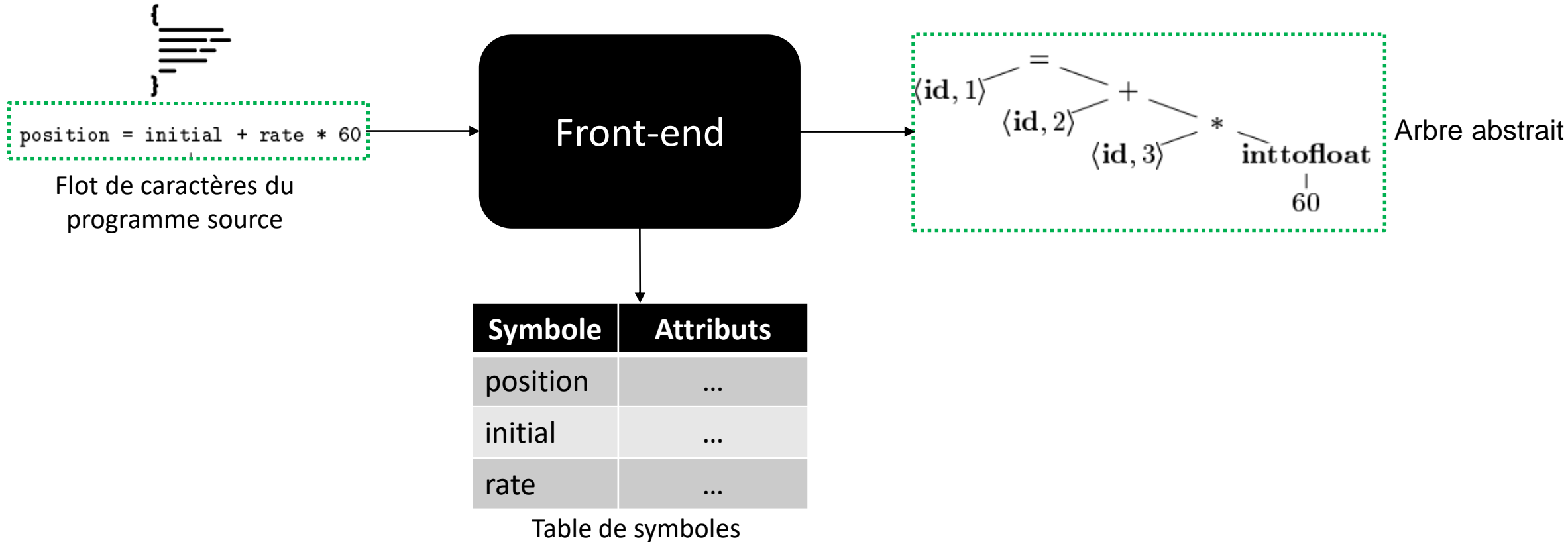


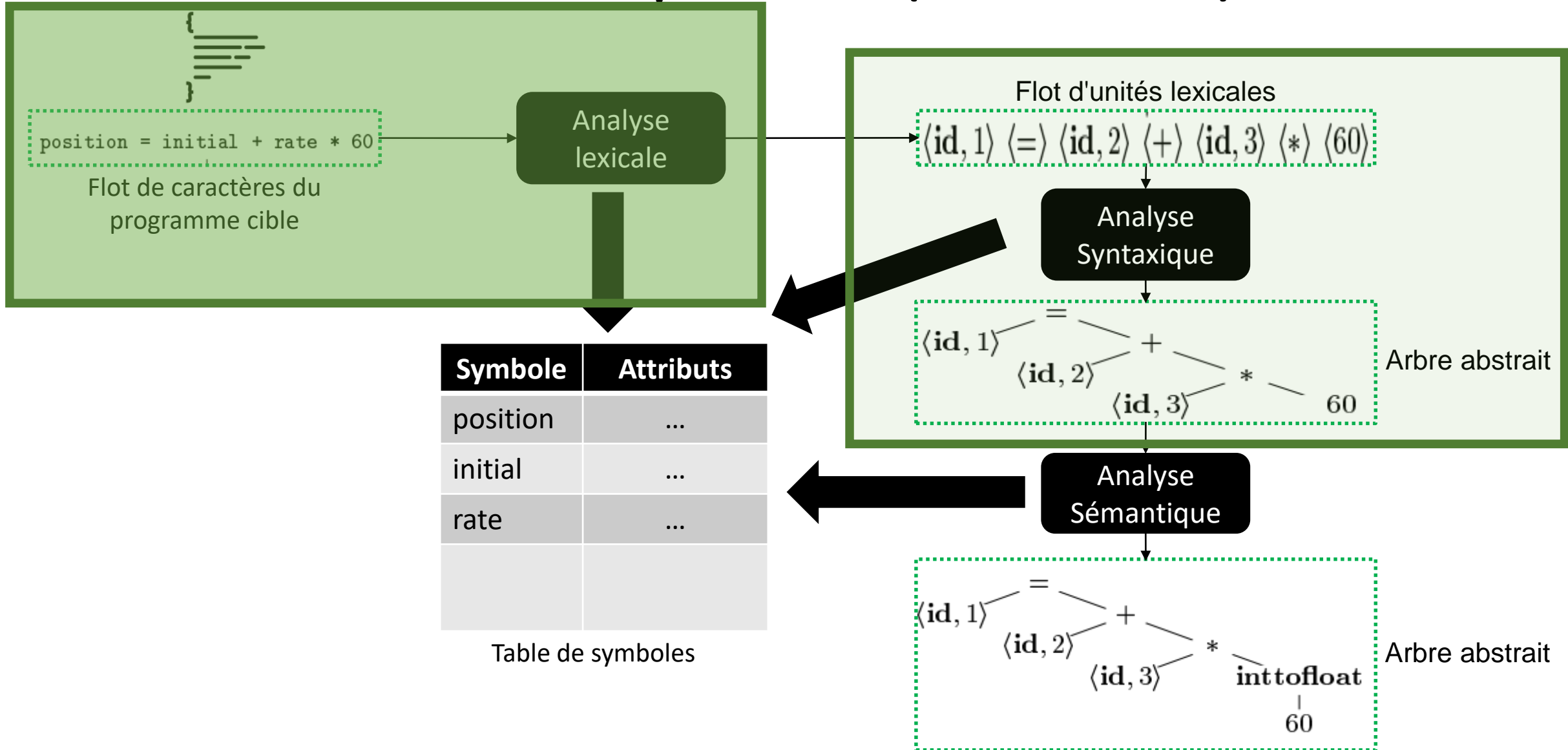
Compilation

Analyse syntaxique

Structure d'un compilateur (Front-end)



Structure d'un compilateur (Front-end)



C'est quoi l'analyse lexicale ? (Exemple)

- Chaine de caractères:

`position = initial + rate * 60`

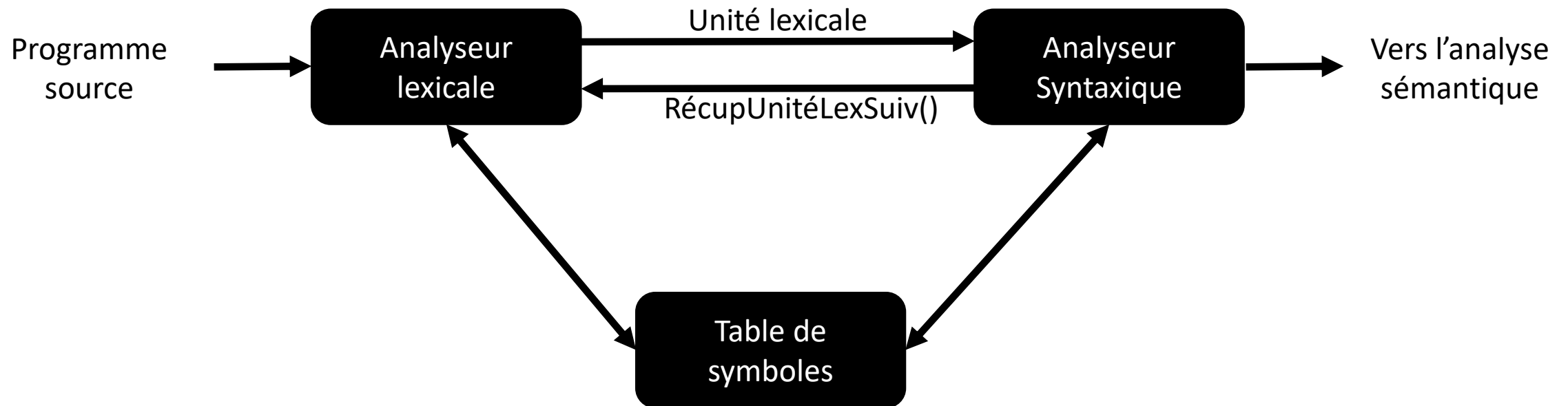
- Chaine des unités lexicales

< id, 1 >, < = >, < id, 2 >, < + >, < id, 3 >, < * >, < nbr, 60 >

- Table de symbols

Symbole	Attributs
position	...
initial	...
rate	...

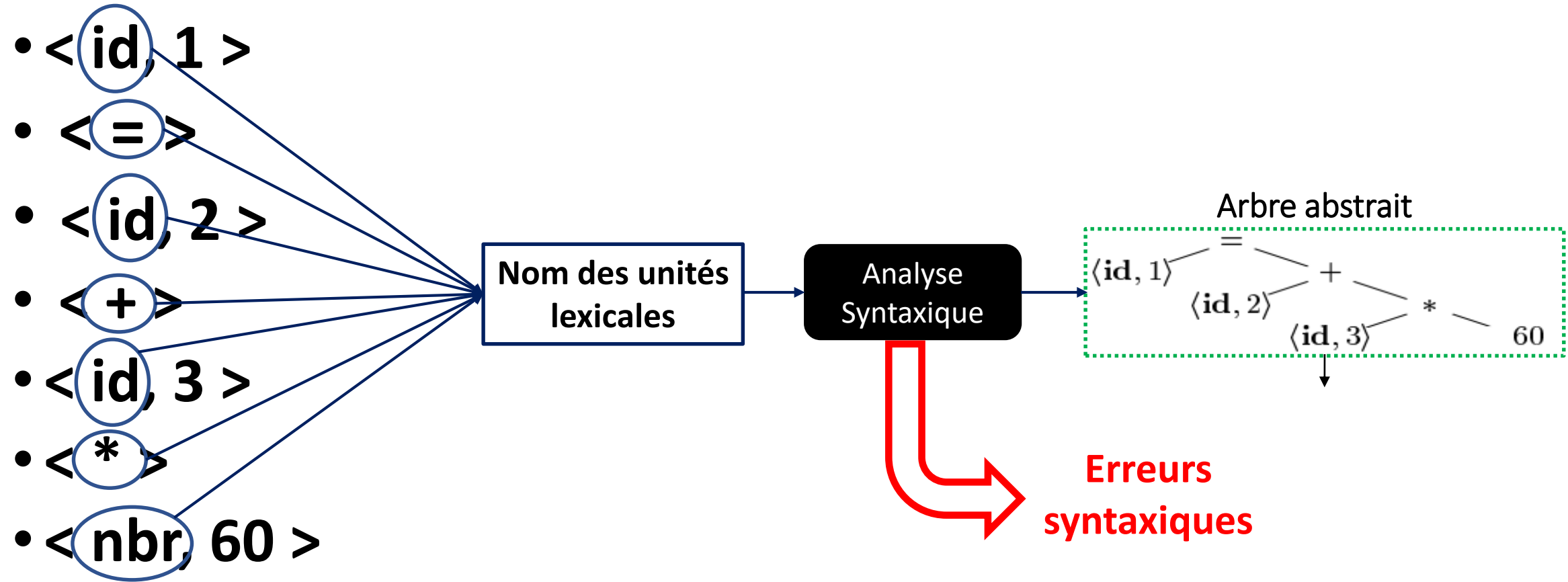
C'est quoi l'analyse lexicale ?



C'est quoi l'analyse syntaxique ?

- Reçoit une chaine d'unité lexicales fournie par l'analyseur lexicale
- Vérifier si cette chaine d'unité lexicales respecte la grammaire du langage source
- Signaler les erreurs syntaxiques
- Construire un arbre d'analyse pour les phases suivantes de compilation

C'est quoi l'analyse syntaxique ?



