

Dérivation formelle de l'algorithme d'analyse
syntaxique d'Earley par abstraction d'une
sémantique des grammaires algébriques

Patrick Cousot
ENS

Journée de présentation des cursus en informatique
ENS Cachan
15 mai 2003

<http://www.lsv-cachan.fr/dptinfo/15mai.php>

Grammaires algébriques

— Exemple :

Grammaire : $A \rightarrow AA \mid a$

formellement $\langle \{A\}, \{a\}, \{(A, AA), (A, a)\}, A \rangle$

— Grammaire

$\langle N, T, R, A \rangle$

$\begin{array}{l} \downarrow \text{axiome } A \in N \\ \downarrow \text{règles } R \subseteq N \times (N \cup T)^* \\ \downarrow \text{terminaux} \\ \downarrow \text{non-terminaux } N \cap T = \emptyset \end{array}$

— Chaînes :

X^* chaînes finies sur l'alphabet X (ϵ chaîne vide)

Points fixes

- $\langle P, \leq \rangle$, treillis complet (ensemble partiellement ordonné tel que tout sous-ensemble $X \subseteq P$ a une borne supérieure $\vee X$)
- $f \in P \rightarrow P$, croissante
- $\text{lfp } f \triangleq \bigwedge \{x \mid f(x) \leq x\}$ est le plus petit point fixe de f — Tarski
- si f préserve les bornes sup. existantes ($f(\vee X) = \vee f(X)$)

$$\text{lfp } f = \bigvee_{n \in \mathbb{N}} f^n(\perp)$$

- $\perp \triangleq \vee \emptyset$ est l'infimum du treillis complet
- $f^0(x) = x$
- $f^{n+1}(x) = f \circ f^n(x)$, $n \in \mathbb{N}$

Exemple : sémantique de Schützenberger d'une grammaire algébrique (langage terminal engendré pour chaque non-terminal)

Exemple : $A \rightarrow AA \mid a$

$$S_G(X) = \{(A, a)\} \cup \{(A, \alpha_1 \alpha_n) \mid (A, \alpha_1) \in X \wedge (A, \alpha_n) \in X\}$$

