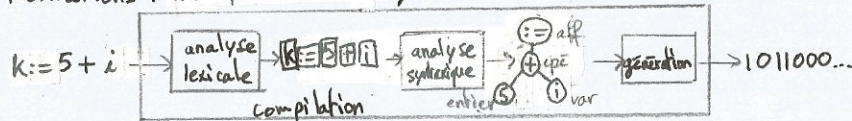


* Introduction: La Compilation.

Motivations: interpréter un code / convertir du LaTeX en HTML / ...



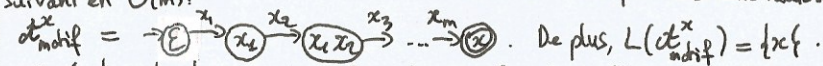
I) Rappels de théorie des langages 1) langages rationnels (RAT)

def 1: un automate est un quintuplet $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ où Q est l'ensemble fini des états, Σ est l'alphabet (ensemble fini), $\delta \subseteq Q \times \Sigma \times Q$ est l'ensemble des transitions, $I \subseteq Q$ les états initiaux, $F \subseteq Q$ les états finaux.

def 2: un chemin dans \mathcal{A} est une suite finie de transitions $q_0 \xrightarrow{x_1} q_1 \xrightarrow{x_2} \dots \xrightarrow{x_n} q_n$, où $q_0 \in I$. On dit que le mot $x_1 \dots x_n$ est l'étiquette du chemin et q_n l'état d'arrivée.

def 3: un mot $x \in \Sigma^*$ est accepté par l'automate \mathcal{A} si il est l'étiquette d'un chemin dans \mathcal{A} d'état d'arrivée $q_n \in F$. On note $L(\mathcal{A})$ le langage des mots acceptés par \mathcal{A} .

ex 4: (automate du motif). Soit $x = x_1 \dots x_m \in \Sigma^*$. On peut construire l'automate suivant en $O(m)$:



ex 5: (automate des occurrences). Soit $x \in \Sigma^*$, P l'ensemble des préfixes de x . L'automate $\mathcal{A}_x = (P, \Sigma, \delta, \{ \epsilon \}, P)$, où $\delta = \{ (u, a, f_x(ua)), u \in P, a \in \Sigma \}$, avec $f_x(v)$ le plus petit suffixe de v préfixe de x , est l'automate minimal du langage Σ^*_x .

Rq 6: l'automate du motif sera utilisé pour l'analyse lexicale; l'automate des occurrences est utile pour la recherche de motif (langage impératif avec label/goto).

def 7: on considère des symboles frais $\epsilon, *, |, +, (,) \notin \Sigma$. Le langage \mathcal{B} des expressions régulières est le plus petit langage contenant $\Sigma \cup \{ \epsilon \}$ tel que:

si $x \in \mathcal{B}, y \in \mathcal{B}$ alors $x^*, (xy)^*, x^+ \in \mathcal{B}$.

def 8: soit $e \in \mathcal{B}$. le langage rationnel $L(e)$ est défini par induction:

- $\forall a \in \Sigma \cup \{ \epsilon \}, L(a) = \{ a \}$
- $\forall x \in \mathcal{B}, L(x^*) = \bigcup_{n \in \mathbb{N}} \{ x^n \}$
- $\forall x \in \mathcal{B}, L(x^+) = \bigcup_{n \in \mathbb{N}^*} \{ x^n \}$
- $\forall x, y \in \mathcal{B}, L(xy) = L(x) \cup L(y)$
- $\forall x, y \in \mathcal{B}, L(x|y) = \{ w, w_2, w_3 \in L(x), w_2 \in L(y) \}$

Rq 9: on omet les parenthèses en cas de non-ambiguïté.

thm 10: (Kleene) Soit $\mathcal{L} \subseteq \Sigma^*$ un langage. on a l'équivalence:

- (i) Il existe un automate \mathcal{A} tq $\mathcal{L} = L(\mathcal{A})$ [\mathcal{L} est reconnaissable par automate].
- (ii) Il existe $e \in \mathcal{B}$ tq $\mathcal{L} = L(e)$ [\mathcal{L} est rationnel].

Rq 11: on dispose d'un algorithme pour passer d'une expression régulière à un automate.