Comment effectuer l'analyse syntaxique ?

- Dans l'analyse lexicale :
 - **Spécification** : Expressions régulières
 - **Implémentation** : Automates finis
- Question

Peut-on utiliser les expressions régulières & les automates finis pour l'analyse syntaxique ???

Exemple

- Une expression est bien parenthésée. si **le nombre de parenthèses ouvrantes est égal au nombre de parenthèses fermantes**

- **(())**

- **((() ()))**

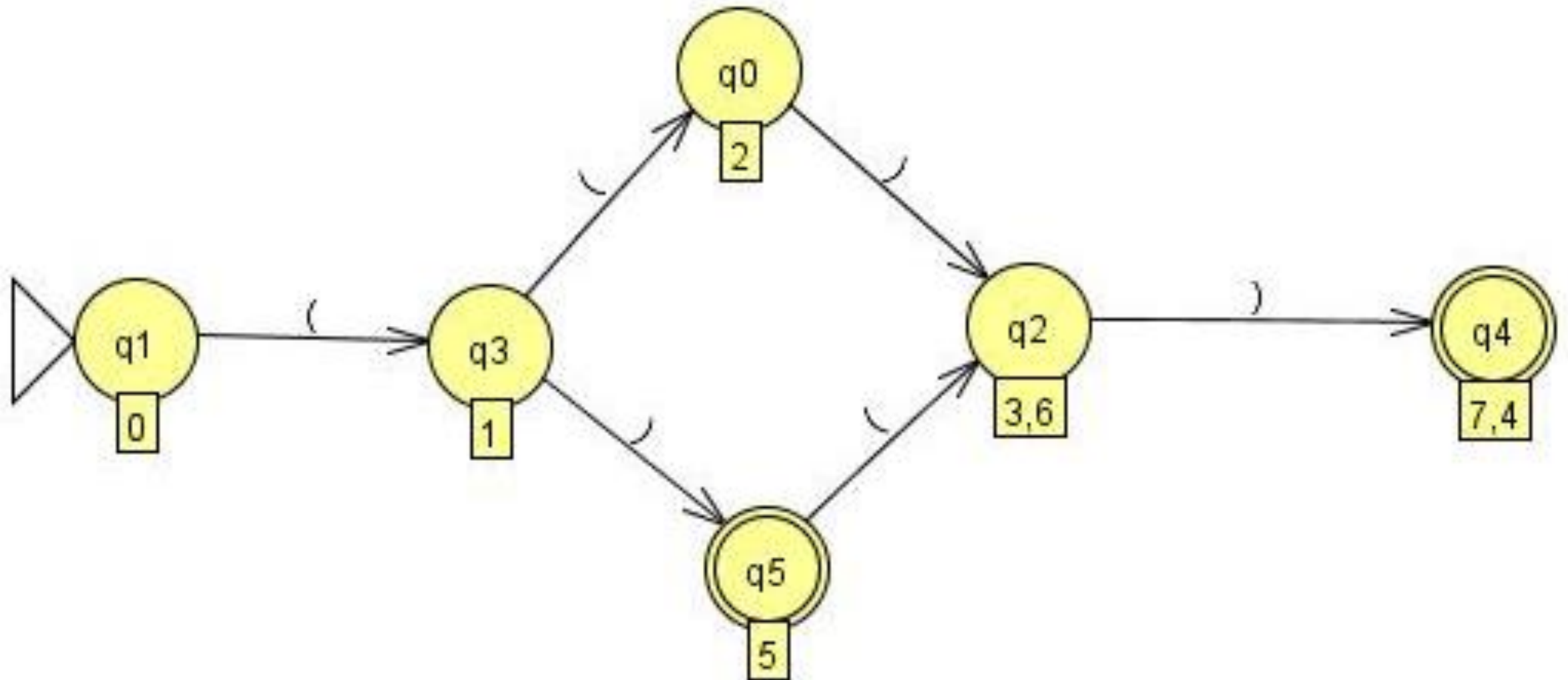


- **) (**

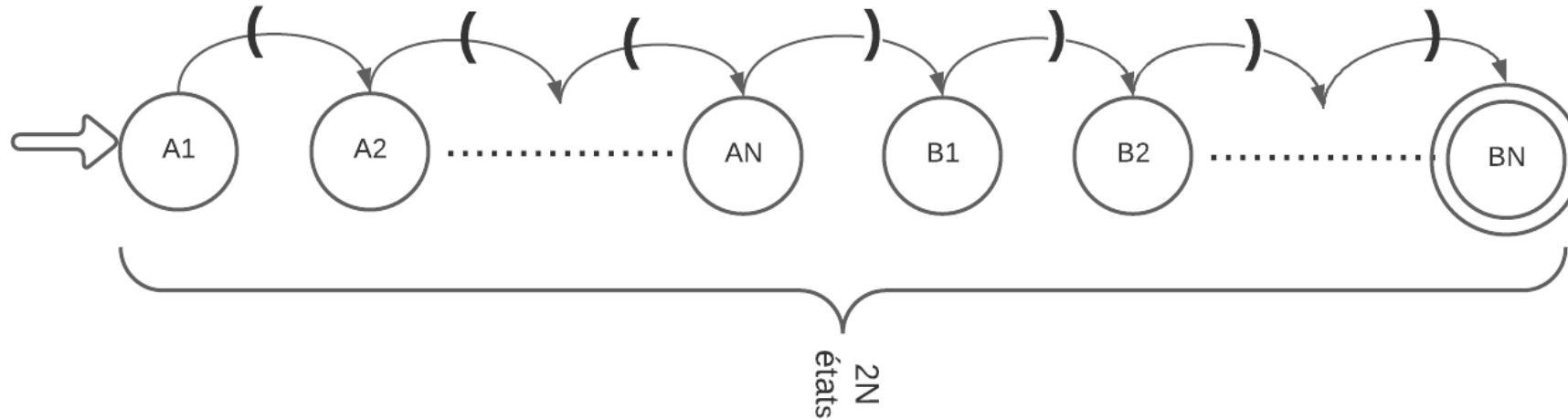
- **((()**



Example



Exemple



- Chemin avec **2N** états pour **N « (» + N «) »**
- Si on ajoute 2 « (» + «) » parenthèses \rightarrow Ajout de 2 états
- Automate n'est pas fini ($N \rightarrow \infty$)

L'automate fini est incapable de compter

Problème des expression régulières et les automate finis

- $L(\text{Expressions régulières}) = L(\text{automates finis})$
- Les automates finis

oublie l'historique



Expressions régulières (Automates finis)

- Incapable de spécifier la règle de parenthésage
- Plusieurs constructions de langages de programmation sont similaires au parenthésage et nécessitant le comptage :
 - **Pascal** : Begin End (nombre(**Begin**) = nombre (**End**))
 - **C**: { } (nombre({) = nombre (}))
- En générale le langage $a^n b^n$ n'est pas un langage régulier

Nécessité d'un outil plus puissant
(Les grammaire non contextuelle)

Grammaire non contextuelle (Context-free grammar (CFG))

- $G = \langle S, T, N, P \rangle$
- $S \in N$: **Symbole non terminal de départ**
- T : **les terminaux** sont les nom des unités lexicales
 - Nbr, Id, =, +, if, then
- N : **Les non terminaux** sont des variables syntaxiques utilisées durant les dérivations des chaines de terminaux
 - A, B, *expression*, *affectation*
- P : **les Productions** définissent les règles de remplacement de non terminaux
 - $N \rightarrow (T \cup N)^*$
 - $A \rightarrow \alpha A$ (**exemple**)

Exemple:

- $G = \langle S, T, N, P \rangle$

- Non terminal de départ :

$$S = Expr$$

- Terminaux:

$$T = \{id, nbr, +, -, *, (,)\}$$

- Non terminaux:

$$N = \{Expr, Op\}$$

- Productions:

- $expr \rightarrow (expr) \mid expr \ op \ expr \mid \mathbf{id} \mid \mathbf{nbr}$
 $op \rightarrow + \mid - \mid *$

Exemple:

- $expr \rightarrow (expr) \mid Expr \ op \ Expr \mid \mathbf{id} \mid \mathbf{nbr}$
 $op \rightarrow + \mid - \mid *$

- Mots acceptés par cette grammaire

- **id**
- **nbr**
- **nbr + nbr**
- **id + nbr**
- **Id + (id + nbr)**

- Mots non acceptés

- **+**
- **-**
- **Id +***
- **)nbr)**

Conventions de notation

- Terminaux
 - Premiers lettres miniscules : a, b, c
 - Opérateurs: $+, *, \dots$
 - Signes de ponctuation: Parenthèses, virgules
 - Chiffres: $0, 1, \dots, 9$
 - Chaine de caractères gras : **id, if, nbr,...** (une chaine représente un seul alphabet)
- Non terminaux
 - Premiers lettres majuscules : A, B, C
 - S quand utilisé représente un **Non terminal de départ**
 - Mots minuscules italique : $expr, instr$
 - D'autre lettres majuscules: E, T, F

Conventions de notation

- **Symboles grammaticaux génériques (Terminale ou Non-terminale):**
 - Derniers lettres majuscules: X, Y, Z (un seul symbole)
- **Chaine de terminaux (possiblement vide):**
 - Derniers lettres minuscules: u, v, \dots, z (chaine de terminaux)
- **Chaine de symbole grammaticaux (Terminaux ou Non-terminaux)**
 - Lettres grecques minuscules: α, β, γ (chaine qui contient tous les symboles)
- Si le **non terminal de départ** n'est pas spécifié, alors c'est la **partie gauche de la première production**.

Dérivation des mots

- **Grammaire :**

- $E \rightarrow E + E \mid E * E \mid E \mid (E) \mid \mathbf{id}$

- Dérivation du mot **$(\mathbf{id} + \mathbf{id}) * \mathbf{id}$**

$$E \xRightarrow{E \rightarrow E * E} E * E$$

$$\xRightarrow{E \rightarrow (E)} (E) * E$$

$$\xRightarrow{E \rightarrow E + E} (E + E) * E$$

$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + E) * E$$

$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{id}) * E$$

$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{id}) * \mathbf{id}$$

Dérivation des mots

- **Grammaire :**

- $E \rightarrow E + E \mid E * E \mid E \mid (E) \mid \mathbf{id}$

- Dérivation du mot **(id + id) * id**

$$E \xRightarrow{E \rightarrow E * E} E * E$$

$$\xRightarrow{E \rightarrow (E)} (E) * E$$

$$\xRightarrow{E \rightarrow E + E} (E + E) * E$$

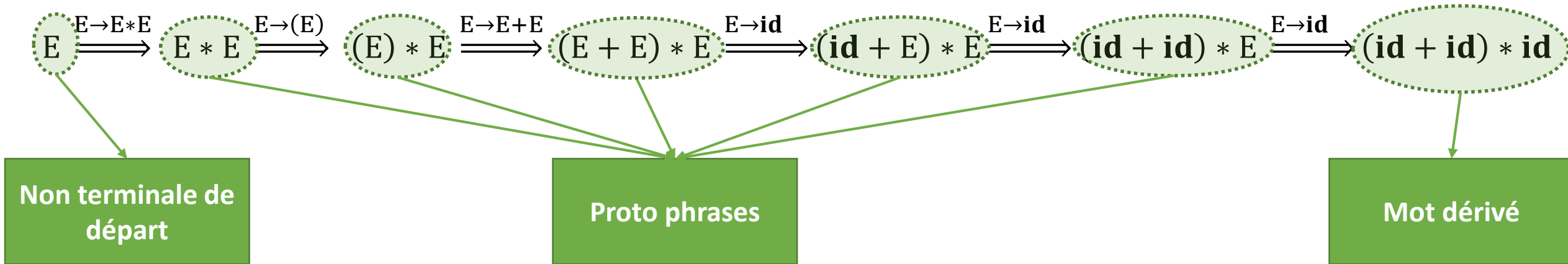
$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + E) * E$$

