $\alpha \rightarrow \beta$ avec $\alpha \in V$, et $\beta \in (\Sigma \cup V)^*$
 - ▶ La désignation d'un des non-terminaux en tant que *symbole de départ* S .
-

Grammaire non contextuelle (suite)

- ▶ les terminaux sont les symboles de base à partir desquels les chaînes sont formées (\Leftrightarrow unités lexicales)
 - ▶ les non-terminaux sont les variables syntaxiques qui dénotent un ensemble de chaînes
 - ▶ l'axiome est un non-terminal particulier. L'ensemble des chaînes que l'on peut obtenir à partir de l'axiome est le langage défini par la grammaire
-

Grammaire non contextuelle (suite)

Une grammaire définit un langage formé par l'ensemble des séquences finies de symboles terminaux qui peuvent être dérivées du symbole initial par des applications successives de règles

Convention :

- ▶ les chiffres, les signes ou les chaînes en **gras** sont des terminaux
 - ▶ un nom en *italique* désigne un non-terminal
-

Dérivations

Pour montrer qu'un mot w est dans le langage L , il faut exhiber une suite finie d'applications de règles de production partant de S et aboutissant à w , appelée dérivation

On peut représenter une dérivation par

$$S = w_0 \xrightarrow{P_1} w_1 \xrightarrow{P_2} \dots \xrightarrow{P_n} w_n = w$$

Cela peut également se représenter sous la forme d'un arbre que l'on nomme (*arbre de dérivation*)

Exemple de grammaire

phrase → *GN GV*

GN → *Article Nom* | *Article Nom Subord*

Subord → *Pron_relatif GV*

GV → *Verbe GN* | *Verbe GN_prep*

GN_prep → *Prep GN*

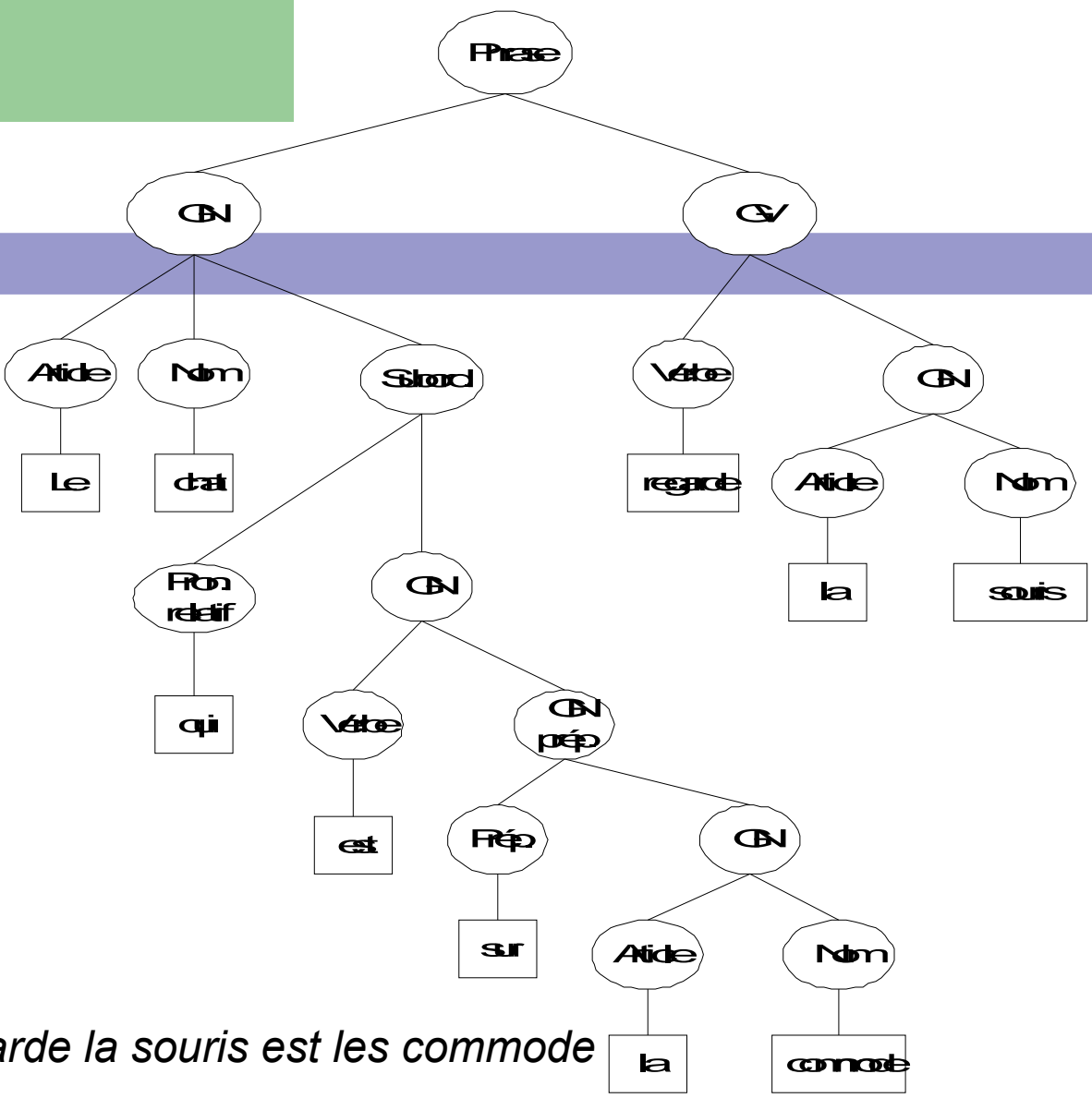
Article → **le** | **la** | **les**

Nom → **chat** | **commode** | **souris**

Verbe → **est** | **regarde**

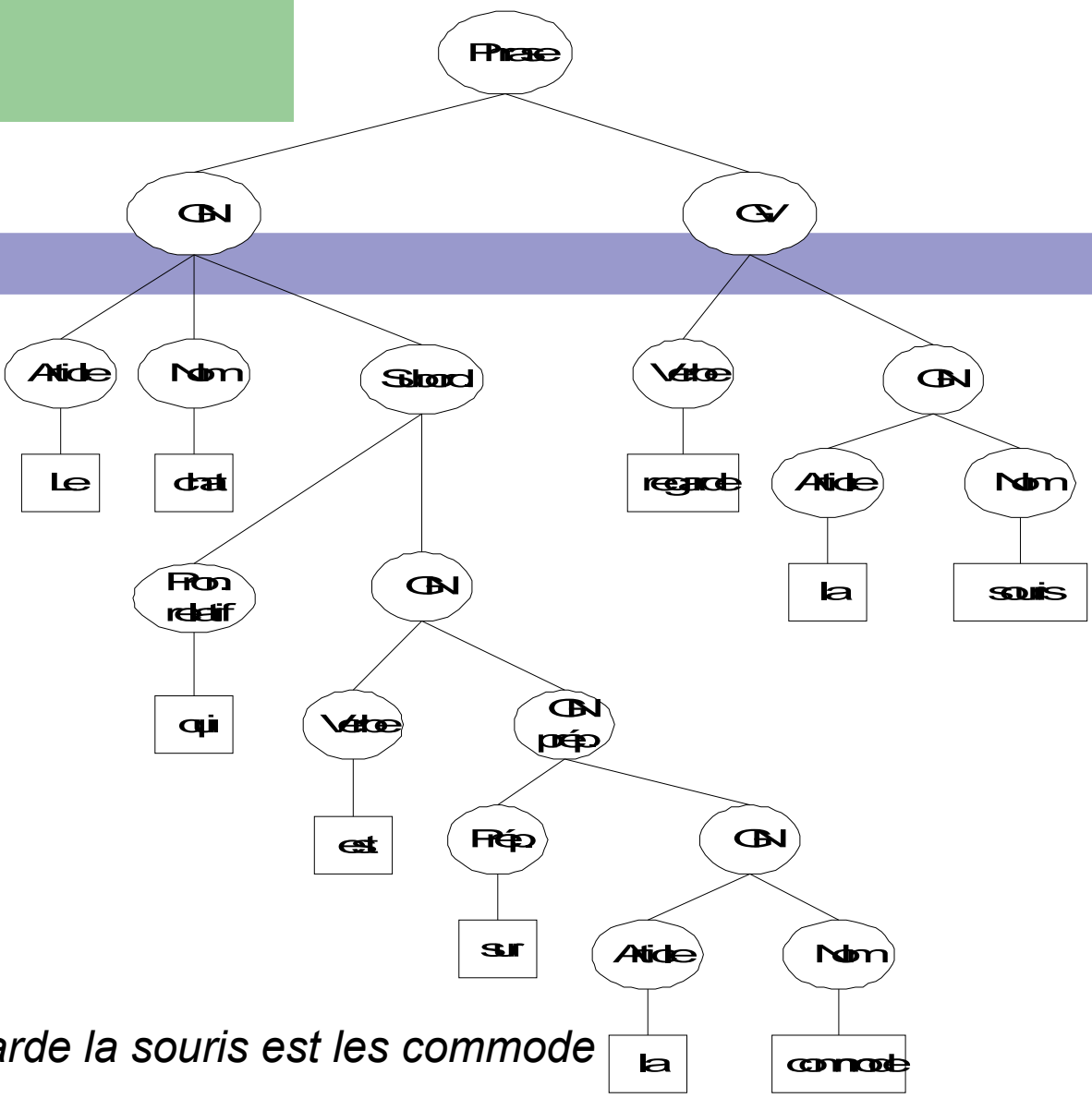
Prep → **sur** | **sous**

Pron_relatif → **qui** | **que** | **dont** | **ou**



le chat qui regarde la souris est les commode

la commode



le chat qui regarde la souris est les commode

la commode

Langage engendré

Le langage engendré par une grammaire G est l'ensemble de terminaux que l'on peut dériver à partir de l'axiome

$$\forall w \in \Sigma^* \quad w \in L(G) \Leftrightarrow S \xrightarrow{+} w$$

Dérivation à gauche – à droite

- ▶ Une dérivation est dite à gauche si le non terminal le plus à gauche est remplacé



- ▶ Une dérivation est dite à droite si le non-terminal le plus à droite est remplacé



Exemple

G :

$E \rightarrow E + E / E * E / (E) / - E / Id$

$-(Id + Id)$ appartient $L(G)$

$E \rightarrow -E \rightarrow -(E) \rightarrow -(E + E)$

1er cas :

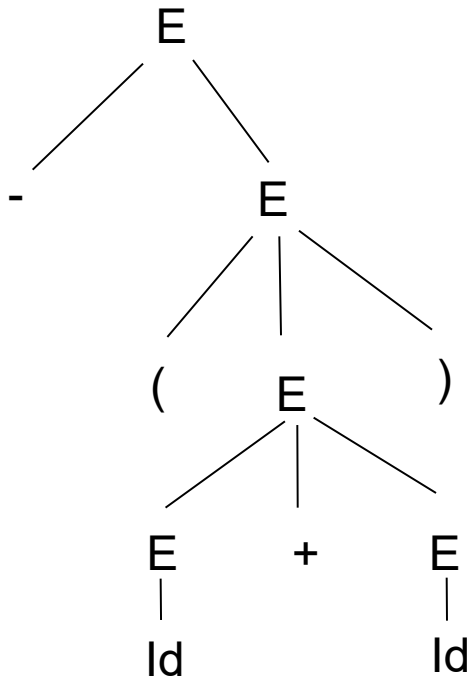
$\rightarrow -(Id + E) \rightarrow -(Id + Id)$

2ème cas :

$\rightarrow -(E + Id) \rightarrow -(Id + Id)$

Arbre de dérivation

Un arbre de dérivation est une représentation graphique d'une dérivation.