2) Langages algébriques (ALG).

def 12: une grammaire algébrique sur un alphabet Σ est un quadruplet $G = (\mathcal{D}, T, S, R)$ où \mathcal{D} est l'ensemble fini des non-terminaux, $T \subseteq \Sigma$ est l'ensemble des terminaux, $S \in \mathcal{D}$ est l'axiome et $R \subseteq \mathcal{D} \times (\mathcal{D} \cup T)^*$ l'ensemble des règles.

def 13: Soient $u, v \in (\mathcal{D} \cup T)^*$. On dit que u se dérive en v ($u \rightarrow v$) si il existe $\alpha, \beta \in (\mathcal{D} \cup T)^*$ et $X \in \mathcal{D}$ tq $\begin{cases} u = \alpha X \beta \\ v = \alpha w \beta \end{cases}$ et $(X \rightarrow w) \in R$.

def 14: un mot $x \in \Sigma^*$ est engendré par une grammaire $G = (\mathcal{D}, T, S, R)$ si $S \xrightarrow{*} x$, où $\xrightarrow{*}$ est la clôture réflexive transitive de \rightarrow .

On note $L(G)$ le langage des mots engendrés par la grammaire G .

ex 15: le langage $L(a^*b)$ est engendré par la grammaire $\begin{cases} S \rightarrow aS \\ S \rightarrow ab \end{cases}$

ex 16: (langage de Dyck sur n paires de parenthèses). $\begin{cases} S \rightarrow b \\ S \rightarrow aS \end{cases}$

La grammaire $G_{Dyck} : \begin{cases} S \rightarrow a_1 S \bar{a}_1 \\ S \rightarrow a_2 S \bar{a}_2 \\ \vdots \\ S \rightarrow a_n S \bar{a}_n \\ S \rightarrow \epsilon \end{cases}$ engendre le langage des mots bien parenthésés avec au plus n types de parenthèses.

ex 17: (expressions arithmétiques préfixes) $G_{pref} : \begin{cases} S \rightarrow +SS \\ S \rightarrow c \end{cases}$

ex 18: (expressions arithmétiques postfixes) $G_{post} : \begin{cases} S \rightarrow SS+ \\ S \rightarrow c \end{cases}$

thm 19: (RAT \subseteq ALG). Tout langage rationnel est algébrique (ie: engendré par une grammaire algébrique)

def 20: un arbre de dérivation est un arbre fini étiqueté par $\mathcal{D} \cup T \cup \{ \epsilon \}$ tq:

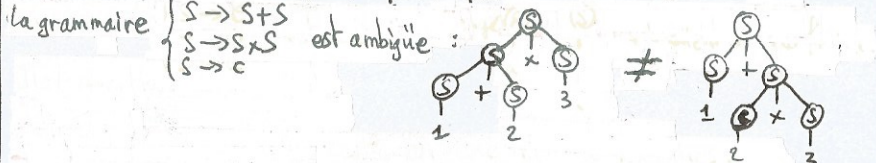
- l'arbre est enraciné en S .
- si T est l'étiquette d'un noeud interne et a_1, \dots, a_n ses fils, alors $(T \rightarrow a_1 \dots a_n) \in R$.

on appellera frontière d'un tel arbre le mot obtenu par concaténation des feuilles de gauche à droite.

def 21: une grammaire est ambiguë si il existe un mot a permettant deux arbres différents.

def 22: un langage algébrique est non-ambiguë si il est engendré par une grammaire non-ambiguë. Dans le cas contraire, il est dit intrinsèquement ambiguë.

ex 23: les grammaires $G_{Dyck}, G_{pref}, G_{post}$ sont non-ambiguës.



II) Analyse lexicale

But: transformer un texte en une liste de lexèmes (terminaux d'une grammaire).

Intérêts: → filtration des commentaires

- mise en page ("1+2" et "1 +2" représentent le même programme)
- repérage des types.
- décoration des lexèmes (nom du fichier, n° de ligne, n° de caractères, ...)
- repérage des erreurs lexicales (motif interdit).

def 24: (règle d'analyse lexicale). Soit T l'ensemble fini des types de lexèmes.

Une règle d'analyse lexicale sur l'alphabet Σ est la donnée:

- d'une expression régulière e sur Σ tq $E \notin L(e)$.
- d'un type de lexème $t \in T$.

Cette règle est notée $e \rightarrow t$.