$$X^0 = \emptyset$$

-- itérés

$$X^1 = S_G(X^0) = \{(A, a)\}$$

$$\dots$$
$$X^n = \{(A, a^k) \mid 1 \leq k \leq 2^{n-1}\}$$

-- hypothèse d'induction

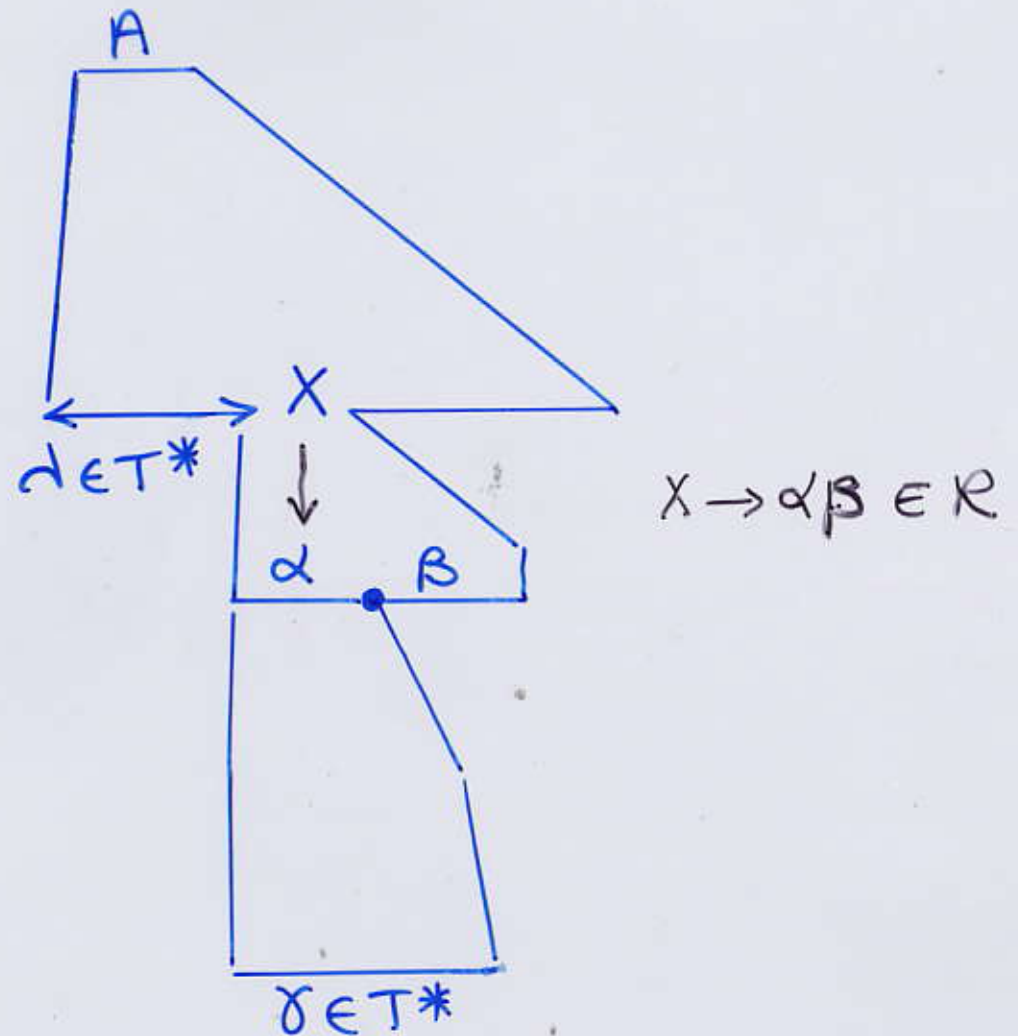
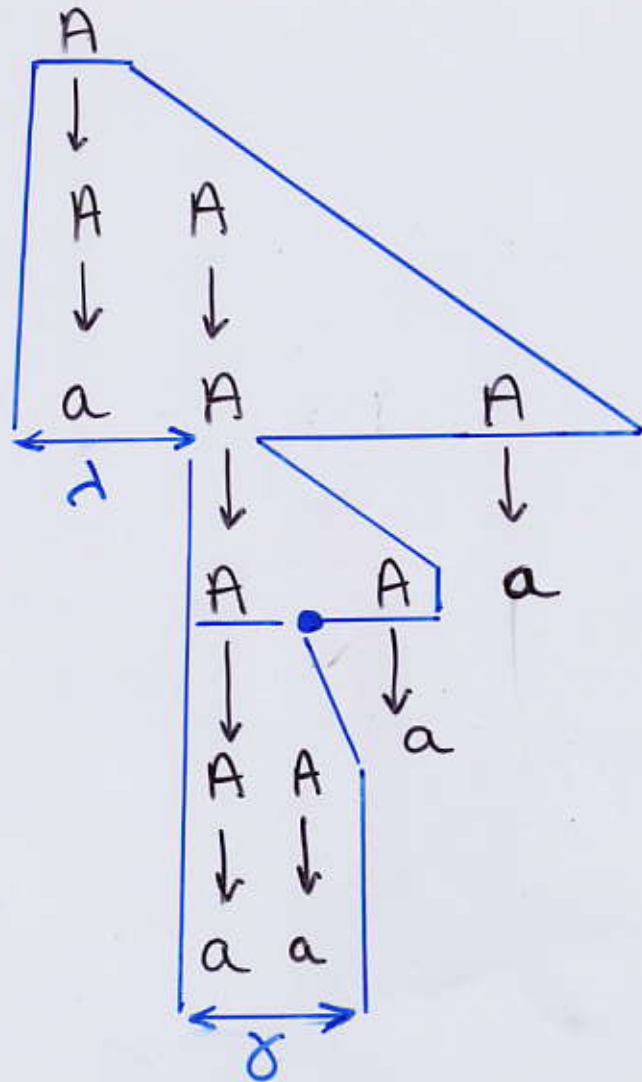
$$X^{n+1} = S_G(X^n)$$

$$= \{(A, a)\} \cup \{(A, a^{k_1+k_2}) \mid 1 \leq k_1+k_2 \leq 2^{n-1}\}$$

$$= \{(A, a^k) \mid 1 \leq k \leq 2^n\}$$

$$\text{lfp } S_G = \bigcup_{n \geq 0} X^n = \{(A, a^k) \mid k \geq 1\}$$

Dérivation : exemple & notation



$[d, X \rightarrow \alpha \cdot \beta, \gamma]$ (item)