A chaque arbre d'analyse, est associé une unique dérivation gauche (ou droite)

Les deux suites différentes de dérivations donnent le même arbre de dérivation.

Ambiguïté



Une grammaire est ambiguë si elle produit plus d'un arbre d'analyse (ou d'une dérivation à gauche) pour une phrase donnée.

Remarque : la grammaire précédente **est** ambiguë

Exemple

Soit la grammaire G

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow -E$

$E \rightarrow \text{id}$

Il existe deux arbres d'analyse syntaxique pour $\text{id} + \text{id} * \text{id}$

$E \rightarrow E + E$

$\rightarrow \text{id} + E$

$\rightarrow \text{id} + E * E$

$\rightarrow \text{id} + \text{id} * E$

$\rightarrow \text{id} + \text{id} * \text{id}$

ou

$E \rightarrow E * E$

$\rightarrow E + E * E$

$\rightarrow \text{id} + E * E$

$\rightarrow \text{id} + \text{id} * E$

$\rightarrow \text{id} + \text{id} * \text{id}$

Exemples : expressions arithmétiques

Les terminaux (lexèmes) sont **NUM**, **ID**, **+** et **-**

$$S \rightarrow E$$

$$E \rightarrow \text{NUM}$$

$$E \rightarrow \text{ID}$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

Pour reconnaître l'expression $1 - 1 + x$, plusieurs dérivations sont possibles (on parle d'ambiguïté)

Dérivations possibles

$$\begin{array}{c}
 S \rightarrow E \rightarrow \left\{ \begin{array}{l} E \rightarrow \left\{ \begin{array}{l} E \rightarrow \text{NUM} \\ - \\ E \rightarrow \text{NUM} \end{array} \right. \\ + \\ E \rightarrow \text{ID} \end{array} \right.
 \end{array}
 \quad
 \begin{array}{c}
 S \rightarrow E \rightarrow \left\{ \begin{array}{l} E \rightarrow \text{NUM} \\ - \\ E \rightarrow \left\{ \begin{array}{l} E \rightarrow \text{NUM} \\ + \\ E \rightarrow \text{ID} \end{array} \right. \end{array} \right.
 \end{array}$$

Exercice

Trouvez une grammaire G non ambiguë engendrant le même langage que G (* plus prioritaire que +)

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * \text{id} \mid \text{id}$$

On obtient une grammaire non ambiguë mais moins lisible

Processus d'analyse

Le but est reconstruire l'arbre syntaxique d'une phrase à partir d'une grammaire non ambiguë. Il existe deux processus :

- ▶ l'analyse descendante. Cette analyse correspond à un parcours descendant gauche (on lit de gauche à droite). Il s'agit de deviner à chaque étape de deviner la règle qui sert à engendrer le mot qu'on lit.
 - ▶ l'analyse ascendante. Dans cette analyse, on essaye d'appliquer le membre droit d'une règle et à le remplacer par le non-terminal gauche correspondant. C'est l'opération inverse de la dérivation.
-

L'analyse descendante

- ▶ Tente de construire un arbre syntaxique pour une chaîne d'entrée de la racine vers les feuilles
ou
 - ▶ Tente de déterminer une dérivation gauche associée à une chaîne d'entrée
-

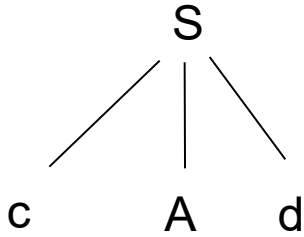
Exemple

$$\left| \begin{array}{l} S \rightarrow cAd \\ A \rightarrow ab \mid a \\ w = cad \end{array} \right.$$

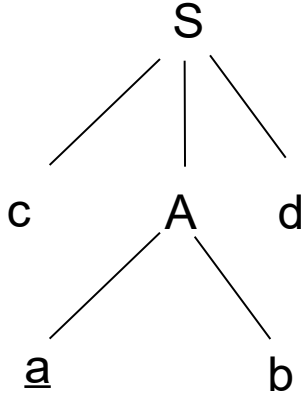
c a d
↑ ↑

mémorisé

c a d



$S \Rightarrow c \underline{A} d$

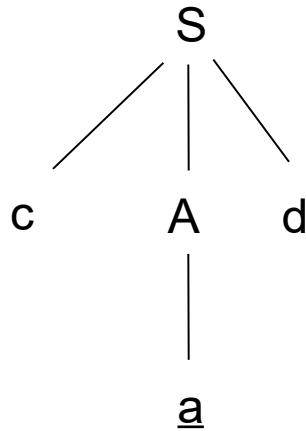


$\Rightarrow c a b d$

échec
retour au choix précédent

Exemple (suite)

c a d
↑
retour



$$S \Rightarrow c A d \Rightarrow c a d$$

Dans cet exemple, nous traitons une forme générale d'analyse descendante où des retours arrières sont possibles avec passages répétés sur le texte source. Dans la pratique, le rebroussement est rarement nécessaire, on peut adapter les grammaires pour effectuer une analyse descendante sans rebroussement

Remarque

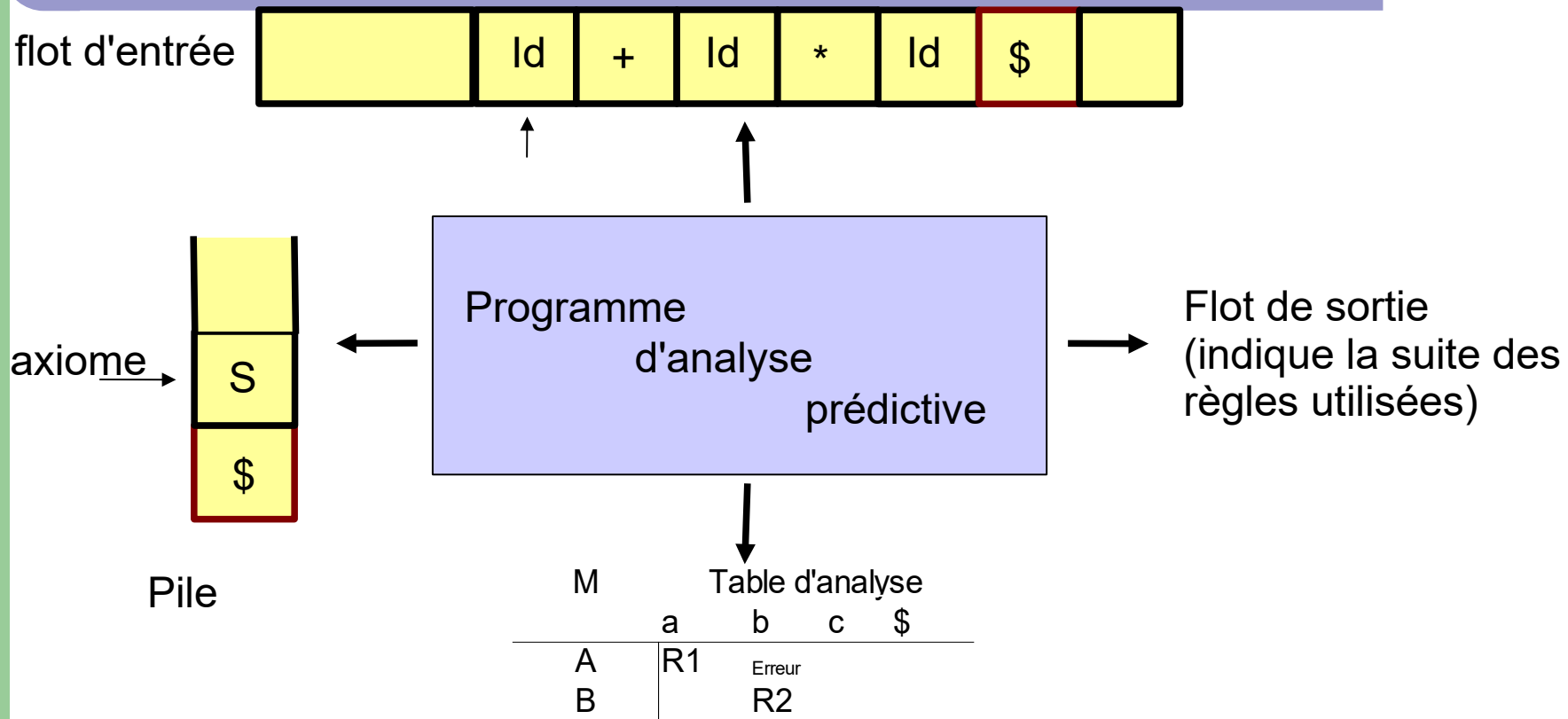
$S \rightarrow aSb | aSc | d$ avec le mot $w = aaaaaaadbbbbb$

- ▶ Pour savoir quelle règle utiliser, il faut cette fois-ci connaître aussi la dernière lettre du mot.
 - ▶ Conclusion : ce qui serait pratique, ça serait d'avoir une table qui nous dit : quand je lis tel caractère et que j'en suis à dériver tel symbole non-terminal, alors j'applique telle règle et je ne me pose pas de questions. Ça existe, et ça s'appelle **une table d'analyse**.
-

Analyse descendante

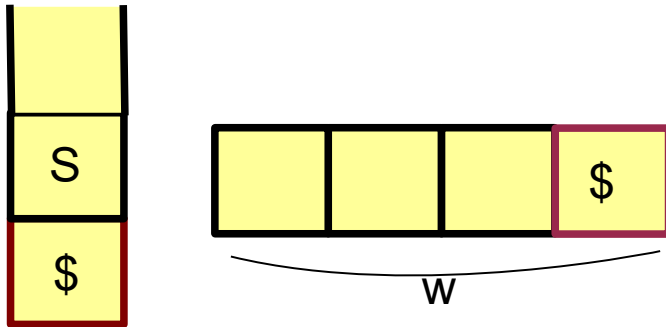
- Analyse prédictive
-

Implémentation à l'aide d'une pile explicite



Algorithme

au départ :



si $X=a$ alors Dépiler(X);
Avancer;

si X est un non terminal

si $M[X,a]$ contient une
production $X \rightarrow UVW$

terminal ou
non-terminal

Remplacer X pour WVU

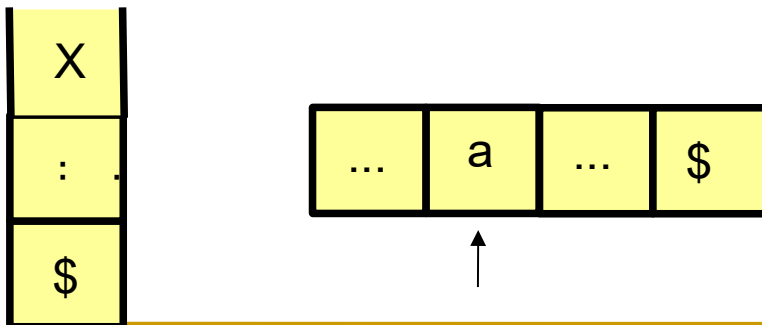
si $M[X,a]$ vide

au fond

Routine d'erreur

si $X=a=\$$ alors succès

en cours :