$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{id}) * E$$

$$\xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{id}) * \mathbf{id}$$

Dérivation des mots

- Dérivation du mot **(id + id) * id**



- $E \Rightarrow E * E \Rightarrow (E) * E \Rightarrow (E + E) * E \Rightarrow (\text{id} + E) * E \Rightarrow (\text{id} + \text{id}) * E \Rightarrow (\text{id} + \text{id}) * \text{id}$
- $E \xRightarrow{+} (\text{id} + \text{id}) * \text{id}$

Dérivation des mots

- Commencer à partir le non terminale de départ
- Remplacer à chaque étape **un non-terminale « A »** par la **partie droite « α »** d'une **production de la grammaire « $A \rightarrow \alpha$ »** :

$$S \xRightarrow{*} \beta \textcolor{red}{A} \gamma \xRightarrow{A \rightarrow \alpha} \beta \textcolor{green}{\alpha} \gamma$$

- Après plusieurs dérivations une proto phrase contenant que des alphabets est produite. Cette proto phrase finale représente un mot « w » reconnu par la grammaire

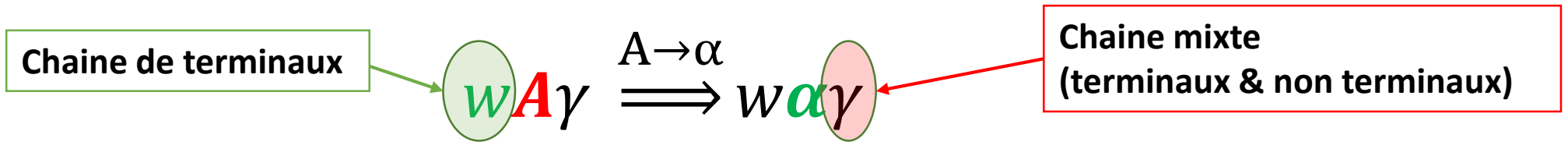
$$S \xRightarrow{+} w$$

- Le langage reconnu par une grammaire « G » inclus tous les mots qui peuvent être dérives par cette grammaire

$$LG(G) = \{w | S \xRightarrow{+} w\}$$

Dérivations gauches (Leftmost derivations)

- Choisir toujours le non-terminal le plus à gauche



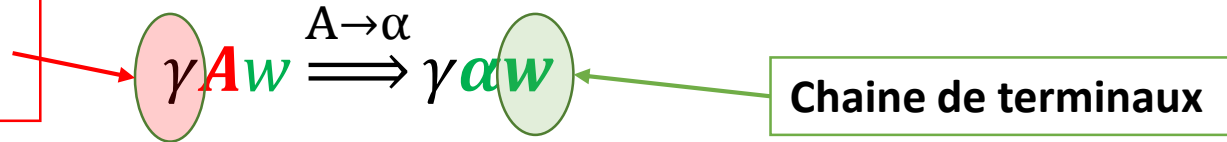
- Exemple:

- $E \Rightarrow \mathbf{E} * E$
 $\Rightarrow (\mathbf{E}) * E$
 $\Rightarrow (\mathbf{E} + E) * E$
 $\Rightarrow (\mathbf{id} + \mathbf{E}) * E$
 $\Rightarrow (\mathbf{id} + \mathbf{id}) * \mathbf{E}$
 $\Rightarrow (\mathbf{id} + \mathbf{id}) * \mathbf{id}$

Dérivations droites (Rightmost derivations)

- Choisir toujours le non-terminal le plus à gauche

Chaine mixte
(terminaux & non terminaux)



Chaine de terminaux

- Exemple:

- $E \Rightarrow E * E$

- $\Rightarrow E * E$

- $\Rightarrow E * id$

- $\Rightarrow (E) * id$

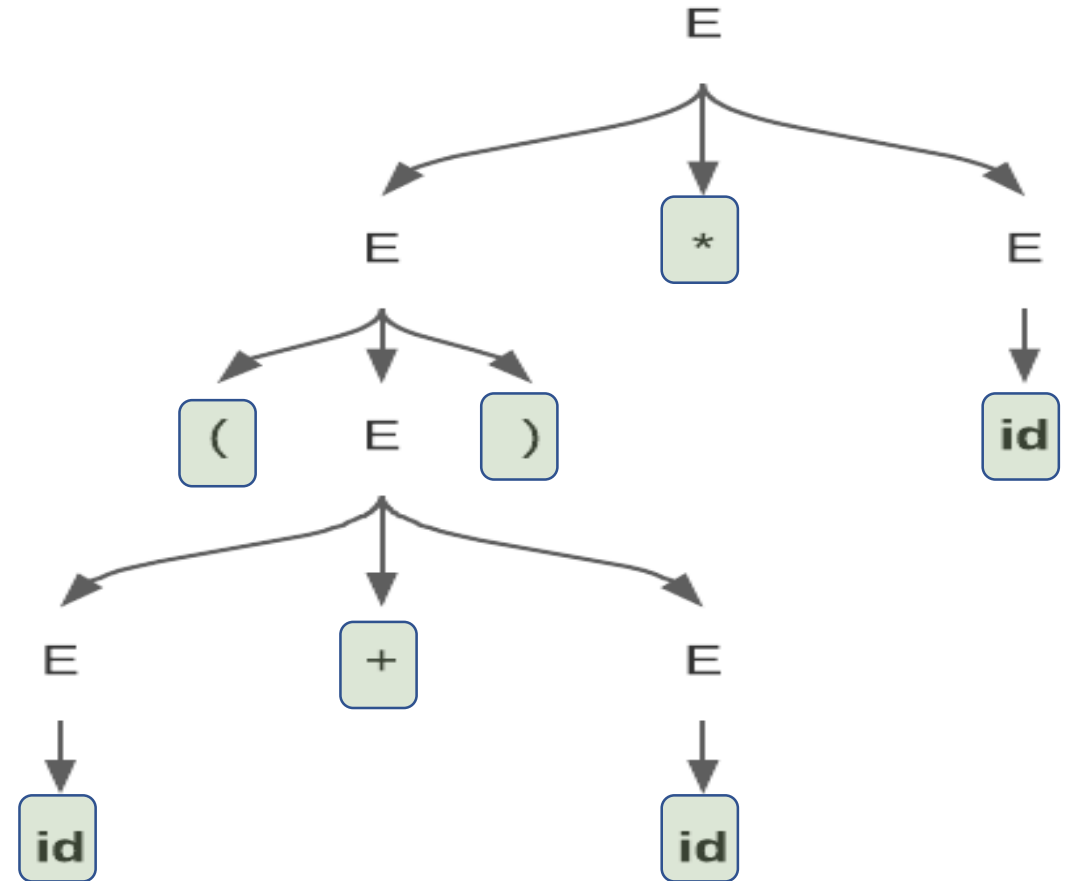
- $\Rightarrow (E + E) * id$

- $\Rightarrow (E + id) * id$

- $\Rightarrow (id + id) * id$

Arbre d'analyse syntaxique

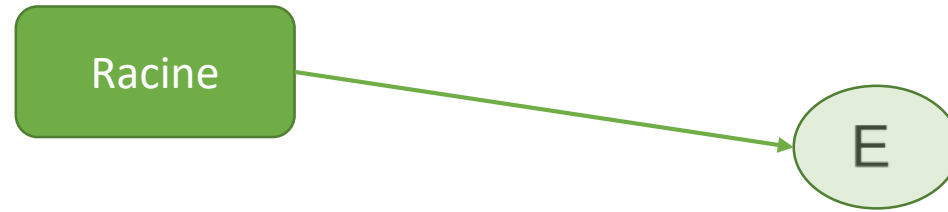
- $E \Rightarrow \mathbf{E} * E$
 $\Rightarrow (\mathbf{E}) * E$
 $\Rightarrow (\mathbf{E} + E) * E$
 $\Rightarrow (\mathbf{id} + \mathbf{E}) * E$
 $\Rightarrow (\mathbf{id} + \mathbf{id}) * \mathbf{E}$
 $\Rightarrow (\mathbf{id} + \mathbf{id}) * \mathbf{id}$



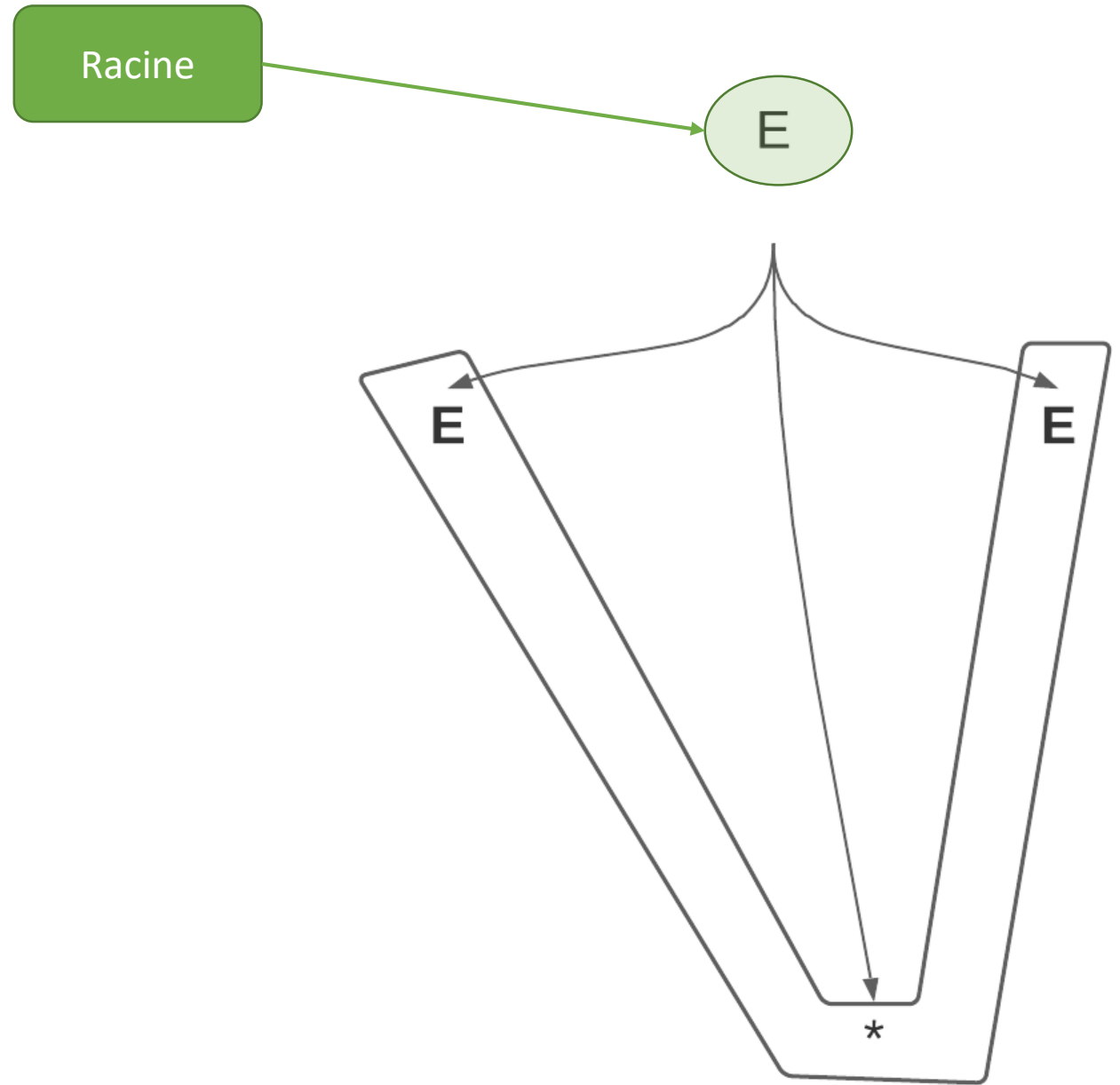
Analyse syntaxique descendante

- Commencer à partir du non terminal de départ (**Racine de l'arbre**) pour arriver aux alphabets du mots (**Feuillies de l'arbre**)
- Remplacer un non terminale « A » de la limite de l'arbre par une partie gauche « α » d'une production « $A \rightarrow \alpha$ »
- Arrêter quand tous les non terminaux sont remplacés
- Construire l'arbre à partir la racine tous en descendant vers les feuillies