Algo 25: Entrée: texte $e \in \Sigma^*$; liste de règles $(e_j \rightarrow t_j)_{j \in \{1 \dots n\}}$

Sortie: liste de lexèmes correspondant au texte.

Analyse_lexicale (texte):

si texte = " " retourner []

sinon:

(lexème, suite) ← Analyse_aux (texte)

retourner lexème :: Analyse_lexicale (suite)

Analyse_aux (texte):

fin ← 0 {position de fin du plus long préfixe connu}

i_aut ← -∞ {indice de l'automate du plus long préfixe}

i_texte ← 1

Tant que (il existe un $ct_{motif}^{e_j}$ non-bloqué) et ($i_{\text{texte}} \leq \text{longueur}(\text{texte})$):

pour $j \leftarrow 1$ à n :

Faire lire $\text{texte}[i_{\text{texte}}]$ à l'automate $ct_{motif}^{e_j}$

Si $ct_{motif}^{e_j}$ est dans un état final:

si fin < i_{texte} :

fin ← i_{texte}

i_aut ← j

incrémenter i_{texte}

Si $i_{\text{aut}} = 0$: erreur ("le texte ne contient aucun préfixe connu")

préfixe ← $\text{texte}[1 \dots \text{fin}]$

lexème ← (préfixe, $t_{i_{\text{aut}}}$)

retourner (lexème, $\text{texte} - \text{préfixe}$)

→ Complexité temporelle: $O(\text{longueur}(\text{texte}))$

ex 26: pour analyser le texte a^n avec les règles $a \rightarrow \text{lettre} - a$
 $a^*b \rightarrow a\text{-iterée} - b$
il faut $\frac{n(n+1)}{2}$ lectures.

Rq 27: A la fin de l'analyse lexicale, on dispose d'une table des symboles dans laquelle on a stocké les informations souhaitées (type, lignes ...).

III) Analyse syntaxique

But: déterminer si un mot est reconnu par une grammaire, et construire une dérivation ou l'arbre syntaxique.

1) Algorithmes génériques.

thm 28: Soit $G = (N, T, S, R)$ une grammaire algébrique sur l'alphabet Σ et $m \in \Sigma^*$. Notons a le nombre maximal de symboles d'une partie droite d'une règle dans R . Alors:

$m \in L(G) \iff$ il existe une dérivation reconnaissant m dans G , de longueur inférieure à $a^{|m| \times |N|}$

→ Il suffit de parcourir les dérivation de longueur $\leq a^{|m| \times |N|}$

def 29: une grammaire est en forme normale de Chomsky si ses règles sont de la forme " $N \rightarrow HL$ " ou " $N \rightarrow a$ ", $N, H, L \in N$, $a \in T$.

prop 30: pour toute grammaire G , il existe une grammaire G' en forme normale de Chomsky tq $L(G') = L(G) \setminus \{\epsilon\}$ (détaille $|G'| = O(|G|^2)$)

algo 31: (CYK) Entrée: une grammaire G en forme normale de Chomsky, un mot $w \in \Sigma^*$.

Sortie: oui si $w \in L(G)$, non sinon

Algorithme, correction, terminaison, complexité ϵ .

DEVI

2) Algorithmes linéaires.

a) Analyse descendante.

Principe 32: on part de l'axiome et on le dérive pour arriver au texte.

ex 33: un algorithme simple permet d'analyser la grammaire G_{pref} [cf Annexe 1]

La lecture du caractère courant permet de choisir la règle à appliquer.

⚠ ce n'est pas le cas pour toutes les grammaires: $G_{\text{pref}}^2 \begin{cases} S \rightarrow TSS \\ T \rightarrow + \\ S \rightarrow c \end{cases} \quad G_{\text{pref}}^3 \begin{cases} S \rightarrow U+SS \\ U \rightarrow E \\ S \rightarrow c \end{cases}$

def 34 : ($\text{premier}(\vec{a})$)
 Soit $G = (\mathcal{N}, T, S, R)$ une grammaire. Pour $\vec{a} \in (\mathcal{N} \cup T)^*$, on définit:

$$\text{premier}(\vec{a}) = \begin{cases} \{a \in T \mid \vec{a} \xrightarrow{*} a\vec{\beta}, \vec{\beta} \in (T \cup \mathcal{N})^*\} & \text{si } \vec{a} \not\xrightarrow{*} \epsilon \\ \{a \in T \mid \vec{a} \xrightarrow{*} a\vec{\beta}, \vec{\beta} \in (T \cup \mathcal{N})^*\} \cup \{\Delta\} & \text{si } \vec{a} \xrightarrow{*} \epsilon \end{cases}$$
 où $\Delta \notin \Sigma$ est un symbole frais.