Exemple

Soit la grammaire suivante :

$$E \rightarrow TE' [r1]$$

$$E' \rightarrow +TE' [r2] \mid \varepsilon [r3]$$

$$T \rightarrow FT' [r4]$$

$$T' \rightarrow *FT' [r5] \mid \varepsilon [r6]$$

$$F \rightarrow (E) [r7] \mid \text{nb} [r8]$$

Table d'analyse

	nb	+	*	()	\$
E	<i>r1</i>			<i>r1</i>		
E'		<i>r2</i>			<i>r3</i>	<i>r3</i>
T	<i>r4</i>			<i>r4</i>		
T'		<i>r6</i>	<i>r5</i>		<i>r6</i>	<i>r6</i>
F	<i>r8</i>			<i>r7</i>		

E	$\underline{2}+3*4\$$	r1
TE'	$2+\underline{3}*4\$$	r4
FT'E'	$2+3*\underline{4}\$$	r8
T'E'	$2+\underline{3}*4\$$	r6
E'	$2+\underline{3}*4\$$	r2
TE'	$2+\underline{3}*4\$$	r4
FT'E'	$2+\underline{3}*4\$$	r8
T'E'	$2+3*\underline{4}\$$	r5
FT'E'	$2+3*\underline{4}\$$	r8
T'E'	$2+3*\underline{4}\$$	r6
E'	$2+3*\underline{4}\$$	r3
0	succès	

Fonctionnement de la table

Elle permet d'expliciter les dérivations à effectuer pour accepter une phrase d'une grammaire

Exemple :

dérivation de $w = 2 + 3 * 4$ \$

symbole de fin

Résumé

$$\begin{aligned}
 E &\Rightarrow \underline{T} E' \Rightarrow \underline{F} T' E' \Rightarrow id T' E' \Rightarrow id E' \Rightarrow \\
 &\Rightarrow id + T E' \Rightarrow id + F T' E' \Rightarrow id + id T' E' \Rightarrow id + id * F T' E' \\
 &\Rightarrow id + id * id T' E' \Rightarrow id + id * id E' \Rightarrow id + id * id
 \end{aligned}$$

Un peu de théorie ☹

les terminaux avant α

grosso modo

définition : pour $\alpha \in (\Sigma \cup V)^*$

$$first(\alpha) = \{x \in \Sigma \mid \alpha \xrightarrow{*} x\beta\} \cup \{\varepsilon \text{ si } \alpha \text{ est une chaîne nullable}\}$$

ensemble des terminaux qui commencent
les chaînes qui se dérivent de α

les terminaux après A

On définit pour $A \in V$:

$$follow(A) = \{x \in \Sigma \mid S \xrightarrow{*} \alpha A \gamma \text{ avec } x \in first(\gamma)\}$$

ensemble des terminaux qui peuvent se trouver immédiatement à droite
de A après dérivation. Si A est le symbole le plus à droite, alors

$\$ \in \sigma_{\cup i \in \alpha \cup \tau}(A)$

Définition de premier

- ▶ si X est un terminal
 $\text{premier}(X) = \{X\}$
- ▶ Si $X \rightarrow \epsilon$
 $\epsilon \in \text{premier}(X)$
- ▶ si $X \rightarrow Y_1 Y_2 \dots Y_k$
 $\text{premier}(Y_1) \subset \text{premier}(X),$

si $Y_1 \xRightarrow{*} \epsilon, \text{ first}(Y_2) \subset \text{first}(x)$

si $Y_2 \xRightarrow{*} \epsilon \dots$

- ▶ Pour une chaîne
 $X_1 X_2 \dots X_n$
 $\text{premier}(X_1 X_2 \dots X_n) = \text{premier}(X_1)$
 si $\epsilon \in \text{premier}(X_1) +$
 $\text{premier}(X_2)$ si
 $\epsilon \in \dots$

Application

$\text{premier}(E) = \{ (, \text{nb} \}$

$\text{premier}(E') = \{ +, \varepsilon \}$

$\text{premier}(T) = \{ (, \text{nb} \}$

$\text{premier}(T') = \{ *, \varepsilon \}$

$\text{premier}(F) = \{ (, \text{nb} \}$

Calcul de suivant

1. s'il y a une production $A \rightarrow \alpha B \beta$, le contenu de $\text{premier}(\beta)$, excepté ε , est ajouté à $\text{suivant}(B)$
 2. mettre $\$$ dans $\text{suivant}(S)$, où S est l'axiome et $\$$ le marqueur de fin de texte
 3. Pour calculer $\text{suivant}(A)$ pour tous les non-terminaux A , appliquer les règles suivantes jusqu'à ce qu'aucun terminal ne puisse être ajouté
 4. s'il existe une production $A \rightarrow \alpha B$ ou $A \rightarrow \alpha B \beta$ tq $\text{premier}(\beta)$ contient ε ($\beta \rightarrow^* \varepsilon$), les éléments de $\text{suivant}(A)$ sont ajoutés à $\text{suivant}(B)$
-

Définition de suivant

Pour l'axiome S $\$ \in \text{follow}(S)$
 si $A \rightarrow \alpha B \beta$ $\text{first}(\beta) - \{\epsilon\} \subset \text{follow}(B)$
 si $A \rightarrow \alpha B$
 ou $A \rightarrow \alpha B \beta$ avec $\beta \Rightarrow^* \epsilon$

} $\text{follow}(A) \subset \text{follow}(B)$

Exemple

$E \rightarrow TE' [r1]$

$E' \rightarrow +TE' [r2] \mid \varepsilon [r3]$

$T \rightarrow FT' [r4]$

$T' \rightarrow *FT' [r5] \mid \varepsilon [r6]$

$F \rightarrow (E) [r7] \mid \text{nb} [r8]$

$\text{premier}(E) = \text{premier}(T) = \text{premier}(F) = \{ (, \text{nb} \}$

$\text{premier}(E') = \{ +, \varepsilon \}$

$\text{premier}(T') = \{ *, \varepsilon \}$

$\text{suivant}(E) = \{ \$,) \}$

$\text{suivant}(T) = \{ +, \$,) \}$

$\text{suivant}(F) = \{ *, +, \$,) \}$

$\text{suivant}(E') = \{ \$,) \}$

$\text{suivant}(T') = \{ +, \$,) \}$

règles $\alpha B \beta$

r1 $B=T$ $\beta=E'$ rajout de +

r2 *idem*

r4 $B=F$ $\beta=T'$ rajout de *

r5 $\alpha=* B=F$ $\beta=T'$ *idem*

r7 $\alpha=(B=E$ $\beta=)$ rajout de)

règles αB ou $\alpha B \beta$ β nullable

r1 $\alpha=T$ $B=E'$ rajout de $\text{suivant}(E)$

r1 $B=T$ $\beta=E'$ rajout de $\text{suivant}(E)$

r2 $\alpha=+$ $B=T$ $\beta=E'$ rajout de $\text{suivant}(E')$

r4 $\alpha=F$ $\beta=T'$ rajout de $\text{suivant}(T)$

r4 $B=F$ $\beta=T'$ rajout de $\text{suivant}(T)$

r5 $\alpha=* F$ $B=T'$

r5 $\alpha=* B=F$ $\beta=T'$ rajout de $\text{suivant}(T')$

Exemple

$S \rightarrow iEtSS' \text{ (r1) } \mid a \text{ (r2)}$

$S' \rightarrow eS \text{ (r3) } \mid \varepsilon \text{ (r4)}$

$E \rightarrow b \text{ (r5)}$

$\text{premier}(S) = \{i, a\}$

$\text{premier}(S') = \{e, \varepsilon\}$

$\text{premier}(E) = \{b\}$

$\text{suivant}(S) = \{\$, e\}$

$\text{suivant}(E) = \{t\}$

$\text{suivant}(S') = \{\$, e\}$

règles $\alpha B \beta$

r1 $B=E$, $\beta=tSS'$ rajout de t

r1 $B=S$, $\beta=S'$ rajout de premier(S')

règles αB ou $\alpha B \beta$ avec $\beta \rightarrow^* \varepsilon$

r1 $\alpha=iEt$ $B=S$ $\beta=S'$

rajout de suivant(S) à suivant(S)

r1 $\alpha=iEtS$ $B=S'$

rajout de suivant(S) à suivant(S')


r3 $\alpha=e$ $B=S$

rajout de suivant(S')

Construction de la table d'analyse

1. Pour chaque règle $r_i \alpha \rightarrow \beta$:
2. Pour chaque $a \in \text{premier}(\beta)$, mettre r_i dans $T[\alpha, a]$
3. Si en plus β est nullable ($\varepsilon \in \text{premier}(\beta)$), pour chaque $a \in \text{suivant}(\alpha)$ ($a \in \Sigma \cup \$$), mettre r_i dans $T[\alpha, a]$

Remarque : nous avons un analyseur prédictif si chaque case la table contient *au plus* une règle de production (pas de conflit)

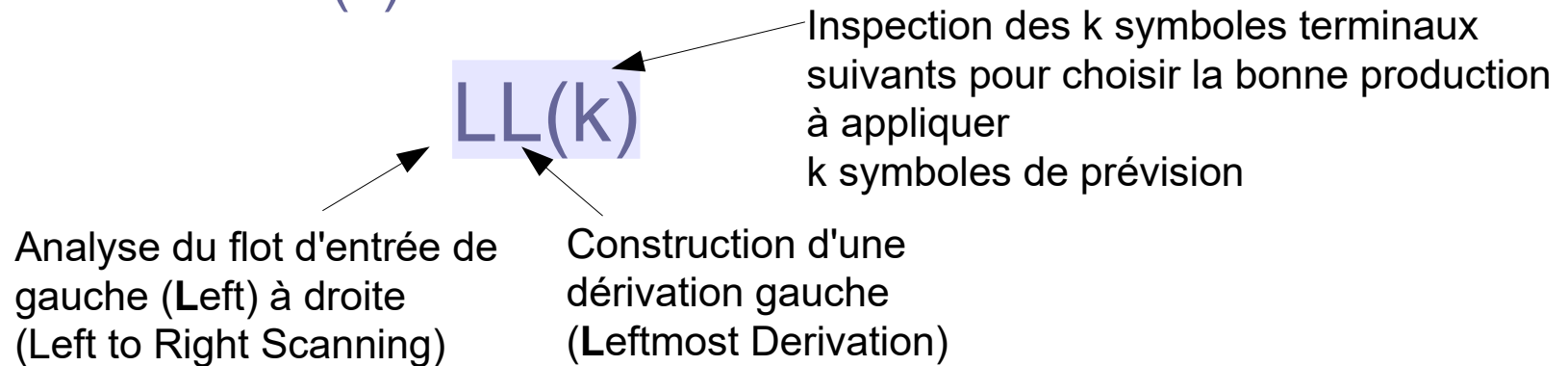


	a	b	e	i	t	\$
S	r2			r1		
S'			r3 r4			r4
E		r5				

ambiguïté de la grammaire

Les grammaires LL(k)

- ▶ Il existe une classe de grammaires adaptée à l'analyse descendante déterministe (\approx sans rebroussement) nommée LL(k)



En pratique, si $k > 1$, la réalisation de l'analyseur correspondant est difficile et peu performante. Si $k = 0$, la grammaire est trop restrictive.

Si $k = 1$, notre analyseur est prédictif

LL(1) \approx aucune case de la table d'analyse n'a plusieurs alternatives

Grammaires prédictives LL(1)

Une grammaire G est LL(1) si et seulement si :

Pour $A \rightarrow \alpha | \beta$

- ▶ il n'existe pas de terminal a tel que :
 $a \in \text{premier}(\alpha)$ et $a \in \text{premier}(\beta)$
- ▶ α et β ne peuvent pas se dériver tous les deux en ϵ
- ▶ Si $\beta \xRightarrow{*} \epsilon$, α ne se dérive pas en une chaîne commençant par un terminal de $\text{suivant}(A)$

On appelle grammaire LL(1) une grammaire pour laquelle la table d'analyse décrite précédemment n'a aucune case définie de façon multiple.