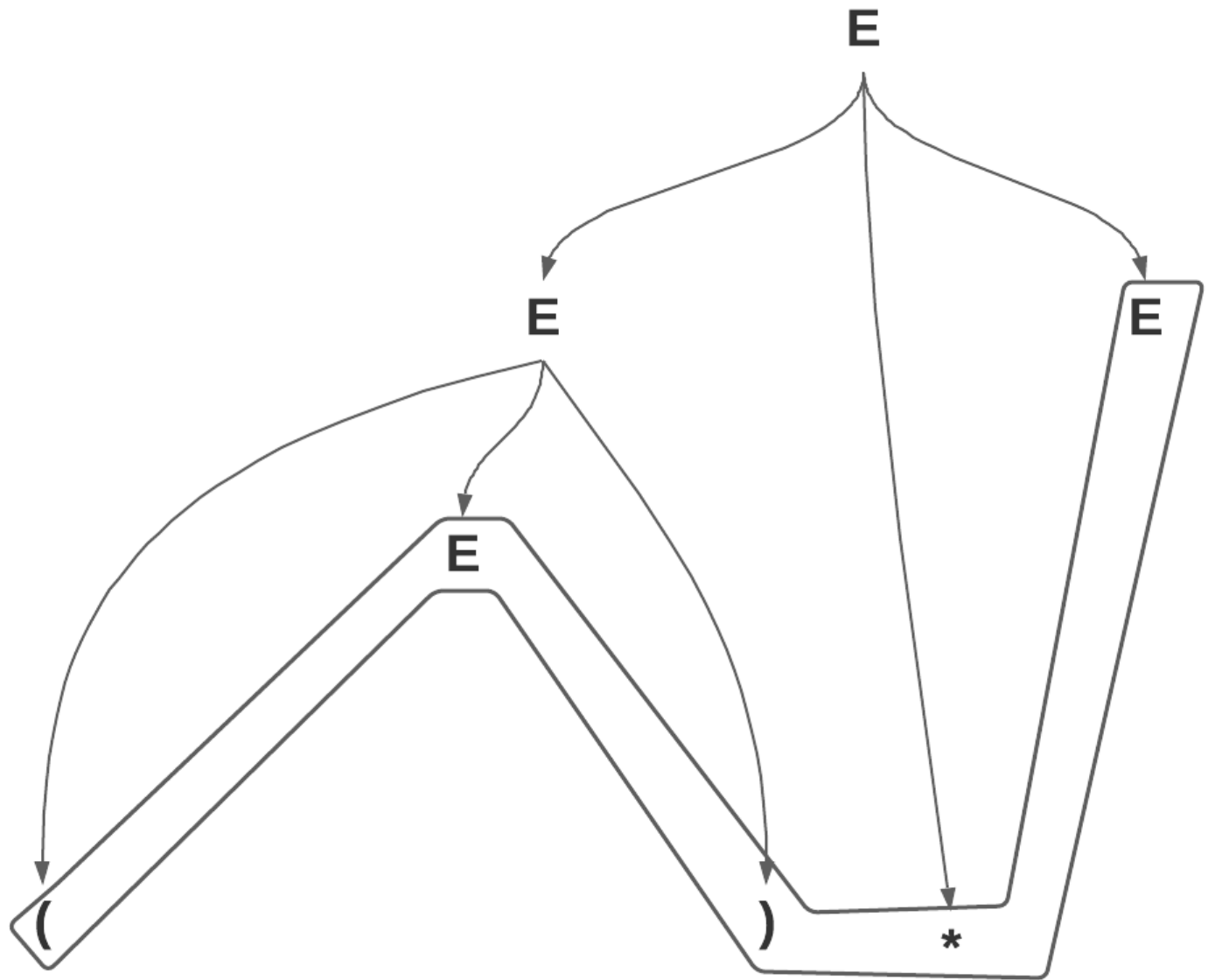
$E \Rightarrow$



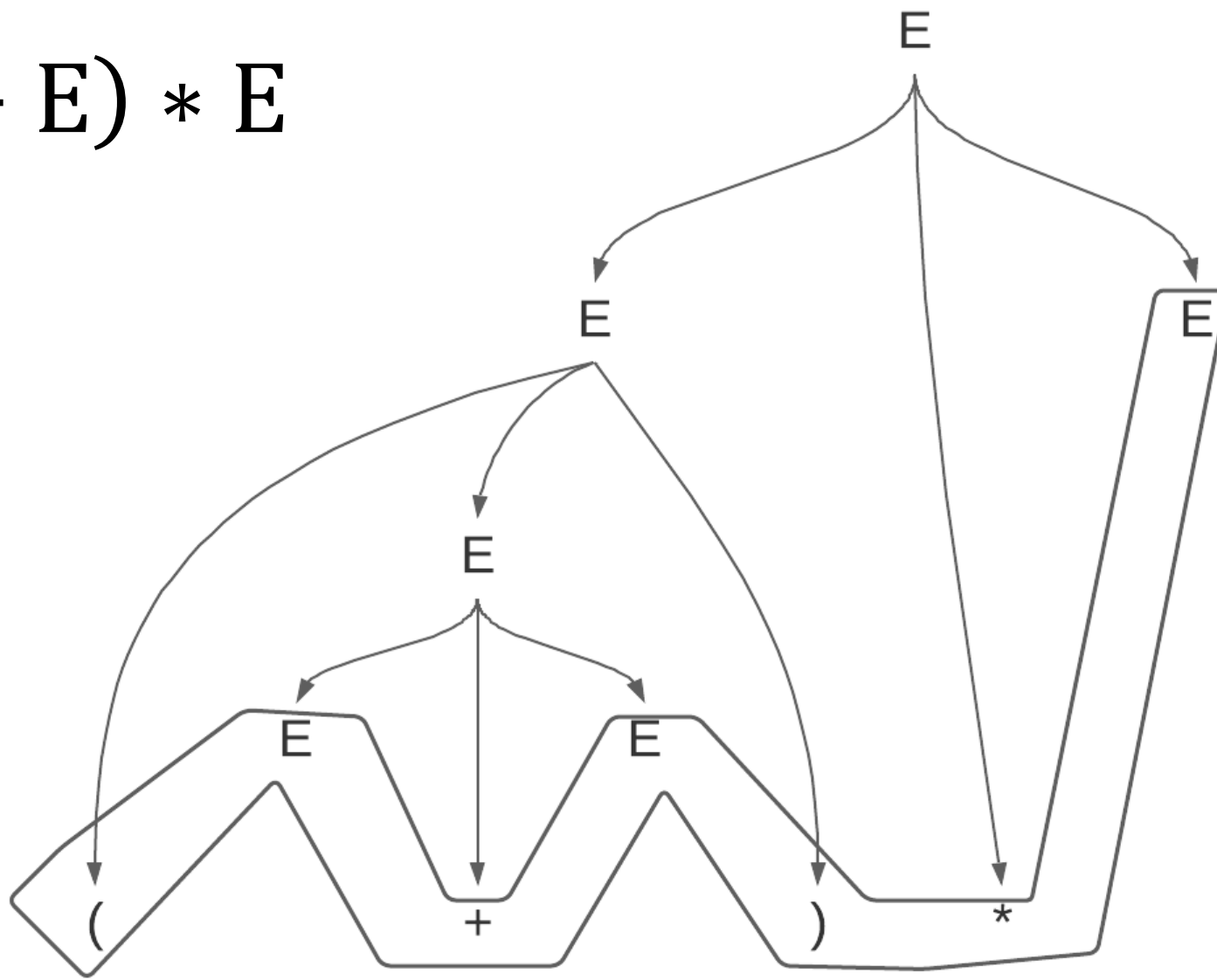
$$E \xrightarrow{E \rightarrow E * E} \mathbf{E} * E$$



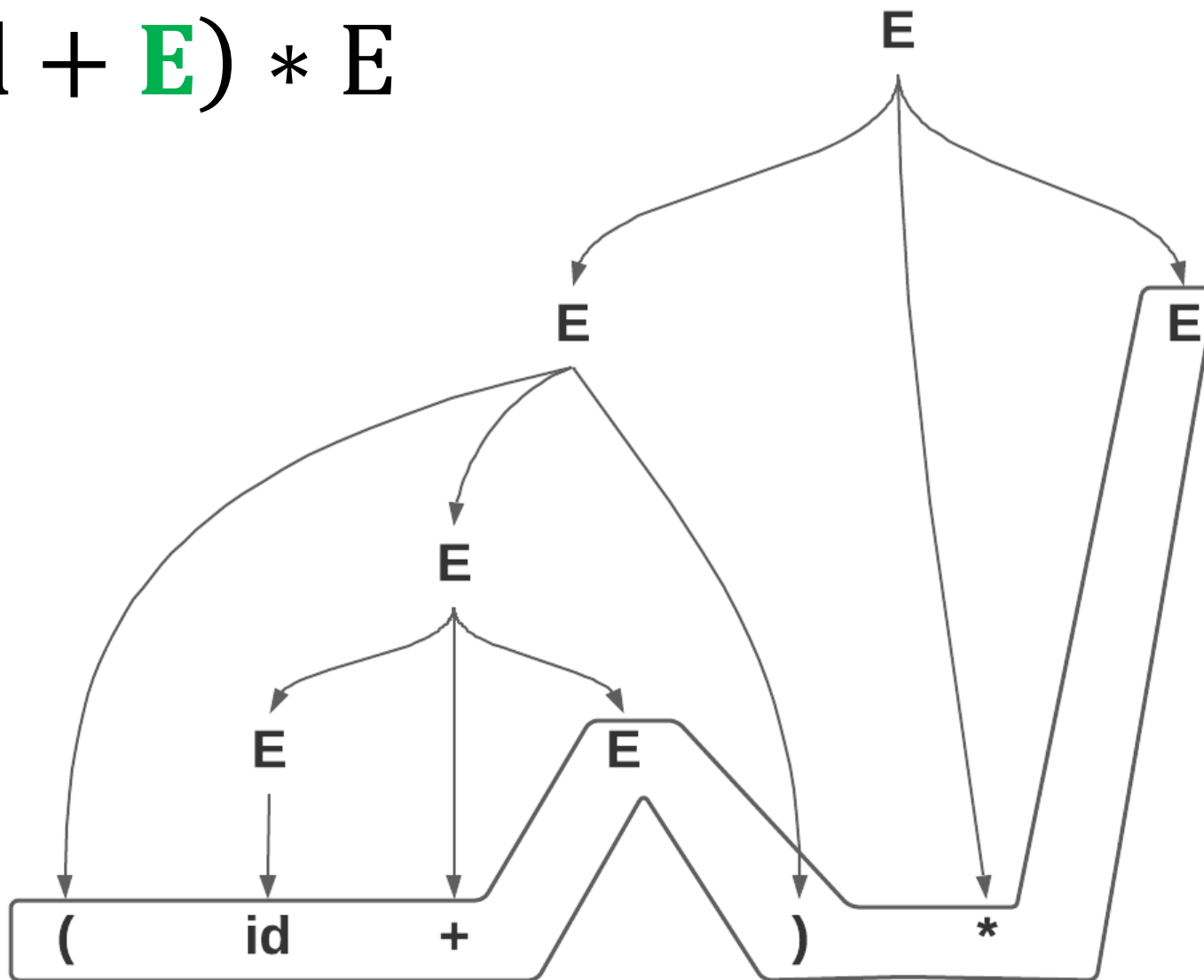
$$\mathbf{E} * E \xRightarrow{E \rightarrow (E)} (\mathbf{E}) * E$$