def 35 : ($\text{suivant}(N)$)
 Soit $G = (\mathcal{N}, T, S, R)$ une grammaire. Pour $N \in \mathcal{N}$, on définit:

$$\text{suivant}(N) = \begin{cases} \{a \in T \mid S \xrightarrow{*} \vec{a}N\vec{\beta}, \vec{a}, \vec{\beta} \in (T \cup \mathcal{N})^*\} & \text{si } S \not\xrightarrow{*} \vec{a}N \\ \{a \in T \mid S \xrightarrow{*} \vec{a}N\vec{\beta}, \vec{a}, \vec{\beta} \in (T \cup \mathcal{N})^*\} \cup \{\Delta\} & \text{si } S \xrightarrow{*} \vec{a}N \end{cases}$$
 où $\Delta \notin \Sigma \cup \{\mathcal{N}\}$ est un symbole frais.

thm 36 : (calcul de premier)
 La fonction premier est la plus petite partie de $(T \cup \mathcal{N})^* \times (T \cup \{\Delta\})^*$ tq
 (i) $\forall a \in T, \text{premier}(a) = \{a\}$
 (ii) $\text{premier}(\epsilon) = \{\Delta\}$
 (iii) $\forall (N \rightarrow \epsilon) \in S, \text{premier}(\epsilon) \subseteq \text{premier}(N)$
 (iv) $\forall (N \rightarrow \alpha_1 \dots \alpha_n) \in S, \text{premier}(\alpha_1 \dots \alpha_n) \subseteq \text{premier}(N)$
 (v) Si $\Delta \notin \text{premier}(\alpha_i), \text{premier}(\alpha_i) \subseteq \text{premier}(\alpha_1 \dots \alpha_n)$
 (vi) Si $\Delta \in \text{premier}(\alpha_i), \text{premier}(\alpha_i) \setminus \{\Delta\} \cup \text{premier}(\alpha_2 \dots \alpha_n) \subseteq \text{premier}(\alpha_1 \dots \alpha_n)$

DEV 2

def 37 : une grammaire G est dite LL(1) si pour tout $N \in \mathcal{N}$, en notant $N \rightarrow \vec{\alpha}_1, \dots, N \rightarrow \vec{\alpha}_n$ les règles dont la partie gauche est N , on a:
 (i) les ensembles $\text{premier}(\vec{\alpha}_1), \dots, \text{premier}(\vec{\alpha}_n)$ sont deux à deux disjoints.
 (ii) Si $N \xrightarrow{*} \epsilon$, alors $\text{suivant}(N)$ est disjoint de chaque $\text{premier}(\vec{\alpha}_i)$

ex 38 : $G_{\text{pref}}, G_{\text{pref}}^2, G_{\text{Dyck}} \in \text{LL}(1), G_{\text{post}} \notin \text{LL}(1)$

def 39 : Un langage est LL(1) s'il existe une grammaire LL(1) qui le reconnaît.

ex 40 : $L(G_{\text{post}}) \in \text{LL}(1)$ [cf Annexe 2]

thm 41 : $\text{RAT} \subseteq \text{LL}(1) \subseteq \{\text{langages non-ambigus}\}$

Analyse-LL(1) (texte) : arbre
 $A, \text{texte}' = \text{Analyse-LL}(1)\text{-aux}(S, \text{texte})$
 si $\text{texte}' \neq \epsilon$: erreur "le mot n'est pas reconnu"
 retourner (A)