Remarque

- ▶ Une grammaire **ambiguë** ou **récursive à gauche** ou **non factorisée à gauche** n'est pas LL(1)
-

Transformation de grammaires

- ▶ Deux grammaires sont équivalentes si elles engendrent le même langage
 - ▶ Selon la méthode d'analyse (ascendante, descendante...) la grammaire doit vérifier certaines contraintes et par conséquent subir éventuellement quelques transformations
-

Suppression des ambiguïtés

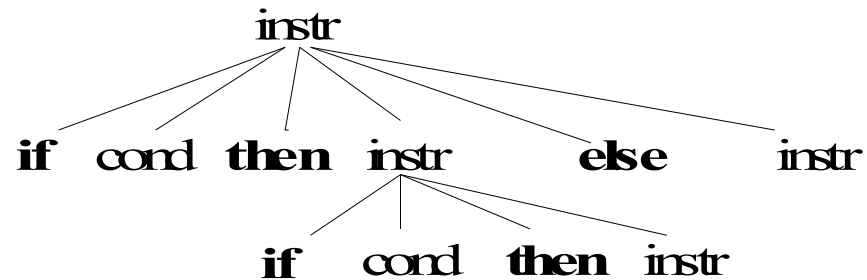
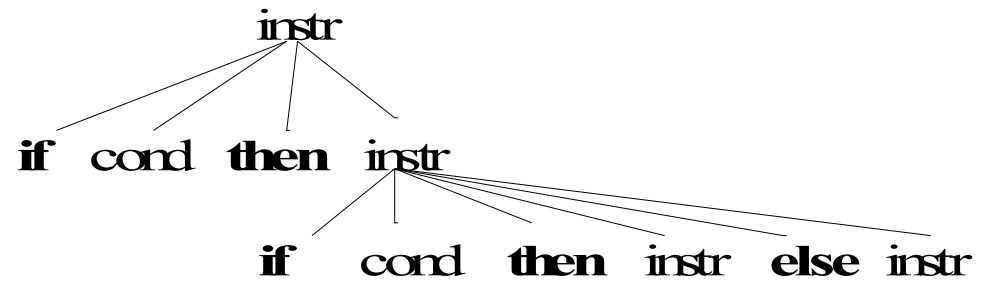
Soit la grammaire suivante :

instr → **si** *expr* **alors** *instr*
| **si** *expr* **alors** *instr* **sinon** *instr*
| *autre*

Cette grammaire est ambiguë car la chaîne :
si E1 alors si E2 alors S1 sinon S2
peut être analysée de deux façons différentes

Exemple

phrase : if cond1 then if E2 then I1 else I2



Désambiguïser une grammaire

Il est parfois possible d'enlever les ambiguïtés d'une grammaire. Dans le cas précédent, il suffit d'ajouter des règles de réécriture permettant d'appliquer le principe : "chaque sinon doit être associé au si le plus près".

instruction → *instr_close* | *instr_non_close*

instr_close → **si** *expr* **alors** *instr_close* **sinon**
instr_close | *autre*

instr_non_close → **si** *expr* **alors** *instr* |
si *expr* **alors** *instr_close* **sinon**
instr_non_close

Exemple en Java

```
if (i>0)
    if (i>5)
        System.out.println("erreur")
else
    System.out.println("négatif")

if (i>0)
{
    if (i>5)
        System.out.println("erreur");
}
else
    System.out.println("négatif")
```

Réversivité à gauche

Une grammaire est réversivité à gauche si elle contient un non terminal A tel qu'il existe une dérivation .

$$A \overset{+}{\Rightarrow} A\alpha$$

Les méthodes d'analyse descendante ne peuvent pas fonctionner avec des grammaires réversives à gauche

$$A \Rightarrow A\alpha \Rightarrow A\alpha\alpha \Rightarrow A\alpha\alpha\alpha \Rightarrow \dots$$

Cas d'une récursivité à gauche immédiate

$$A \rightarrow A\alpha|\beta \quad \longrightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A'|\epsilon \end{array}$$

La récursivité à gauche est transformée en récursivité à droite

Ex 1 :

$$A \rightarrow Aa|b$$

$$A \Rightarrow Aa \Rightarrow Aaa \Rightarrow \dots \Rightarrow b \dots aaa$$

$$\begin{array}{l} \left\{ \begin{array}{l} A \rightarrow bA' \\ A' \rightarrow aA'|\epsilon \end{array} \right. \\ A \Rightarrow bA' \Rightarrow baA' \Rightarrow baaA' \Rightarrow \dots \Rightarrow b \dots aa \end{array}$$

Récurtivité à gauche

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

remplacé par :

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \varepsilon$$

Cas général

$$A \stackrel{+}{\Rightarrow} A \alpha$$

1. Classifier les non-terminaux A_1, A_2, \dots, A_n

2. Pour i de 1 à n

Pour j de 1 à $i - 1$

*Remplacer $A_i \rightarrow A_j \gamma$ par $A_i \rightarrow \delta_1 \gamma | \delta_2 \gamma | \dots | \delta_k \gamma$
où $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$*

Éliminer les récursivités gauches immédiates des A_i

Exemple

$$S \rightarrow Aa|b$$

$$A \rightarrow Ac|Sd|\epsilon$$

- S est récursif à gauche $\underline{S} \Rightarrow Aa \Rightarrow \underline{S}da$
- A est immédiatement récursif à gauche

1) classer $A_1=S$, $A_2=A$

2) pour $A_1=S$

- vide
- S n'est pas immédiatement récursif

pour $A_2=A$

- remplacer $A_1=S$ dans $A_2=A$

$$A \rightarrow Ac \quad \textit{inchangé}$$

$$A \rightarrow Aad|bd$$

$$A \rightarrow \epsilon \quad \textit{inchangé}$$

Éliminer la récursivité gauche
immédiate

$$A \rightarrow bdA'|A'$$

$$A' \rightarrow cA'|adA'|\epsilon$$

On obtient G'

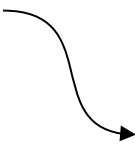
$$S \rightarrow Aa|b$$

$$A \rightarrow bdA'|A'$$

$$A' \rightarrow cA'|adA'|\epsilon$$

Ex 2 : expressions arithmétiques

$$\begin{cases} E \rightarrow E + T \mid T \\ T \rightarrow T * Id \mid Id \end{cases}$$


$$\begin{cases} E \rightarrow TE' \\ E' \rightarrow +TE' \mid \epsilon \\ T \rightarrow IdT' \\ T' \rightarrow *IdT' \mid \epsilon \end{cases}$$

Réversivité à gauche

Il faut enlever les ambiguïtés

$\underline{E} \rightarrow \underline{E} + T \quad | \quad T$ boucle infinie
réversivité à gauche

$A \rightarrow A\alpha \quad | \quad \beta\bar{x}\bar{x}A = E, \alpha = + T$

à remplacer par

$A \rightarrow \beta R$ $A \rightarrow T R$

$R \rightarrow \alpha R \quad | \quad \varepsilon\bar{x}\bar{x}R \rightarrow + T \quad | \quad \varepsilon$

Factorisation à gauche

C'est une transformation utile pour avoir une grammaire permettant l'analyse prédictive. L'idée est que pour développer un non-terminal A , quand il n'est pas évident de choisir l'alternative à utiliser, on ré-écrit les règles de productions

Exemple :

instr \rightarrow **si** *expr* **alors** *instr* **sinon** *instr*
| **si** *expr* **alors** *instr*

← quelle alternative choisir après lecture du si ?

Exemple

instr → **si** *expr* **alors** *instr-suite*

instr-suite → **sinon** *instr* | ε

$$A \rightarrow \alpha \beta_1 | \alpha \beta_2 \quad \longrightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 | \beta_2 \end{array}$$

Remarque : la factorisation ne supprime pas l'ambiguïté

Factorisation à gauche (suite)

Si les règles de production sont de la forme :

$$A \rightarrow \alpha\beta_1|\alpha\beta_2|\dots|\alpha\beta_n|\gamma$$

on transforme en

$$\begin{aligned} A &\rightarrow \alpha A' |\gamma \\ A' &\rightarrow \beta_1|\beta_2|\dots|\beta_n \end{aligned}$$

Problèmes rencontrés

- ▶ On enlève les ambiguïtés
- ▶ On élimine les récursivités à gauche
- ▶ On factorise à gauche

On espère obtenir une grammaire LL(1)

Note : les grammaires ambiguës ou récursives à gauche ne sont pas LL(1)

Si la grammaire est ambiguë, on peut choisir arbitrairement la règle à utiliser en cas de conflit mais on affecte le langage

Implémentation à l'aide de procédures récursives

On considère que chaque symbole non-terminal du langage représente une procédure, et on implémente directement cette procédure.

La règle appelée correspond à l'élément présent dans l'ensemble premier

Grammaire LL(1)

$T \rightarrow R \mid aTc$

$R \rightarrow \varepsilon \mid bR$

$\text{premier}(T) = \{a, b, \varepsilon\}$ $\text{premier}(R) = \{\varepsilon, b\}$

$\text{suivant}(T) = \{\$, c\}$ $\text{suivant}(R) = \{\$, c\}$

	a	b	c	\$
T	r2	r1	r1	r1
R		r4	r3	r3

Analyse de aabbbcc\$

Exemple

```
fonction parseR()  
si next = 'b' alors  
    consomme('b') ; parseR()  
sinon si next = 'c' ou next = '$' alors  
    (* ne rien faire *)  
sinon ErreurSyntaxique()
```

```
fonction parseT()  
si next = 'a' alors  
    consomme('a') ; parseT(); consomme('c')  
sinon si next = 'b' ou next = 'c' ou next = '$' alors  
    parseR()  
sinon ErreurSyntaxique()
```

Conclusions

En transformant la grammaire de façon à ce qu'elle soit :

- ▶ non ambiguë. Il n'y a pas de méthodes. Une grammaire ambiguë est une grammaire qui a été mal conçue
- ▶ factorisée à gauche
- ▶ sans récursivité à gauche

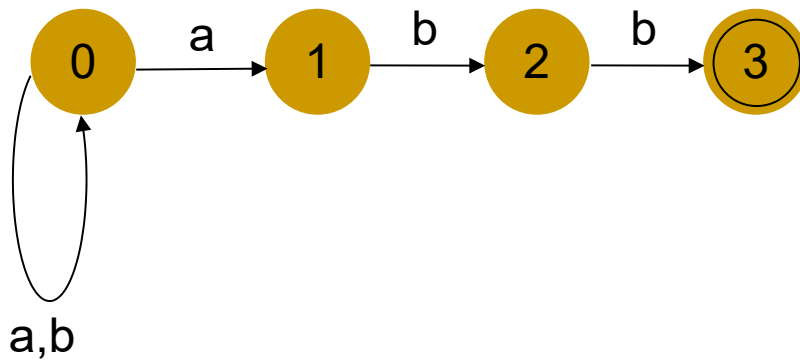
Il ne reste plus qu'à espérer que ça soit LL(1). Sinon, il faut concevoir une autre méthode pour l'analyse syntaxique

On parle d'analyseur prédictif.

LR \subset LNC

A partir d'un AFN, on peut facilement définir une grammaire non contextuelle engendrant le même langage

Ex : $L=(a|b)^*abb$



$$A_0 \rightarrow aA_0 | bA_0 | aA_1$$

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon$$

- Note :
- LR = Langages réguliers
- LNC = Langages non contextuels
- LC = Langages contextuels

LNC $\not\subset$ LR ?