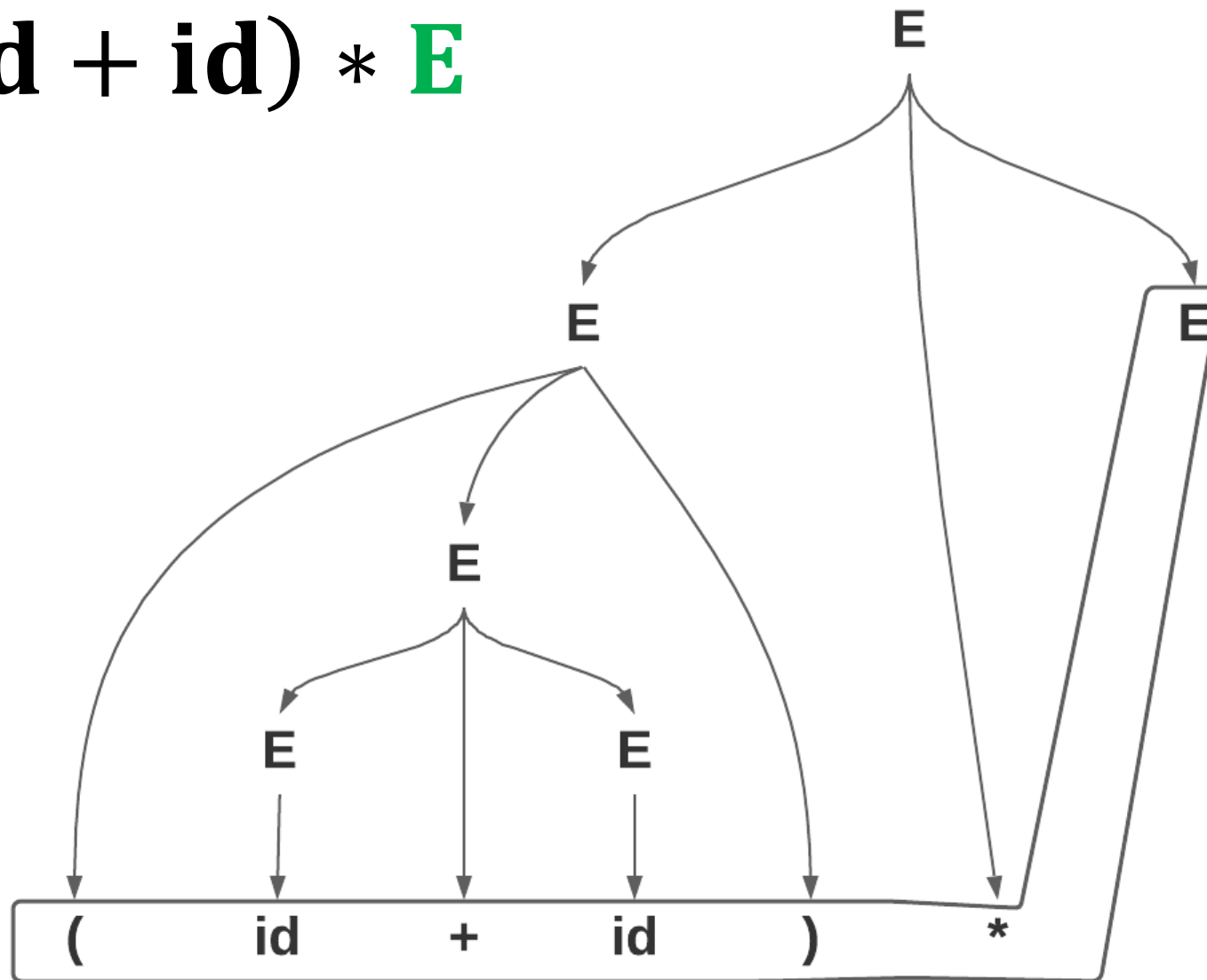
$$(\mathbf{E}) * E \xRightarrow{E \rightarrow E + E} (\mathbf{E} + E) * E$$



$$(\mathbf{E} + E) * E \xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{E}) * E$$



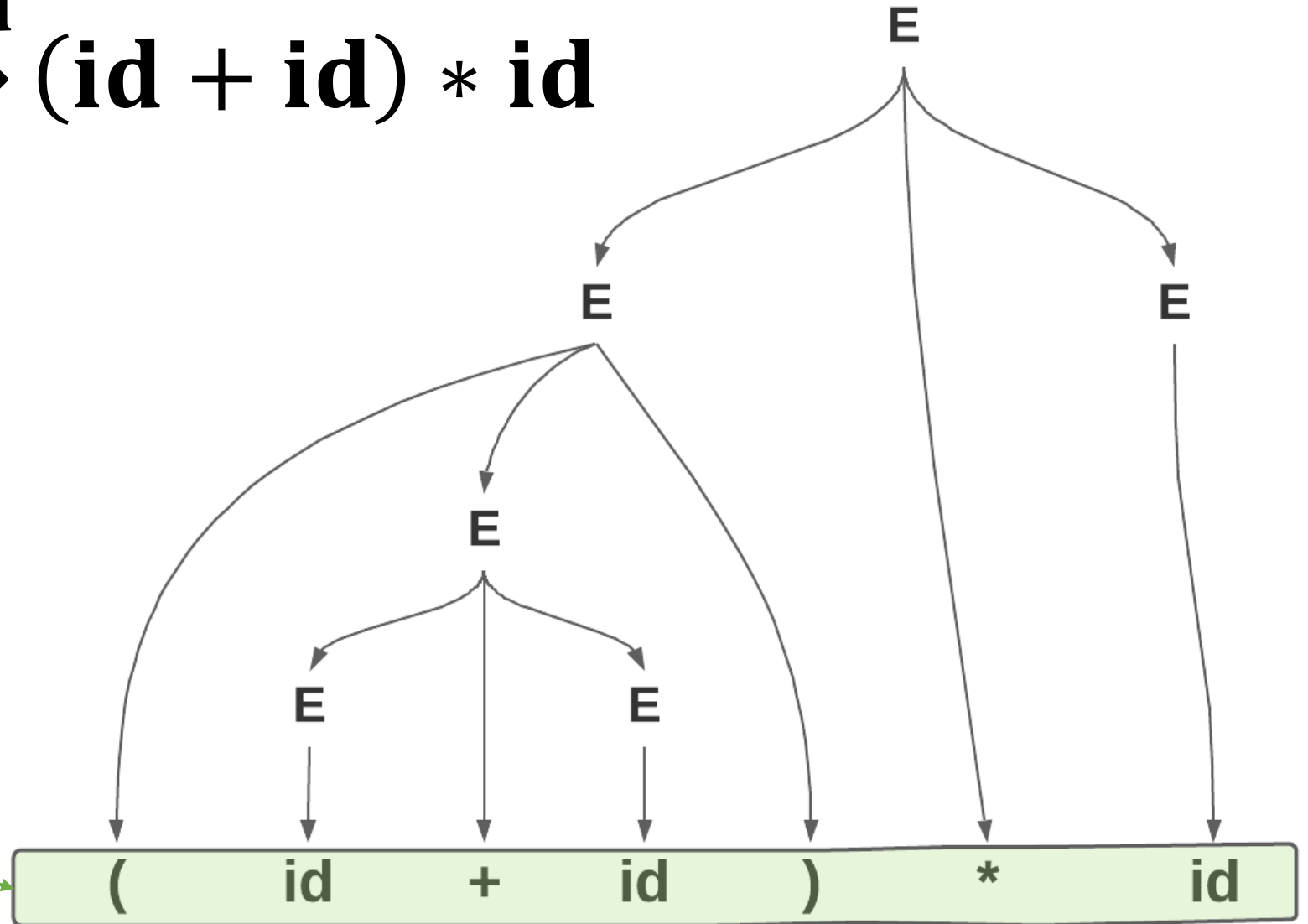
$$(\mathbf{id} + \mathbf{E}) * E \xRightarrow{E \rightarrow \mathbf{id}} (\mathbf{id} + \mathbf{id}) * \mathbf{E}$$



$$(\text{id} + \text{id}) * \textcolor{green}{E} \xRightarrow{E \rightarrow \text{id}} (\text{id} + \text{id}) * \text{id}$$



Mot dans les feuilles



Analyse syntaxique ascendante

- Commencer à partir le mot à reconnaître (**Feuillies de l'arbre**) pour arriver à la (**Racine de l'arbre**)
- Remplacer une partie droite « α » d'une production « $A \rightarrow \alpha$ » par son non terminale du partie gauche « A » (**Réduction**)
- Arrêter quand le non terminale de départ est atteint
- Construire l'arbre à partir les feuilles tous en montant vers la racine

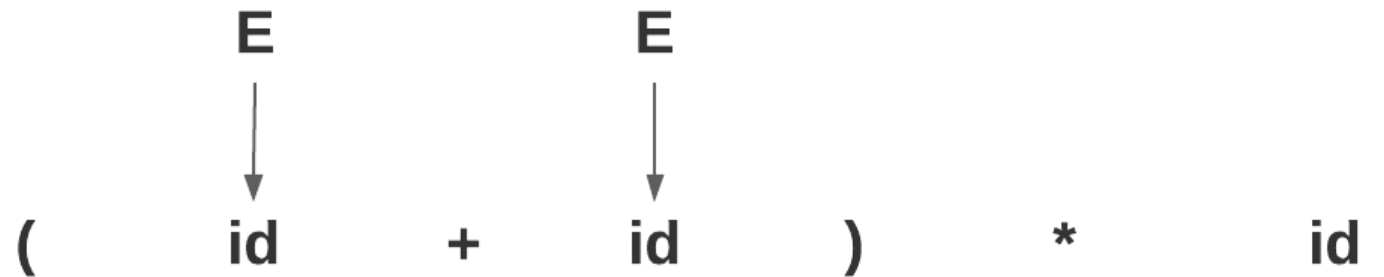
Mot dans les feuilles

(id + id) * id

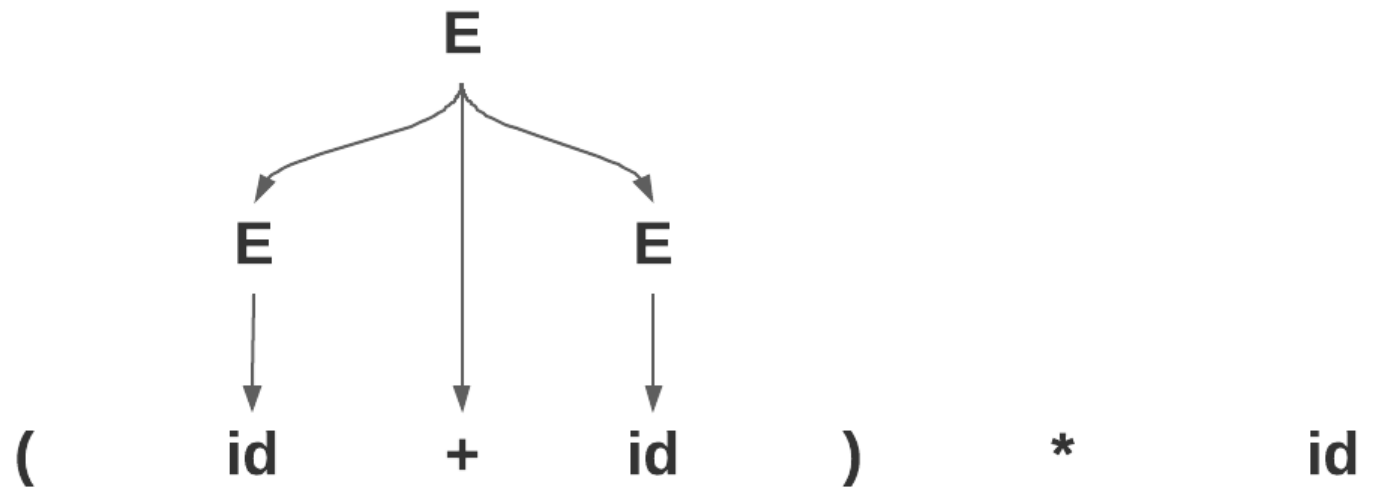
$$(\mathbf{E} + \mathbf{id}) * \mathbf{id} \Rightarrow (\mathbf{id} + \mathbf{id}) * \mathbf{id}$$

$$\begin{array}{c} \mathbf{E} \\ \downarrow \\ (\mathbf{id} + \mathbf{id}) * \mathbf{id} \end{array}$$