Analyse-LL(1) - aux(α, texte) : arbre x texte

si $\alpha \in T$:
 | si $\text{texte}[1] = \alpha$: retourner ($\boxed{\text{texte}[1]}$, $\text{texte}[2 \dots]$)
 | sinon : erreur ("le mot n'est pas reconnu")
 # les règles sont de la forme $\alpha \rightarrow \vec{\beta}_1, \dots, \alpha \rightarrow \vec{\beta}_n$
 Si $\text{texte} = \epsilon$:
 | si $\Delta \in \text{premier}(\alpha)$ et $\Delta \in \text{suivant}(\alpha)$:
 | | $i \leftarrow$ l'indice tq $\Delta \in \text{premier}(\vec{\beta}_i)$
 | | sinon : erreur ("le mot n'est pas reconnu")
 Sinon :
 | si ($\Delta \in \text{premier}(\alpha)$) et ($\text{texte}[1] \in \text{suivant}(\alpha)$) :
 | | $i \leftarrow$ l'indice tq $\Delta \in \text{premier}(\vec{\beta}_i)$
 | | sinon $i \leftarrow$ l'indice tq $\text{texte}[1] \in \text{premier}(\vec{\beta}_i)$
 si $\vec{\beta}_i = \epsilon$: retourner ($\begin{matrix} \textcircled{\alpha} \\ \boxed{\epsilon} \end{matrix}$, texte)
 Sinon : $A_1, \text{texte}'_1 \leftarrow \text{Analyse-LL}(1)\text{-aux}(\vec{\beta}_i[1], \text{texte})$
 $A_2, \text{texte}'_2 \leftarrow \text{Analyse-LL}(1)\text{-aux}(\vec{\beta}_i[2], \text{texte}'_1)$
 \vdots
 Retourner ($\begin{matrix} \alpha \\ \circlearrowleft \\ A_1 \quad A_2 \quad A_3 \quad \dots \quad A_n \end{matrix}$, texte'_n)

b/Analyse ascendante

Principe : on part du texte et on remonte les règles pour obtenir l'axiome.

def 3 : une situation LR(0) d'une grammaire G , notée $N \rightarrow \alpha_1 \dots \alpha_j \alpha_{j+1} \dots \alpha_n$ est :

- une règle $(N \rightarrow \alpha_1 \dots \alpha_n) \in R$ avec $N \in \mathcal{N}, \alpha_1 \dots \alpha_n \in \mathcal{N} \cup T$.
- un entier $j \in \{0, n\}$. \rightarrow on note \mathcal{S}_G l'ensemble des situations.

def 4 : la fermeture de $E \in \mathcal{S}_G$ est la plus petite partie $\text{Cl}_G(E)$ de \mathcal{S}_G tq :

- $E \subseteq \text{Cl}_G(E)$
- $(N \rightarrow \alpha_1 \uparrow \beta) \in \mathcal{S}_G(E)$ et $(L \rightarrow \gamma) \in R \Rightarrow (L \rightarrow \gamma \beta) \in \text{Cl}_G(E)$.

def 45 : l'automate LR(0) d'une grammaire G est défini par :

- $Q = \{\text{Cl}_G(E), E \in \mathcal{S}_G\}$
- $I = \{\text{Cl}_G(\{S \rightarrow \vec{a}\}) \mid (S \rightarrow \vec{a}) \in R\}$
- $\delta : (E, \alpha) \mapsto \text{Cl}_G(\{N \rightarrow \vec{\gamma} \alpha \uparrow \vec{\beta} \mid N \rightarrow \vec{\gamma} \alpha \vec{\beta} \in E\})$

ex 46 : Automate LR(0) de la grammaire G_{post} [cf Annexe 3]

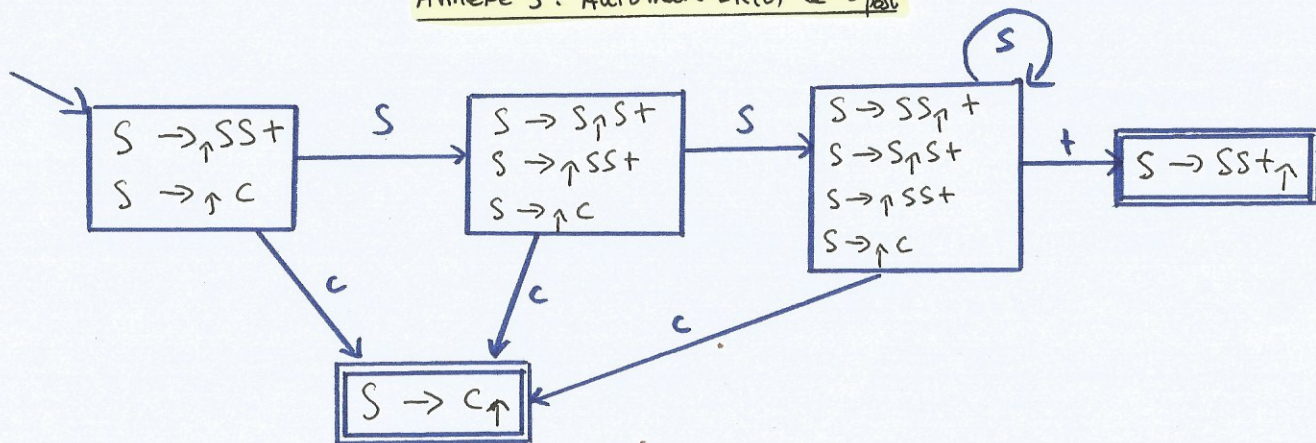
def 47 : une grammaire G est dite LR(0) si tout état de son automate qui contient une situation de la forme $N \rightarrow \alpha_1 \dots \alpha_n \uparrow$ ne contient que celle-ci