Un AFND ne sait pas compter

$$L = \{ a^n b^n \mid n \geq 1 \}$$
$$G \parallel A \rightarrow aAb \mid ab$$

Impossible de construire un AFND qui « compte » le nombre de a pour avoir autant de b après.

LC ?

Certains langages ne peuvent être engendrés par des grammaires non contextuelles

$$L_1 = \{ w c w \mid w \in (a|b)^* \}$$

Ce type de langage pourrait contrôler que les identificateurs sont déclarés avant d'être utilisés

$$L_2 = \{ a^n b^m c^n d^m \mid n \geq 1, m \geq 1 \}$$

Cette phase est laissée à l'analyse sémantique

Ce type de langage pourrait vérifier que le nombre de paramètres formels correspond au nombre de paramètres effectifs (pour deux procédures)

Remarque : des langages très proches de L_1 et L_2 peuvent être engendrés par une grammaire non contextuelle

LC ? (suite)

$$L'_1 = \{ w o w^R \mid w \in (a|b)^* \}$$

$$A \rightarrow aAa \mid bAb \mid c$$

$$L'_2 = \{ a^n b^m c^m d^n \mid n \geq 1 \ m \geq 1 \}$$

$$\begin{cases} A \rightarrow aAd \mid aA'd \\ A' \rightarrow bA'c \mid bc \end{cases}$$

Analyse ascendante

- Analyse par décalage réduction
-

Principes

- ▶ Tente de construire un arbre d'analyse par une chaîne d'entrée des feuilles (bas) vers la racine (haut)
- ▶ Tente de déterminer une dérivation droite en « sens inverse »

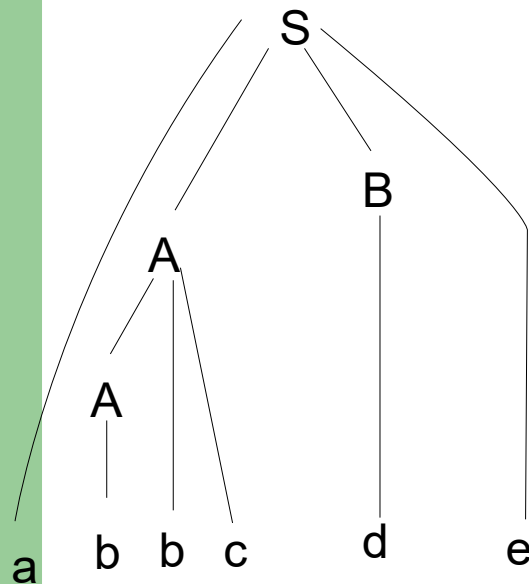
Il s'agit de « réduire » une chaîne w vers l'axiome de la grammaire. On regarde si dans la chaîne à analyser, il existe une sous-chaîne correspondant à la partie droite d'une production et on remplace par le non-terminal de la partie gauche de cette production

Remarques

- ▶ Malgré les apparences, les analyseurs ascendants sont plus efficaces que les analyseurs descendants. En outre, ils acceptent une classe de langage plus large
 - ▶ Le principal inconvénient de ces analyseurs est qu'ils nécessitent des tables qu'il est extrêmement difficile de construire à la main
-

Exemple

$$\left\{ \begin{array}{l} S \rightarrow aABe \\ A \rightarrow Abc|b \\ B \rightarrow d \end{array} \right. \quad \text{et } w = abcde$$



$$S \Rightarrow \underline{aABe} \Rightarrow aA \underline{d} e \Rightarrow a \underline{Abc} de \Rightarrow a \underline{b} bcde$$

On recherche les sous-chaînes correspondants à la partie droite d'une règle de production

On construit ainsi un arbre d'analyse des feuilles vers la racine

Exemple (suite)

La suite des réductions lue à l'envers représente une dérivation droite de l'axiome

$$S \Rightarrow aA \underline{B} e \Rightarrow a \underline{A} de \Rightarrow a \underline{A} bc de \Rightarrow abbcde$$

Remarque : le choix de la sous-chaîne à réduire ne doit pas être quelconque. Dans notre exemple :

$$aA \underline{A} cde \Rightarrow aA \underline{b} cde \Rightarrow a \underline{b} bcde$$

bloqué

Définition

On appelle manche de chaîne, une sous-chaîne correspondant à la partie droite d'une production et dont la réduction vers le non-terminal de la partie gauche représente une étape le long de la dérivation droite inverse.

$$\left| \begin{array}{l} \text{si } S \xrightarrow{d}^* \alpha A w \xrightarrow{d} \alpha \gamma w \\ \text{avec } \left\{ \begin{array}{l} w \in \Sigma^* \\ A \rightarrow \gamma \in P \end{array} \right. \end{array} \right.$$

alors γ est un manche de $\alpha\gamma w$

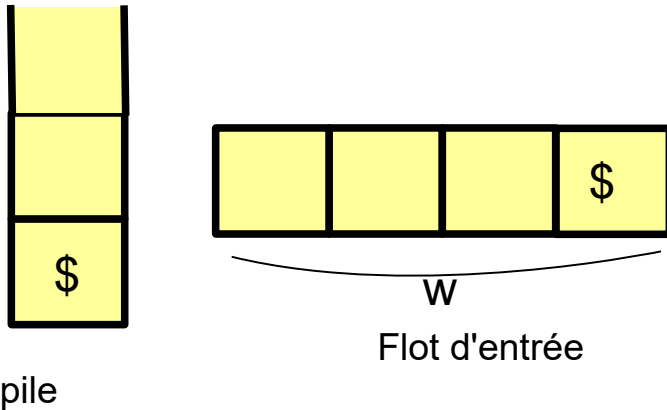
Si la grammaire est ambiguë, il peut exister plusieurs manches dans une chaîne



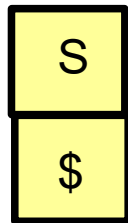
Mise en place de l'analyse par décalage - réduction

Implémentation à l'aide d'une pile

au départ :



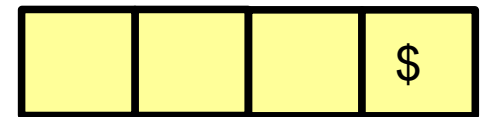
On recommence ces « décalages-réductions » jusqu'à ce que l'on détecte une erreur



en cours :

On décale de l'entrée vers la pile 0, 1 ou plusieurs symboles jusqu'à ce qu'un manche γ se trouve en sommet de pile. On le remplace par la partie gauche de la production.

ou un succès :



Exemple

$$\begin{cases} E \rightarrow E + E \\ E \rightarrow E * E \\ E \rightarrow (E) \\ E \rightarrow id \end{cases}$$

$$w = id + id * id$$

Cas 1

			<i>id</i>	<i>]</i>	<i>E</i>	
		*	*	*		
<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	succès
\$	\$	\$	\$	\$	\$	

			<i>id</i>	<i>]</i>	<i>E</i>
		+	+	+	
<i>id</i>	<i>]</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>
\$	\$	\$	\$	\$	\$

Cas 2

			<i>id</i>	<i>]</i>	<i>E</i>	
		*	*	*		
<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	
+	+	+	+	+		
<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	<i>E</i>	
\$	\$	\$	\$	\$	\$	succès

Note : comme la grammaire est ambiguë,
il existe 2 séquences possibles pour l'analyseur

Conflits possibles au cours de l'analyse

- ▶ Conflit décaler/réduire : en cours d'analyse, on ne peut pas décider s'il faut réduire le manche présent en sommet de pile ou décaler encore quelques symboles pour faire apparaître un autre manche
(exemple grammaire précédente)
 - ▶ Conflit réduire/réduire : le manche au sommet de pile peut être réduit par plusieurs productions
(exemple : grammaire où l'appel des procédures avec paramètres à la même syntaxe que les tableaux ou leur indice)
-

Exemple

```
Proc → id(liste_paramètres)
liste_paramètres → liste_paramètres,paramètre
paramètre → id
tab → id(liste_exprs)
liste_exprs → liste_exprs,expr
expr → id
```

- A l'analyse de l'entrée $A(i,j)$
- le flot d'unités lexicales est $ID(ID,ID)$
- Au cours de l'analyse
 - ID représente-t-il un indice de tableau ou un paramètre effectif ?

Pour éviter ces conflits, nous allons définir un classe de grammaire $LR(k)$ pour lesquelles l'analyseur décalage/réduction se déroule correctement

Exemple

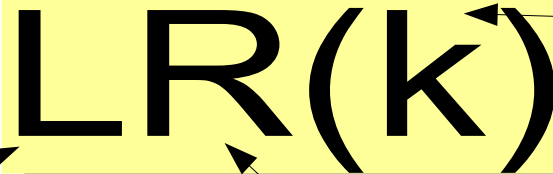
```
Proc → id(liste_paramètres)
liste_paramètres → liste_paramètres,paramètre
paramètre → id
tab → id(liste_exprs)
liste_exprs → liste_exprs,expr
expr → id
```

- A l'analyse de l'entrée $A(i,j)$
- le flot d'unités lexicales est $ID(ID,ID)$
- Au cours de l'analyse
 - ID représente-t-il un indice de tableau ou un paramètre effectif ?

Pour éviter ces conflits, nous allons définir un classe de grammaire $LR(k)$ pour lesquelles l'analyseur décalage/réduction se déroule correctement

Les analyseurs LR

Analyseurs par décalage réduction sans rebroussement

The diagram shows the notation LR(k) in a yellow box. An arrow points from the 'L' to the text 'Analyse du flot d'entrée de Gauche à Droite (Left to Right scanning)'. Another arrow points from the 'R' to the text 'Construction d'une dérivation Droite inverse (Rightmost derivation in reverse)'. A third arrow points from the '(k)' to the text 'Indique le nombre de symboles de pré-vision pour prendre les décisions d'analyse'.