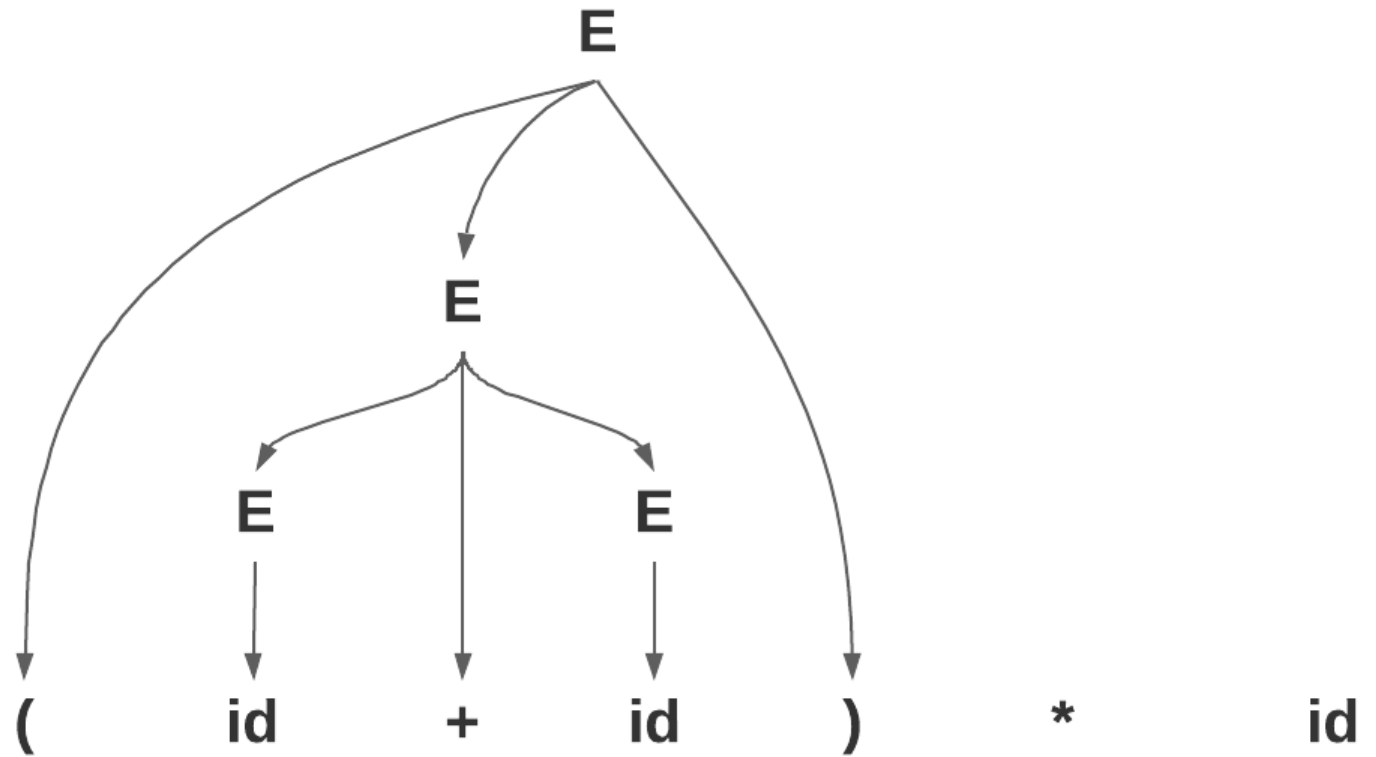
$$(E + \mathbf{E}) * \mathbf{id} \Rightarrow (\mathbf{E} + \mathbf{id}) * \mathbf{id}$$



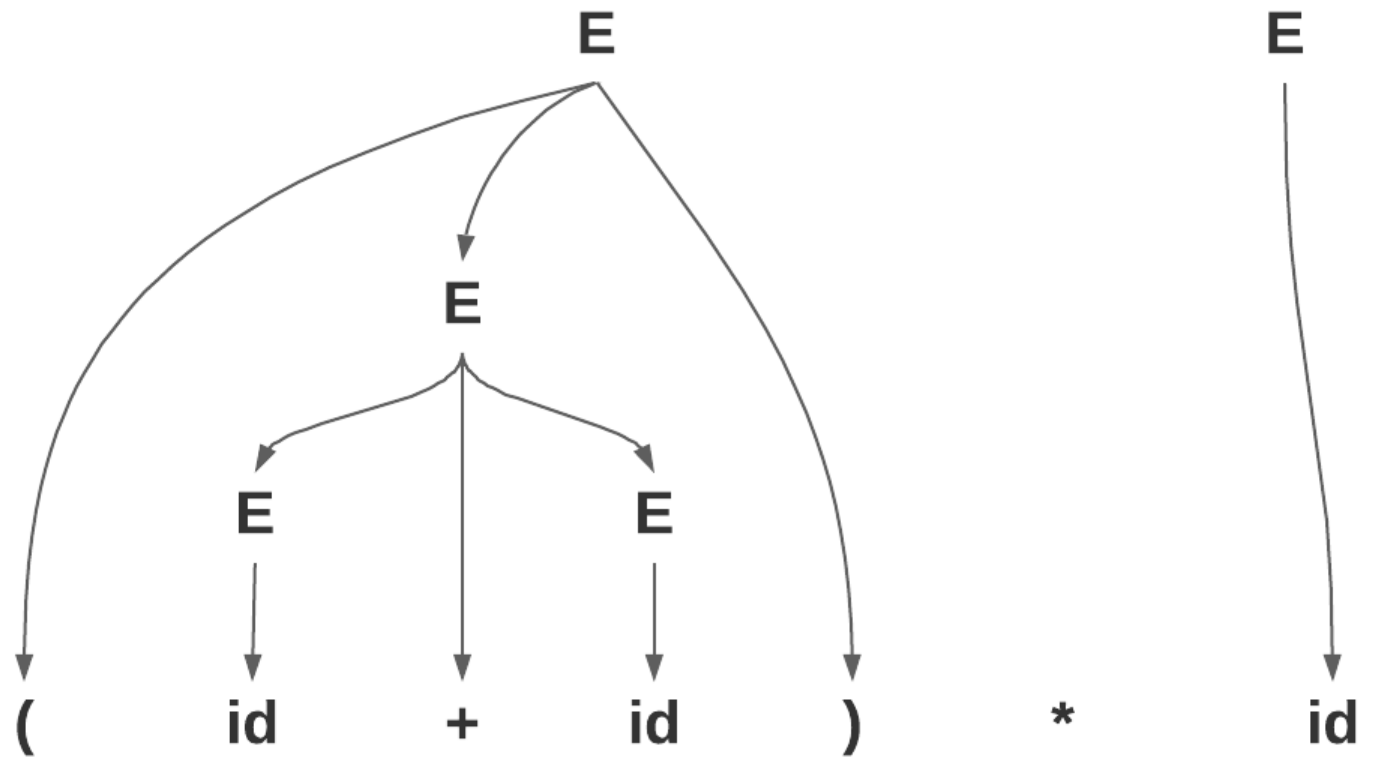
$$(\mathbf{E}) * \mathbf{id} \Rightarrow (E + \mathbf{E}) * \mathbf{id}$$



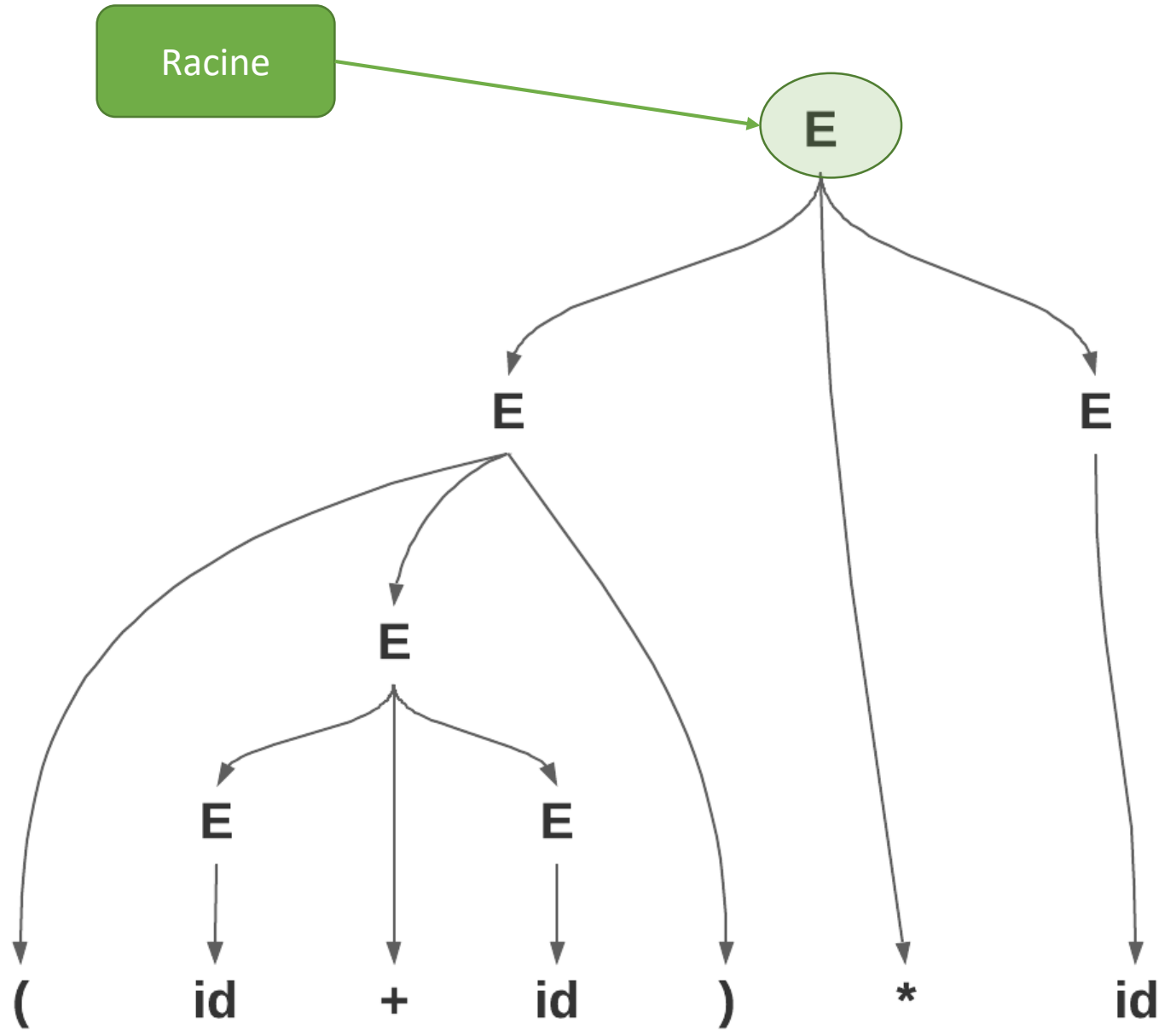
$$\mathbf{E} * \mathbf{id} \Rightarrow (\mathbf{E}) * \mathbf{id}$$



$$E * E \Rightarrow E * id$$



$$E \Rightarrow E * E$$



Ambiguïté

- Pour un mot donné il existe **plus d'un arbre d'analyse**
- Grammaire des expressions arithmétiques
$$E \rightarrow E + E \mid E * E \mid E \mid (E) \mid \text{id}$$