prop 48 : schéma des inclusions $\text{RAT}, \text{LL}(1), \text{LR}(0), \text{ALG}$ [cf Annexe 4]

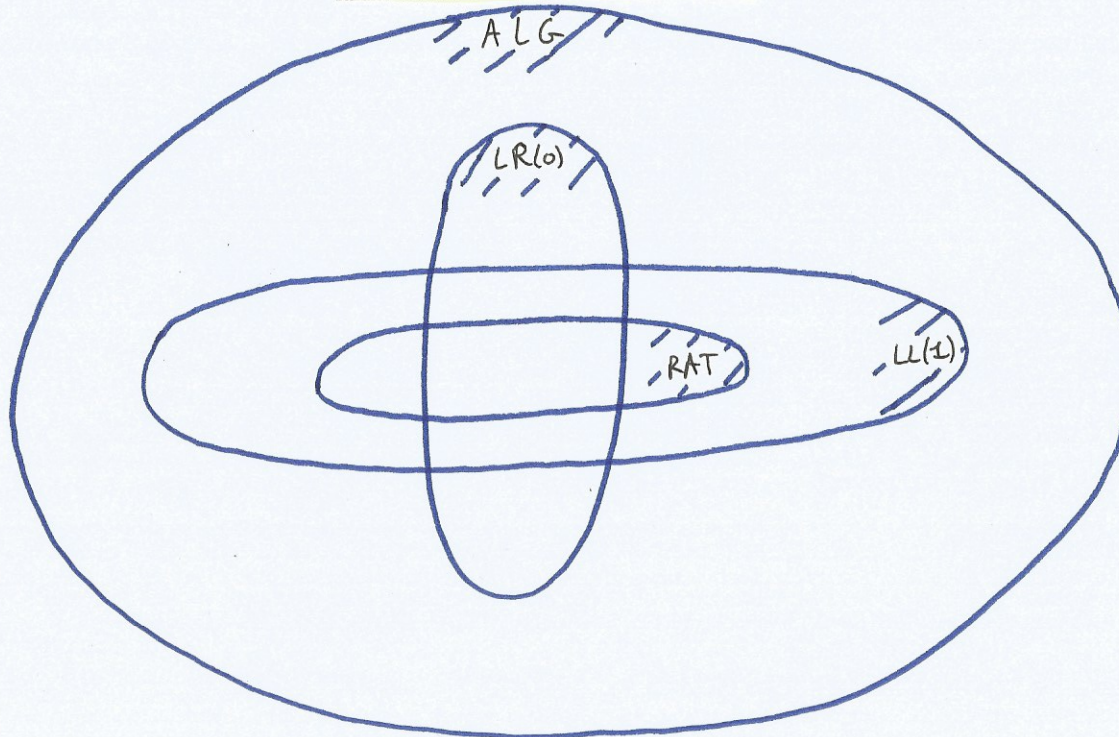
Annexe 1 : Analyse G_{pref}

mot lu	arbre en cours de construction
↑****	
+↑****	
++↑****	
++↑1****	
++↑12*	
++↑123↑	

Annexe 3 : Automate LR(0) de G_{post}



Annexe 4 : schéma des inclusions -



Annexe 2 : une grammaire LL(1) reconnaissant L(G_{post})

$$G_{\text{post}}^{LL(1)} : \begin{cases} S \rightarrow cT \\ T \rightarrow S+T \\ T \rightarrow \epsilon \end{cases}$$