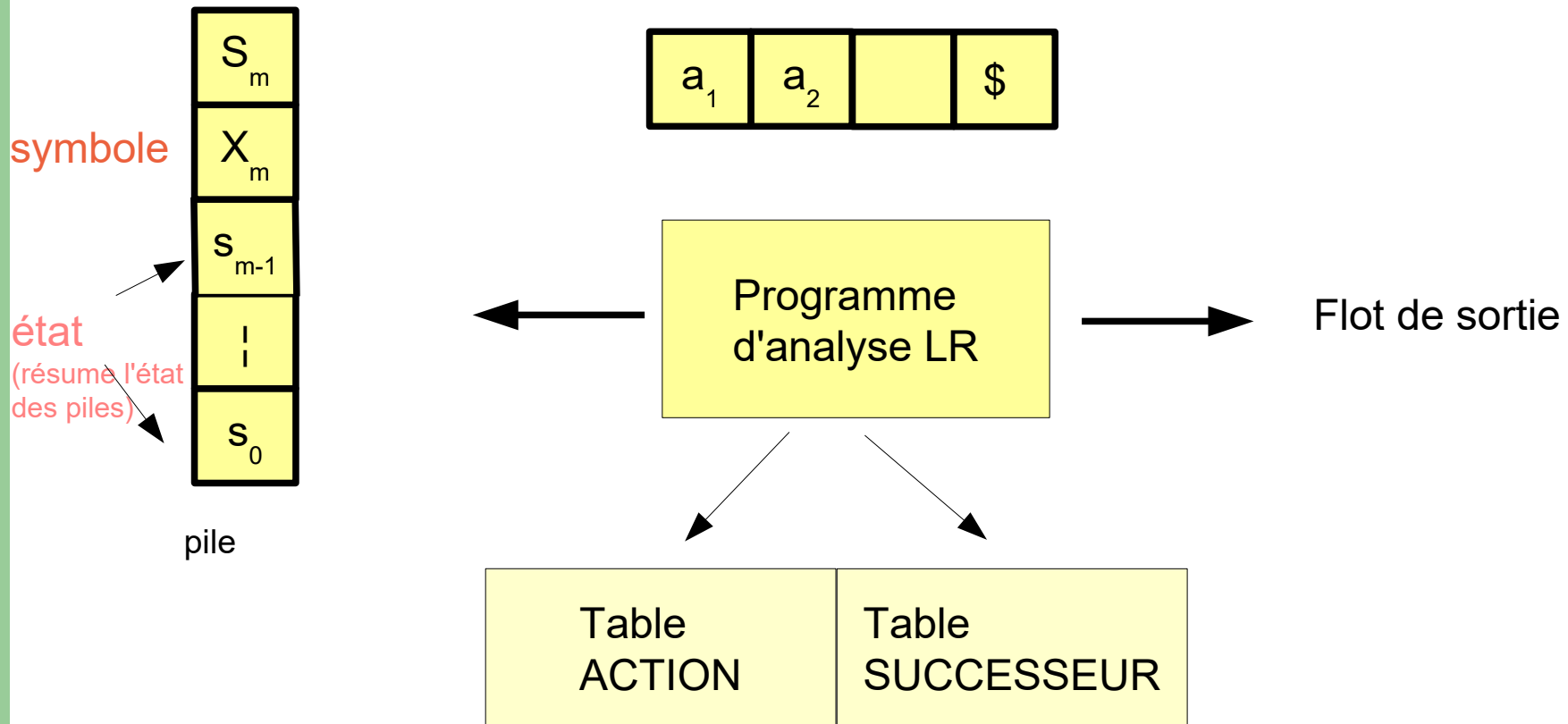
Indique le nombre de symboles de pré-vision pour prendre les décisions d'analyse

Analyse du flot d'entrée de Gauche à Droite
(Left to Right scanning)

Construction d'une dérivation
Droite inverse
(Rightmost derivation in reverse)

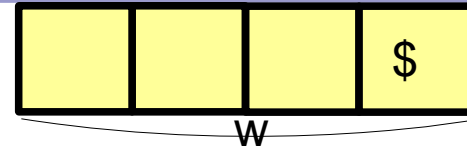
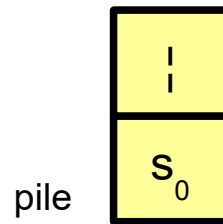
Note : quand k est omis, il est supposé être égal à 1

Algorithme d'analyse LR



Principe

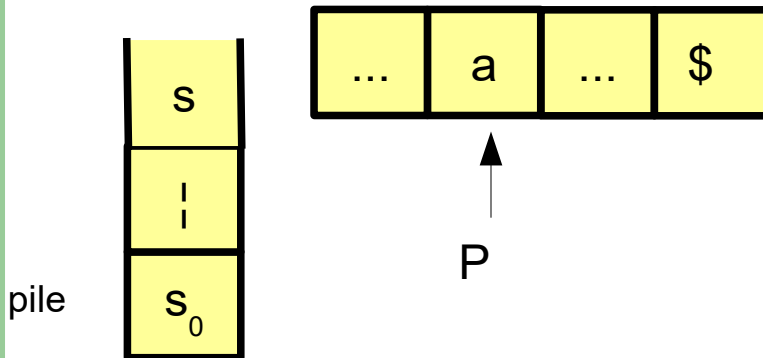
Au départ :



Si $Action[s,a] = \text{Décaler } s'$

[
empiler(a)
empiler(s')
avancer(p)

En cours :



si $Action[s,a] = \text{Réduire } A \rightarrow \beta$

[
Dépiler $2 * \beta$ symboles
Soit s' l'état au sommet
empiler(A)
empiler($Succ[s', A]$)

si $Action[s,a] = \text{Accepter}$
 \Rightarrow Succès

si $Action[s,a] = \text{Erreur}$
 \Rightarrow Routine de récupération
 sur erreur

Exemple : grammaire des expressions arithmétiques

notation :

- ▶ d_i décaler et empiler l'état i
- ▶ r_j réduire avec la règle j
- ▶ acc accepter => succès
- ▶ 0 erreur

Tables d'analyse

Etat	ACTION						SUCC		
	Id	+	*	()	\$	E	T	F
0	d5			d4			1	2	3
1		d6				acc			
2		r2	d7		r2	r2			
3		r4	r4		r4	r4			
4	d5			d4			8	2	3
5		r6	r6		r6	r6			
6	d5			d4				9	3
7	d5			d4					10
8		d6			d11				
9		r1	d7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			

$$E \rightarrow E + T \quad (1)$$

$$E \rightarrow T \quad (2)$$

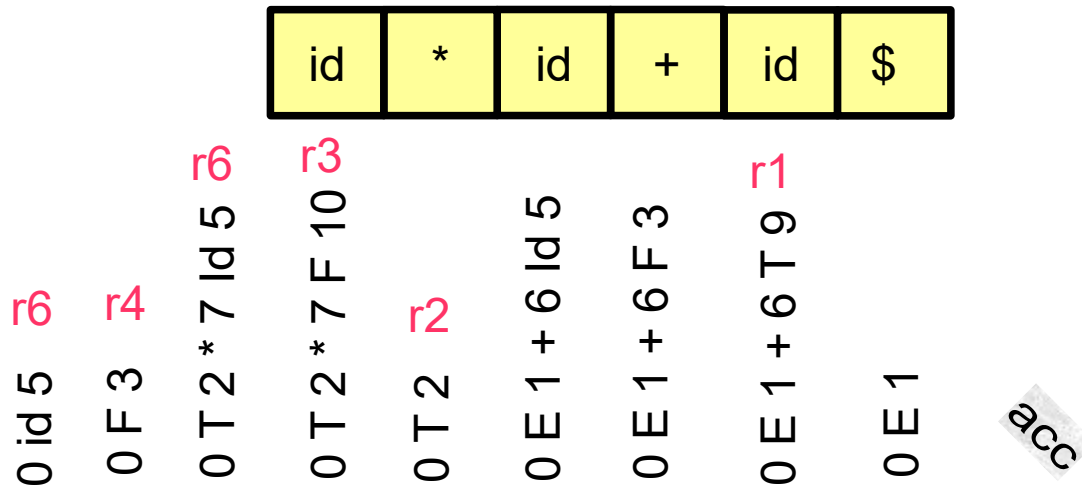
$$T \rightarrow T * F \quad (3)$$

$$T \rightarrow F \quad (4)$$

$$F \rightarrow (E) \quad (5)$$

$$F \rightarrow id \quad (6)$$

Exemple (suite)



$$\begin{aligned}
 E &\Rightarrow E + T' \Rightarrow E + F \Rightarrow E + Id \Rightarrow T + Id \\
 T + Id &\Rightarrow T * F + Id \Rightarrow T * Id * Id \Rightarrow F * Id + Id \Rightarrow Id * Id + Id
 \end{aligned}$$

Conclusion sur les grammaires LR

- ▶ Une grammaire pour laquelle nous pouvons construire les tables Action et Successeur d'analyse LR(k) est une grammaire LR(k)
Il existe des grammaires non contextuelles qui ne sont pas LR, mais tous les langages de programmation classique sont reconnus par des grammaires LR
 - ▶ Les analyseurs LR sont plus puissants que les analyseurs prédictifs, mais s'il est possible de « programmer à la main » un analyseur prédictif, nous devons utiliser un constructeur d'analyseurs pour les analyseurs LR (par exemple yacc ou bison) car la quantité de travail est trop importante
-

Construction des tables d'analyse

Il existe 3 méthodes pour construire ces tables d'analyse :

1. **SLR** (simple LR) : la plus facile mais la moins puissante
 2. **LR** : la plus puissante et la plus coûteuse
 3. **LALR** (Look Ahead LR) : puissante et coût intermédiaire entre les 2 précédentes
-

Définition

- ▶ Un item est une production avec un point repérant une position dans la partie droite

Soit A une règle de production avec $X, Y, Z \in (\Sigma \cup V)$

4 items possibles

$$\left| \begin{array}{l} A \rightarrow .XYZ \\ A \rightarrow X.YZ \\ A \rightarrow XY.Z \\ A \rightarrow XYZ. \end{array} \right.$$

interprétation : $A \rightarrow X.YZ$

- ▶ Nous venons de lire une chaîne dérivée de X
 - ▶ Nous attendons une chaîne dérivée de YZ
-

Construction des tables SLR

L'idée est de construire un AFD pour reconnaître les préfixes possibles pour une analyse par décalage – réduction.

Méthode :

1. on augmente la grammaire de la production $S' \rightarrow .S$ et S' devient le nouvel axiome (à la place de S)
2. pour construire l'automate, on définit 2 fonctions
 $Fermeture(I)$ I étant un ensemble d'items

- $I \subset Fermeture(I)$
- si $A \rightarrow \alpha . B \beta \in Fermeture(I)$ $B \rightarrow . \gamma \in Fermeture(I)$
 et $B \rightarrow \gamma$ est une production

$Transition(I, X) = \Delta(I, X)$ $X \in (T \cup V)$

est la fermeture de l'ensemble des items de la forme
 $(A \rightarrow \alpha X \beta)$ tel que $(A \rightarrow \alpha . X \beta) \in I$

Exemple de fermeture

► Soit la grammaire :

$$G \left\{ \begin{array}{l} E \rightarrow E + T \quad (1) \\ E \rightarrow T \quad (2) \\ T \rightarrow T * F \quad (3) \\ T \rightarrow F \quad (4) \\ F \rightarrow (E) \quad (5) \\ F \rightarrow Id \quad (6) \end{array} \right.$$

► Et soit l'ensemble d'items

$$I = \{ T \rightarrow T * . F, E \rightarrow E . + T \}$$

La fermeture de cet ensemble d'items est :

$$\{ T \rightarrow T * . F, E \rightarrow E . + T, F \rightarrow . Id, F \rightarrow . (E) \}$$

Exemple de transition

► Sur l'exemple de la
grammaire

$$G \left\{ \begin{array}{l} E \rightarrow E + T \quad (1) \\ E \rightarrow T \quad (2) \\ T \rightarrow T * F \quad (3) \\ T \rightarrow F \quad (4) \\ F \rightarrow (E) \quad (5) \\ F \rightarrow Id \quad (6) \end{array} \right.$$

► On aura :

$$\Delta(I, F) = \{T \rightarrow T * F.\}$$

$$\Delta(I, +) = \{E \rightarrow E + .T, T \rightarrow .T * F, T \rightarrow .F, F \rightarrow .Id, F \rightarrow .(E)\}$$

Soit l'ensemble d'items :

$$\left\{ \begin{array}{l} T \rightarrow T * .F, E \rightarrow E . + T, \\ F \rightarrow .Id, F \rightarrow .(E) \end{array} \right\}$$

Construction des tables SLR (suite)

On part de l'état $I_0 = Fermeture(\{S' \rightarrow \cdot S\})$

Pour chaque symbole X de la grammaire, on effectue la transition (I_0, X) en créant ainsi de nouveaux états I_1, I_2, \dots

On recommence ainsi tant que tous les états n'ont pas été étudiés.