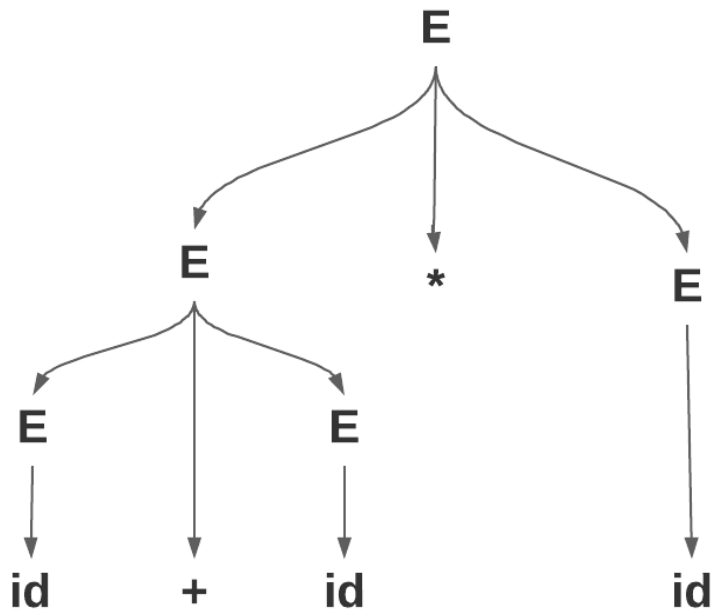
$$w = a + b * c$$

Exemple

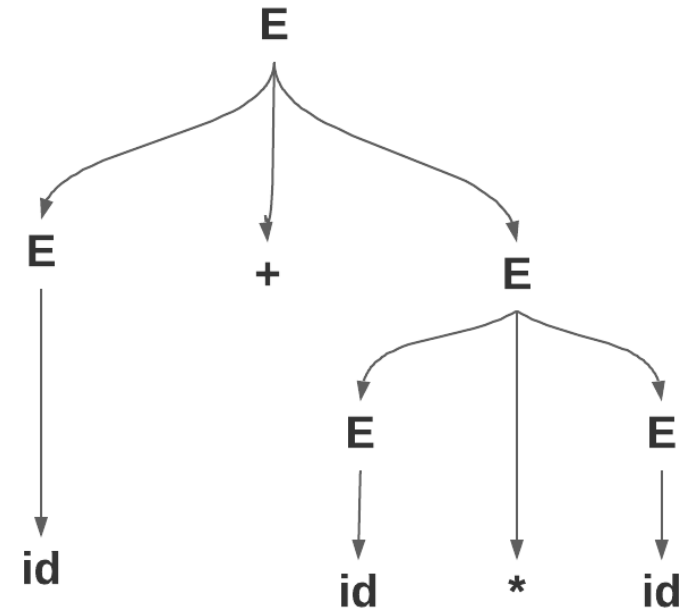
- Grammaire des expressions arithmétiques

$$E \rightarrow E + E \mid E * E \mid E \mid (E) \mid \text{id}$$

$w = \text{id} + \text{id} * \text{id}$



**Grammaire ne
considère pas la priorité
des opérateurs**

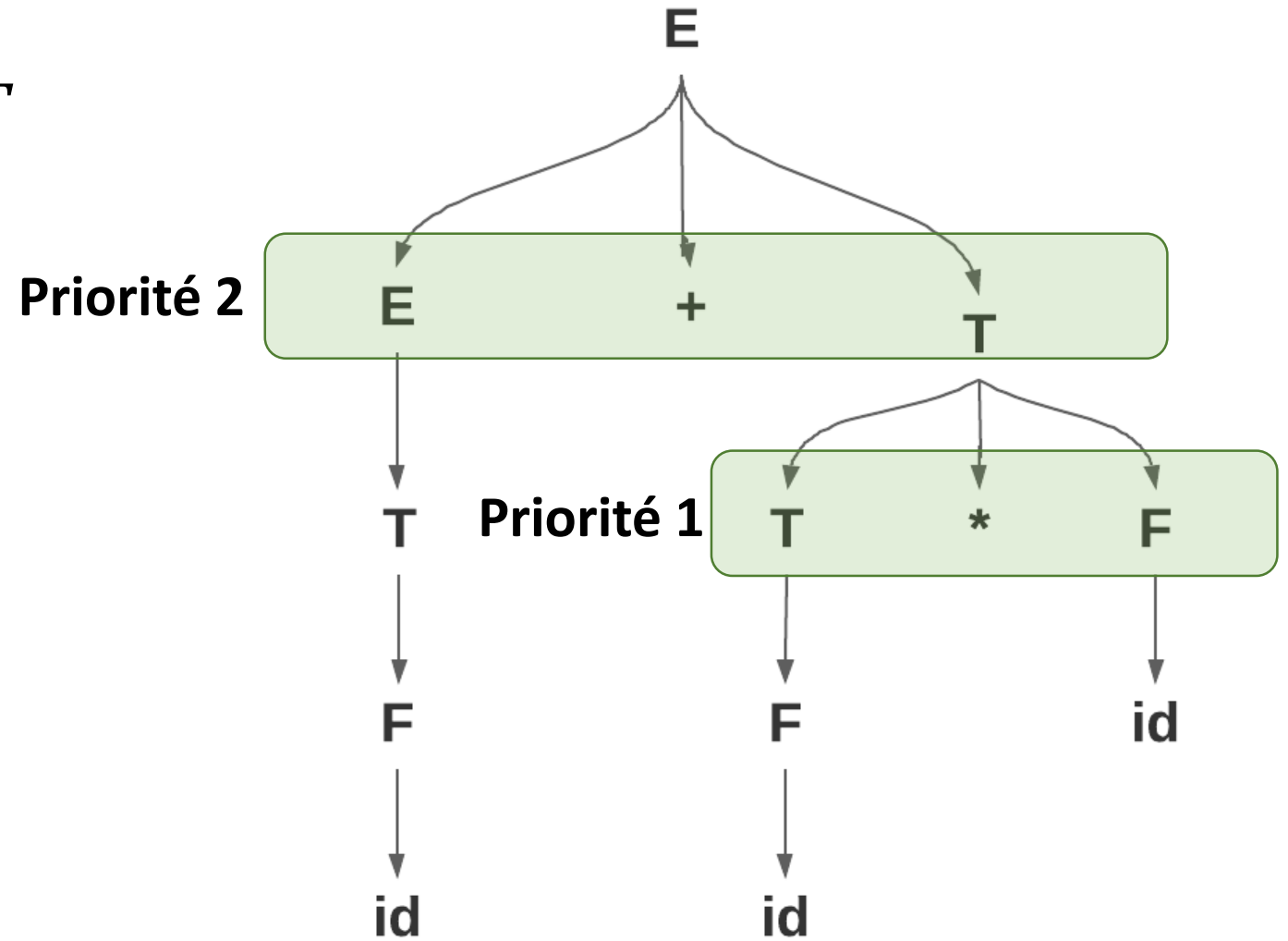


Enlever l'ambigüité

$$\bullet E \rightarrow E + T \mid E - T \mid T$$

$$\bullet T \rightarrow T * F \mid F$$

$$\bullet F \rightarrow (E) \mid \text{id}$$



Enlever l'ambiguïté

- Il n'existe pas un algorithme pour enlever l'ambiguïté de toutes les grammaires ambiguës
- **Un langage n'est pas ambigu** peut avoir **une grammaire ambiguë**
 - L'ambiguïté de la grammaire et n'est pas du langage
 - On peut enlever l'ambiguïté
- Il existe des langages **intrinsèquement ambigus**
 - Toutes les grammaires de ce type de langage est ambiguës
 - On ne peut jamais enlever l'ambiguïté

Exemple:

- $L = \{a^l b^m c^n \mid m = l \text{ ou } m = n\}$ est un langage **intrinsèquement ambigu**

Analyse syntaxique descendante

- Tester toutes les productions pour un non terminal
- Si une production échoue
 - Retour en arrière
 - Essayer une autre production
- Si une production réussit
 - Si on est dans le non terminal de départ alors le mot est reconnu

Analyse syntaxique descendante

A ()

Begin

temp = pos /*position courante le mot d'entrée*/

Pour chaque production $P_i \mid A \rightarrow X_1 X_2 \dots X_k$ /* tester toutes les productions*/

Début

pos \leftarrow temp

Pour i allant de 1 à k faire /* parcourir une production*/

Début

Si (X_i est un non terminal) Alors

$X_i(\quad)$

Sinon si (X_i est un terminal **ET** $X_i = \text{mot}[\text{pos}]$) Alors

pos \leftarrow pos+1

Sinon

break /*production échouée*/

Fin

si (i == k)

break /* production réussie*/

Fin

si (pos == temp)Alors

Ecrire('Analyse échouée')

Quitter le programme

sinon si(A est le non terminal de départ **ET** pos =longueur(mot)) Alors

Ecrire('mot reconnu')

Fin

Enlever la récursivité à gauche

- La **récursivité gauche** peut conduire à une **boucle infinie** quand les **dérivations gauches** sont utilisées

- Récursivité gauche directe

$$A \rightarrow A\alpha \mid \beta$$

- $A \rightarrow A\alpha \rightarrow A\alpha\alpha \rightarrow A\alpha\alpha\alpha \rightarrow \dots \rightarrow A\alpha^* \rightarrow \beta\alpha^*$

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \varepsilon \end{aligned}$$

Example

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

Factorisation de la grammaire

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \alpha\beta_3 \mid \dots \mid \alpha\beta_n \mid \gamma$$

Après factorisation

$$\begin{aligned} A &\rightarrow \alpha B \mid \gamma \\ B &\rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots \mid \beta_n \end{aligned}$$

Exemple

- $instr \rightarrow \mathbf{si} \ expr \ \mathbf{Alors} \ instr \ \mathbf{sinon} \ instr \mid \mathbf{si} \ expr \ \mathbf{Alors} \ instr$
- $instr \rightarrow \mathbf{si} \ expr \ \mathbf{Alors} \ instr \ sinon_instr$
- $sinon_instr \rightarrow \mathbf{sinon} \ instr \mid \varepsilon$

Analyse descendante LL(1)

- Elle analyse un mot d'entrée de gauche à droite (*Left to right* en anglais)
- Construit une dérivation à gauche (*Leftmost derivation* en anglais).
- Un seul caractère de prédiction

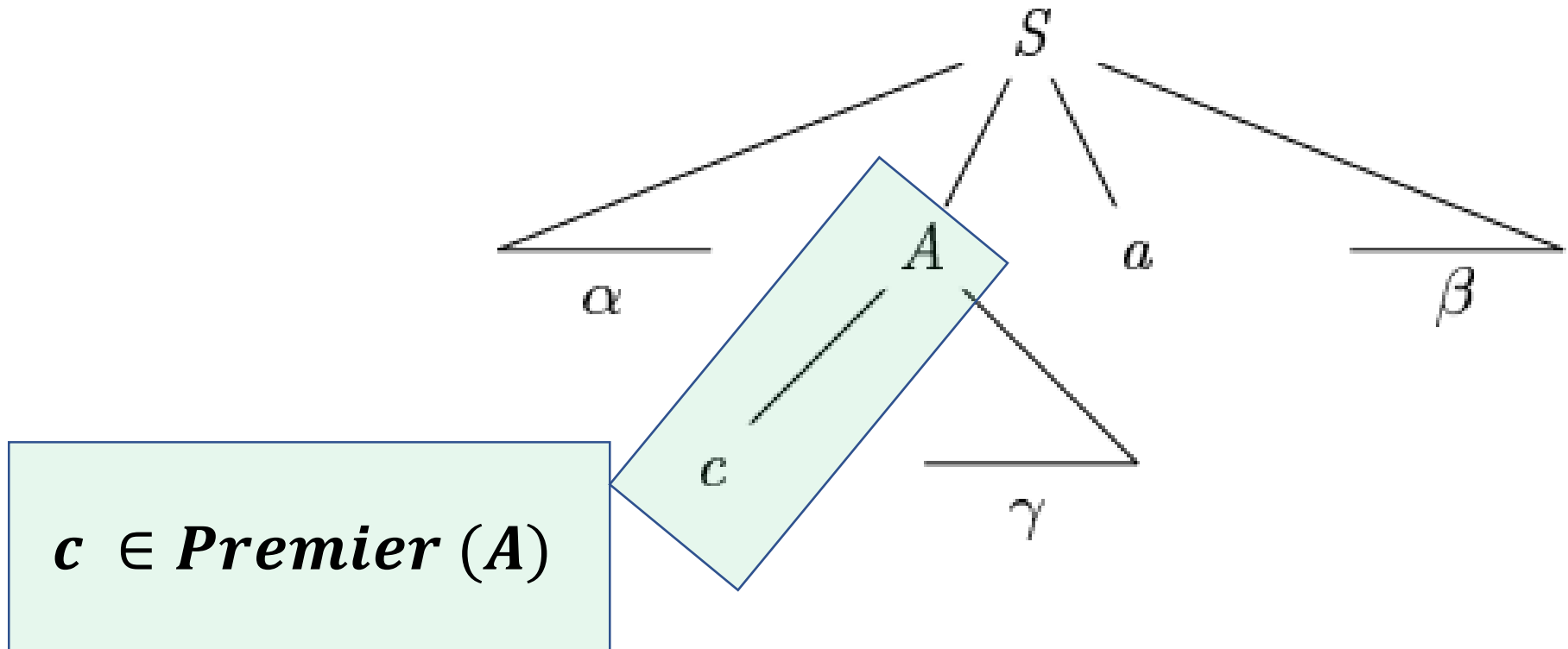
Ensemble Premier

- $Premier(\alpha) = \{a \in T \cup \{\varepsilon\} \mid \alpha \xRightarrow{+} a\beta\}$
- L'ensemble **Premier** contient des **terminaux** ou ε
- L'ensemble des terminaux qui peuvent apparaître **en premier (en début)** d'une proto-phrased après **une ou plusieurs dérivations** commençant de la chaîne α

$$A \xRightarrow{+} \textcolor{green}{a}B$$

Ensemble Premier

- $Premier(\alpha) = \{a \in T \cup \{\varepsilon\} \mid \alpha \xRightarrow{+} a\beta\}$