Application :

$$G \left\{ \begin{array}{l} E \rightarrow E + T \quad (1) \\ E \rightarrow T \quad (2) \\ T \rightarrow T * F \quad (3) \\ T \rightarrow F \quad (4) \\ F \rightarrow (E) \quad (5) \\ F \rightarrow Id \quad (6) \end{array} \right.$$

Construction des tables SLR (suite)

$$I_0: \left| \begin{array}{l} S' \rightarrow . E \\ E \rightarrow . E + T \\ E \rightarrow . T \\ T \rightarrow . T * F \\ T \rightarrow . F \\ F \rightarrow . (E) \\ F \rightarrow . id \end{array} \right.$$

$$I_4: \left| \begin{array}{l} F \rightarrow (. E) \\ E \rightarrow . E + T \\ E \rightarrow . T \\ T \rightarrow . T * F \\ T \rightarrow . F \\ F \rightarrow . (E) \\ F \rightarrow . id \end{array} \right.$$

$$I_7: \left| \begin{array}{l} T \rightarrow T * . F \\ F \rightarrow . (E) \\ F \rightarrow . id \end{array} \right.$$

$$I_8: \left| \begin{array}{l} F \rightarrow (E .) \\ E \rightarrow E . + T \end{array} \right.$$

$$I_1: \left| \begin{array}{l} S' \rightarrow E . \\ E \rightarrow E . + T \end{array} \right.$$

$$I_5: F \rightarrow id .$$

$$I_9: \left| \begin{array}{l} E \rightarrow E + T . \\ T \rightarrow T . * F \end{array} \right.$$

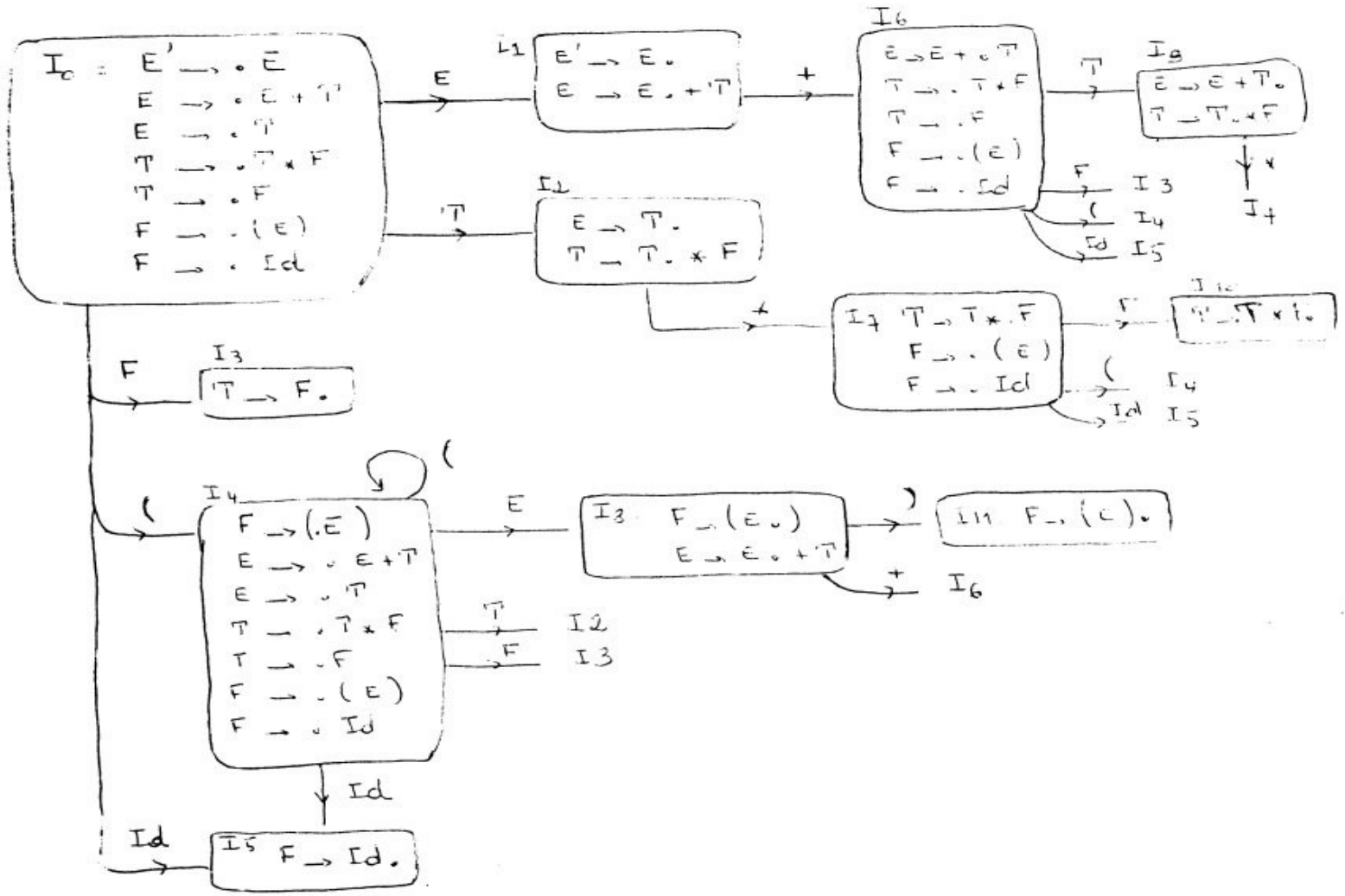
$$I_2: \left| \begin{array}{l} E \rightarrow T . \\ T \rightarrow T . * F \end{array} \right.$$

$$I_6: \left| \begin{array}{l} E \rightarrow E + . T \\ T \rightarrow . T * F \\ T \rightarrow . F \\ F \rightarrow . (E) \\ F \rightarrow . id \end{array} \right.$$

$$I_{10}: T \rightarrow T * F .$$

$$I_3: T \rightarrow F .$$

$$I_{11}: T \rightarrow (E) .$$



Construction de la table d'analyse

1. Construire la collection d'items $\{I_0, I_1, \dots, I_n\}$
 2. L'état i est construit à partir de I_i :
 1. Pour chaque $\Delta(I_i, a) = I_j$: mettre décaler j dans $M[i, a]$
 2. Pour chaque $\Delta(I_i, A) = I_j$: mettre aller en j dans $M[i, A]$
 3. Pour chaque $A \rightarrow \alpha$. (sauf $A = S'$) contenu dans I_i :
mettre réduire $A \rightarrow \alpha$ dans chaque $M[i, a]$ où
 $a \in \sigma\upsilon\iota\omega\alpha\nu\tau(A)$
 4. Si $S' \rightarrow S. \in I_i$: mettre accepter dans $M[i, \$]$
-

Construction des tables SLR (suite)

3. Les tables ACTION et SUCCESSEUR

$si A \rightarrow \alpha . a \beta \in I_i$
 $et transition(I_i, a) = I_j$
 $avec a \in T$

$$\left. \vphantom{\begin{array}{l} si A \rightarrow \alpha . a \beta \in I_i \\ et transition(I_i, a) = I_j \\ avec a \in T \end{array}} \right\} \underline{Action}[i, a] = \text{décaler } j$$

$si A \rightarrow \alpha . \in I_i$
 $et A \neq S'$

$$\left. \vphantom{\begin{array}{l} si A \rightarrow \alpha . \in I_i \\ et A \neq S' \end{array}} \right\} \underline{Action}[i, a] = \text{réduire par } A \rightarrow \alpha$$

pour $A \in \text{suivant}(A)$

$si S' \rightarrow S . \in I_i$

$$\underline{Action}[i, \$] = \text{Accepter}$$

Toutes les cases d'ACTION non définies sont à « Erreur »

$\forall A \in V$ si $Transition(I_i, A) = I_j$ alors $\underline{Successeur}[i, a] = j$

	ACTION						SUCCESSEUR		
	ld	+	*	()	\$	E	T	F
0	d5			d4			1	2	3
1		d6							
2		$E \rightarrow T$	d7		$E \rightarrow T$	$E \rightarrow T$			
3		$T \rightarrow F$	$T \rightarrow F$		$T \rightarrow F$	$T \rightarrow F$			
4	d5			d4			8	2	3
5		$F \rightarrow ld$	$F \rightarrow ld$		$F \rightarrow ld$	$F \rightarrow ld$			
6	d5			d4				9	3
7	d5			d4					10
8		d6			d11				
9		$E \rightarrow E+T$	d7		$E \rightarrow E+T$	$E \rightarrow E+T$			
10		$T \rightarrow T^*F$	$T \rightarrow T^*F$		$T \rightarrow T^*F$	$T \rightarrow T^*F$			
11		$F \rightarrow (E)$	$F \rightarrow (E)$		$F \rightarrow (E)$	$F \rightarrow (E)$			

avec $\text{suivant}(E) = \{+,), \$\}$

$\text{suivant}(T) = \{*, +,), \$\}$

$\text{suivant}(F) = \{*, +,), \$\}$

Remarques

- ▶ Des conflits peuvent apparaître dans la table d'analyse (grammaire non SLR dans ce cas)
 - ▶ Pour certaines grammaires, l'analyse SLR n'est pas assez puissante pour se rappeler suffisamment le contexte gauche pour décider de l'action à faire. Pour un item $A \rightarrow \alpha.$, on place cette production dans la table ACTION pour chaque terminal a de $\text{suivant}(A)$. Or, il se peut que le manche situé sur la pile à cet état de l'automate ne puisse jamais être suivi de a . L'analyse LR résout ce problème. Elle précise après chaque item, l'ensemble des terminaux qui peuvent effectivement suivre le manche en sommet de pile.
-

Construction des tables LR

La méthode est similaire à la construction des tables SLR

Fermeture(I)

$$\begin{array}{l}
 I \subset \text{Fermeture}(I) \\
 \left. \begin{array}{l}
 \text{si } [A \rightarrow \alpha . B \beta, a] \in I \\
 \quad B \rightarrow \gamma \\
 \quad b \in \text{Premier}(\beta a)
 \end{array} \right\} \Rightarrow [B \rightarrow . \gamma, b] \in \text{Fermeture}(I)
 \end{array}$$

Transition(I,X) est la fermeture des Items de la forme

$$[A \rightarrow \alpha X . \beta, a] \text{ tel que } [A \rightarrow \alpha . X \beta, a] \in I$$

Pour l'automate, on part de $I_0 = \text{Fermeture}([S' \rightarrow . S, \$])$

Exemple

$$G \left| \begin{array}{l} S \rightarrow CC \\ C \rightarrow cC \mid d \end{array} \right.$$

$$\text{first}(S) = \{c, d\}$$

$$\text{first}(C) = \{c, d\}$$

Fermeture de $([S' \rightarrow .S, \$])$

identification avec $[A \rightarrow \alpha . B \beta, a]$

$A = S' \quad \alpha = \epsilon \quad B = S \quad \beta = \epsilon \quad \text{et} \quad a = \$$

Ajout de $[B \rightarrow . \gamma, b]$

$\beta \rightarrow \gamma$ est $S \rightarrow CC$

on ajoute donc $[S \rightarrow . CC, \$]$

$$I_0 \left\{ \begin{array}{l} S' \rightarrow . S, \$ \\ S \rightarrow . CC, \$ \\ C \rightarrow . cC, c|d \\ C \rightarrow . d, c|d \end{array} \right.$$

identification de $[S \rightarrow . CC, \$]$ avec $[A \rightarrow \alpha . B \beta, a]$

$A = S \quad \alpha = \epsilon \quad B = C \quad \text{et} \quad a = \$$

puisque C ne dérive pas en la chaîne vide, $\text{first}(C \$) = \text{first}(C)$

on ajoute donc $[C \rightarrow . cC, c], [C \rightarrow . cC, d], [C \rightarrow . d, d]$

Calcul de transition(I_0, X)

pour $X = S$ nous fermons $\{[S' \rightarrow S., \$]\}$

$I_1: S' \rightarrow S., \$$

pour $X = C$ nous fermons $\{[C \rightarrow C.C, \$]\}$

$$I_2 \left| \begin{array}{l} S \rightarrow C.C, \$ \\ C \rightarrow .cC, \$ \\ C \rightarrow .d, \$ \end{array} \right.$$

pour $X = c$ nous fermons $\{[C \rightarrow c.C, c|d]\}$

$$I_3 \left| \begin{array}{l} C \rightarrow c.C, c|d \\ C \rightarrow .cC, c|d \\ C \rightarrow .d, c|d \end{array} \right.$$

pour $X = d$

$I_4: C \rightarrow d., c|d$

Exemple (suite)

- ▶ Aucun nouvel ensemble sur I_1
- ▶ I_2 a des transitions sur C, c et d

$$I_5: S \rightarrow CC., \$$$

- ▶ Sur c nous fermons $\{[C \rightarrow c.C, \$]\}$

$$I_6 \left| \begin{array}{l} C \rightarrow c.C, \$ \\ C \rightarrow .cC, \$ \\ C \rightarrow .d, \$ \end{array} \right.$$

Exemple (suite)

- ▶ Transition (I_2, d)

$$I_7: C \rightarrow d., \$$$

- ▶ Transition $(I_3, c) = I_3$, Transition $(I_3, d) = I_4$, Transition $(I_3, C) = I_8$

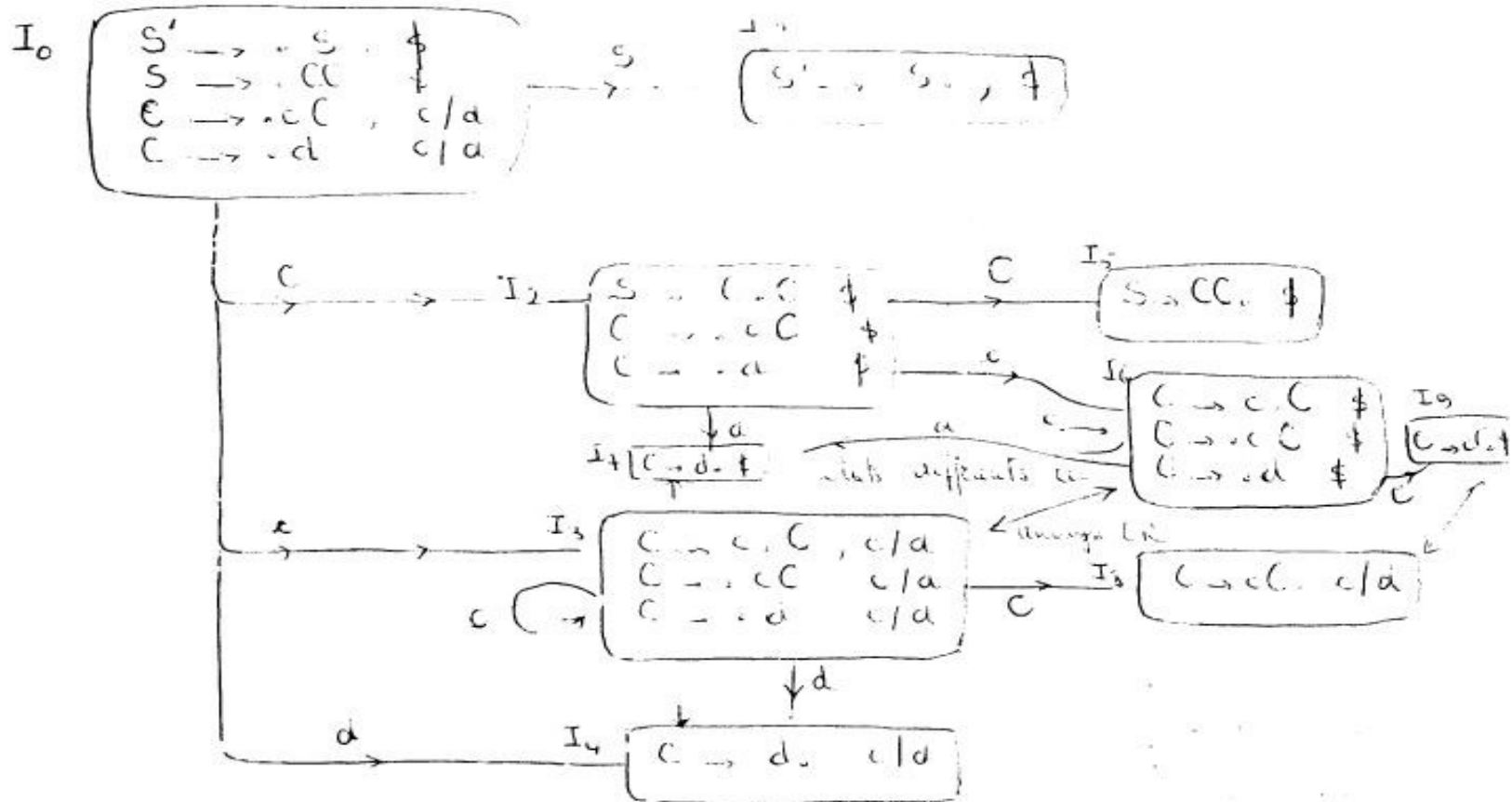
$$I_8: C \rightarrow cC., c|d$$

- ▶ Transition pour I_6 sur c et d sont I_6 et I_7 et

Transition (I_6, C)

$$I_9: C \rightarrow cC., \$$$

Graphe de la fonction Transition



Remplissage de la table ACTION

$si [A \rightarrow \alpha., a] \in I_i, A \neq S' \} Action[i, a] = \text{réduire } A \rightarrow \alpha$

État	ACTION		SUCCESSEUR		
	c	d	\$	S	C
0	d3	d4		1	2
1			Acc		
2	d6	d7			5
3	d3	d4			8
4	$C \rightarrow d$	$C \rightarrow d$			
5			$S \rightarrow CC$		
6	d6	d7			9
7			$C \rightarrow d$		
8	$C \rightarrow cC$	$C \rightarrow cC$			
9			$C \rightarrow cC$		

Remarques

- ▶ Des conflits peuvent encore apparaître ! Même pour des grammaires non ambiguës. Dans ce cas, la grammaire n'est pas LR (rare pour les langages)
 - ▶ Toute grammaire SLR est LR mais le nombre d'états de l'analyseur est supérieur au nombre d'états de l'analyseur SLR (la grammaire de l'exemple est SLR et ne nécessite que 7 états dans un analyseur SLR)
 - ▶ Dans la pratique, un langage comme PASCAL engendre plusieurs centaines d'états en SLR et plusieurs milliers en LR !
 - ▶ L'analyse LALR évite certains conflits sans engendrer autant d'états que l'analyse (Quelques centaines quand même)
-

Construction des tables LALR

- ▶ LALR : Look Ahead LR
 - ▶ On construit l'automate d'analyse LR et on regroupe les états ayant un « cœur » commun (les mêmes items indépendamment des terminaux de pré-vision)
 - ▶ Si aucun conflit ne se produit, la grammaire est LALR
-

Exemple

Dans l'exemple précédent, on fusionne :

- ▶ I_4 et I_7 en I_{47}
- ▶ I_3 et I_6 en I_{36}
- ▶ I_8 et I_9 en I_{89}

	c	d	\$	S	C
0	d36	d47		1	2
1			Acc		
2	d36	d47			5
36	d36	d47			89
47	$C \rightarrow d$	$C \rightarrow d$	$C \rightarrow d$		
5			$S \rightarrow CC$		
89	$C \rightarrow cC$	$C \rightarrow cC$	$C \rightarrow cC$		

Cette grammaire est LALR

Il existe un algorithme pour construire efficacement les tables d'analyse LALR (sans passer par les tables LR)

Un générateur d'analyseur syntaxique

- ▶ `ply.yacc` module AS de `ply`
- ▶ `import ply.yacc as yacc`
- ▶ Il est nécessaire de disposer d'une grammaire au format BNF

Affectation ::= NOM EGAL expr

Expr ::= expr PLUS terme
| expr MINUS terme
| terme

terme ::= terme FOIS facteur
| terme DIVISE facteur
| facteur

facteur ::= NOMBRE

Exemple

```
import ply.yacc as yacc
```

```
import mylexer # on importe le lexer
```

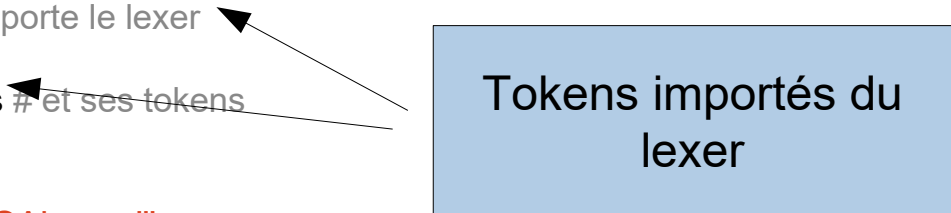
```
tokens = mylexer.tokens # et ses tokens
```

```
def p_affectation(p):  
    "affectation : NOM EGAL expr"
```

```
def p_expr(p):  
    "expr : expr PLUS terme  
    | expr MOINS terme  
    | terme"
```

```
def p_terme(p):  
    "terme : terme FOIS facteur  
    | terme DIVISE facteur  
    | facteur"
```

```
def p_facteur(p):  
    "facteur : NOMBRE"
```



Tokens importés du
lexer

Exemple

```
import ply.yacc as yacc

import mylexer # on importe le lexer

tokens = mylexer.tokens # et ses tokens

def p_affectation(p):
    """affectation : NOM EGAL expr"""

def p_expr(p):
    """expr : expr PLUS terme
    | expr MOINS terme
    | terme """

def p_terme(p):
    """terme : terme FOIS facteur
    | terme DIVISE facteur
    | facteur"""

def p_facteur(p):
    """facteur : NOMBRE"""
```

Les règles de
grammaire
correspondent à des
fonctions

Seul critère : les noms
doivent commencer par
p_

Exemple

```
import ply.yacc as yacc

import mylexer # on importe le lexer

tokens = mylexer.tokens # et ses tokens


def p_affectation(p):
    "affectation : NOM EGAL expr"

def p_expr(p):
    "expr : expr PLUS terme
    | expr MOINS terme
    | terme"

def p_terme(p):
    "terme : terme FOIS facteur
    | terme DIVISE facteur
    | facteur"

def p_facteur(p):
    "facteur : NOMBRE"
```

Ce sont les docstrings
qui contiennent les
règles au format BNF



Construction du parser

- ▶ Parser = `yacc.yacc()`
 - ▶ Construction par introspection
-

Lancement de l'analyse syntaxique

```
parser = yacc.yacc()
```

```
data = 'a = 5 +3*2 '
```

```
resultat = parser.parse(data)
```

```
print(resultat)
```

- Analyse par le parser par invocation des règles
-

Fonctionnement

- ▶ PLY utilise l'analyse LR (comme on a vu dans le cours)
 - ▶ Plus exactement LALR(1)
 - ▶ C'est la technique la plus utilisée pour l'analyse syntaxique des langages
 - ▶ Décalage – réduction
 - ▶ On remplace la partie droite par la partie gauche
 - ▶ Analyse ascendante
-