Calcule Premier

- Si X est un terminal

$$\text{Premier}(X) = \{X\}$$

- Si X est un non-terminal et $X \rightarrow Y_1 Y_2 \dots Y_k$

$$\text{Premier}(X) = \text{Premier}(Y_1) \text{ Si } \varepsilon \notin \text{Premier}(Y_1)$$

$$\text{Premier}(X) = (\text{Premier}(Y_1) \cup \text{Premier}(Y_2)) - \{\varepsilon\}$$

$$\text{Si } \varepsilon \in \mathbf{\text{Premier}(Y_1)} \text{ et } \varepsilon \notin \mathbf{\text{Premier}(Y_2)}$$

$$\text{Premier}(X) = \text{Premier}(Y_1) \cup \dots \cup \text{Premier}(Y_j) - \{\varepsilon\}$$

$$\text{Si } \varepsilon \in \mathbf{\text{Premier}(Y_1)} \text{ et } \varepsilon \in \mathbf{\text{Premier}(Y_2)} \dots \varepsilon \in \mathbf{\text{Premier}(Y_{j-1})} \text{ et } \varepsilon \notin \mathbf{\text{Premier}(Y_j)} : j \leq k$$

$$\varepsilon \in \text{Premier}(X)$$

$$\text{Si } \varepsilon \in \mathbf{\text{Premier}(Y_1)} \text{ et } \varepsilon \in \mathbf{\text{Premier}(Y_2)} \dots \varepsilon \in \mathbf{\text{Premier}(Y_k)}$$

Calcule Premier

- Après avoir calculer les ensembles **Premier** pour tous les symboles, on peut calculer **Premier** pour chaque n'importe chaine **mixte** α

$$\text{Si } \alpha = Y_1 Y_2 \dots Y_k$$

- $\text{Premier}(\alpha) = \text{Premier}(Y_1) \cup \dots \cup \text{Premier}(Y_j) - \{\varepsilon\}$
Si $\varepsilon \in \text{Premier}(Y_1)$ et $\varepsilon \in \dots \varepsilon \in \text{Premier}(Y_{j-1})$ et $\varepsilon \notin \text{Premier}(Y_j)$
- $\varepsilon \in \text{Premier}(X)$
Si $\varepsilon \in \text{Premier}(Y_1)$ et $\varepsilon \in \text{Premier}(Y_2) \dots \varepsilon \in \text{Premier}(Y_k)$

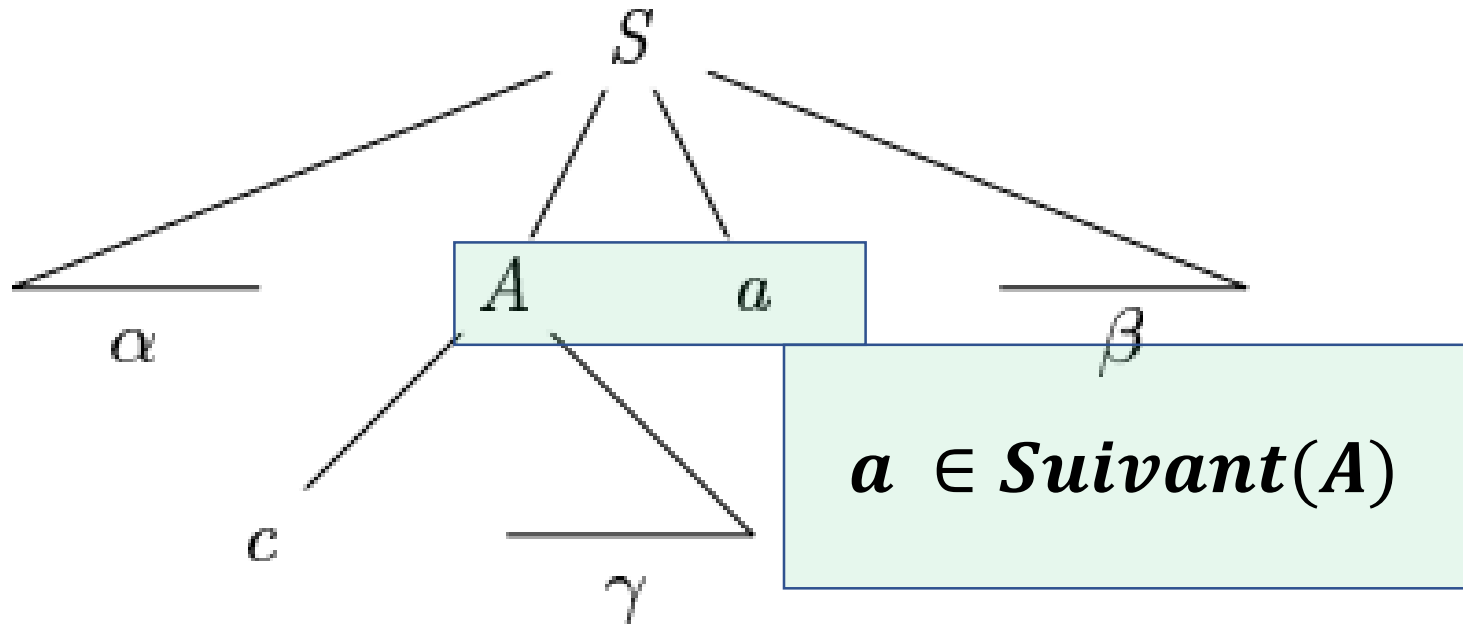
Ensemble Suivant

- $Suivant(\alpha) = \{a \in Premier(\gamma) \mid \alpha \xRightarrow{+} \beta\alpha\gamma\}$
- L'ensemble **Suivant** contient les **terminaux** qui peuvent apparaître **juste après un symbole** dans une suite de dérivations

$$\alpha \xRightarrow{+} \beta\alpha\gamma$$

Ensemble Suivant

- $Suivant(\alpha) = \{a \in Premier(\gamma) \mid \alpha \xRightarrow{+} \beta\alpha\gamma\}$



Calcule Suivant

- Mettre **\$** dans **Suivant(S)**
 - Si S est un non-terminale de départ
- Mettre **Premier(β)- $\{\epsilon\}$** dans **Suivant(B)**
 - Si il existe une production $A \rightarrow \alpha B \beta$
- Mettre **Suivant(A)** dans **Suivant(B)**
 - Si il existe une production $A \rightarrow \alpha B \beta$
 - Si ϵ est dans **Premier(β)** ou $\beta = \epsilon$

Analyseur syntaxique descendant prédictif

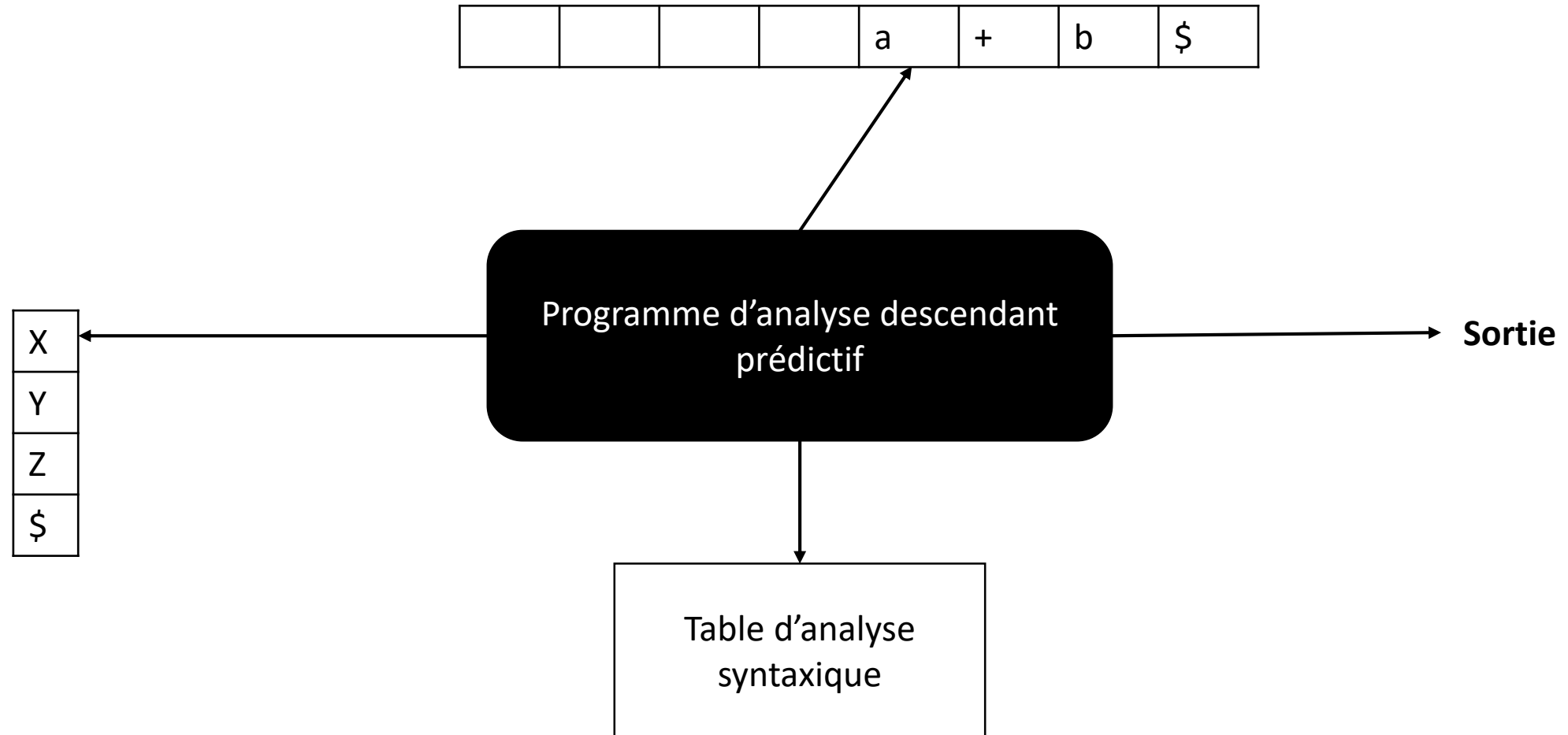


Table d'analyse LL(a)

Si $a_i \in \text{Premier}(\alpha)$

	a_0	a_i		\$
A_0					
A_i			$A_i \rightarrow \alpha$		
A_n					

Si $\varepsilon \in \text{Premier}(\alpha)$ et $a_i \in \text{Suivant}(A_i)$

Table d'analyse LL(a)

Si $a_i \in \text{Premier}(\text{Premier}(\alpha))$
 $\text{Suivant}(A_i)$

	a_0	a_i		\$
A_0					
A_i			$A_i \rightarrow \alpha$		
A_n					

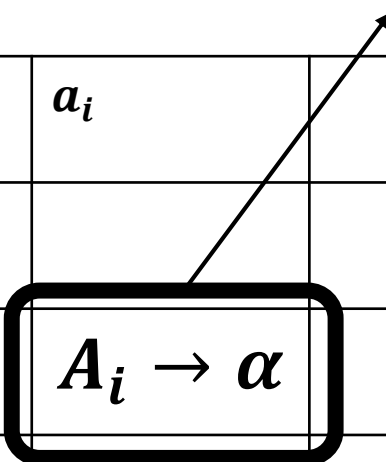


Table d'analyse LL(a)

Accéder à une Case
vide \rightarrow ERREUR

Deux dérivations possibles \rightarrow
Grammaire n'est pas LL(1)

	a_0	a_i		\$
A_0					
A_i			$A_i \rightarrow \alpha$ $A_i \rightarrow \beta$		
A_n					

Reconnaissance d'un mot

- Dans le TD