Sémantique de dérivation d'une grammaire

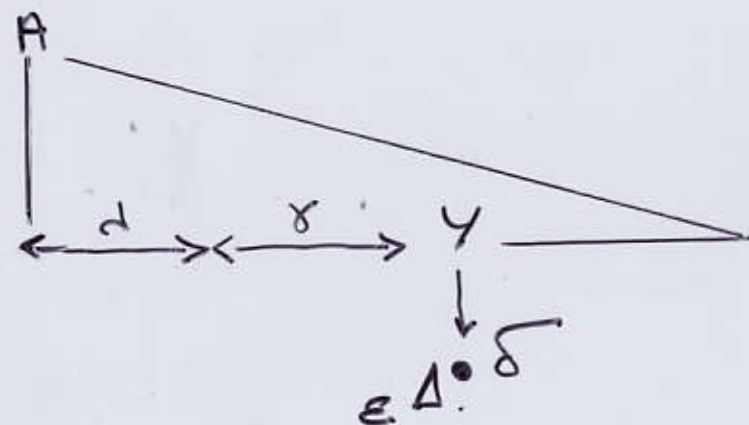
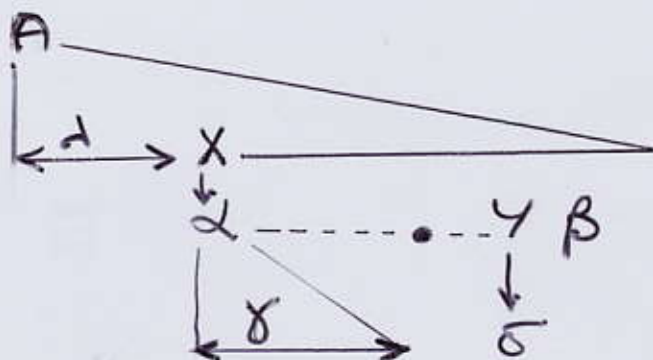
- Schéma d'axiome (initialisation)

$$[\varepsilon, A \rightarrow \cdot \beta, \varepsilon] \quad \text{ssi} \quad A \rightarrow \beta \in R$$

- Schémas de règles :

• Dérivation ($X \rightarrow \alpha \gamma \beta, \gamma \rightarrow \delta \in R$)

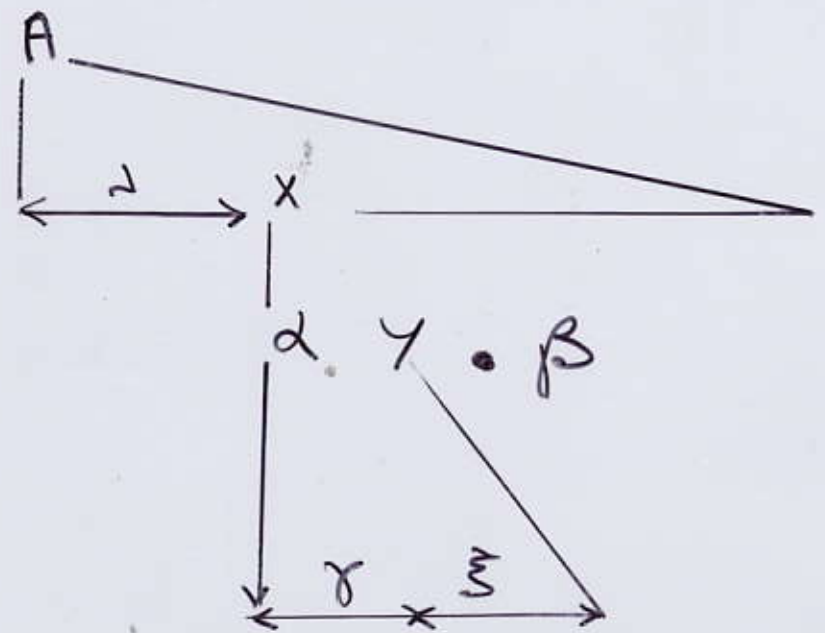
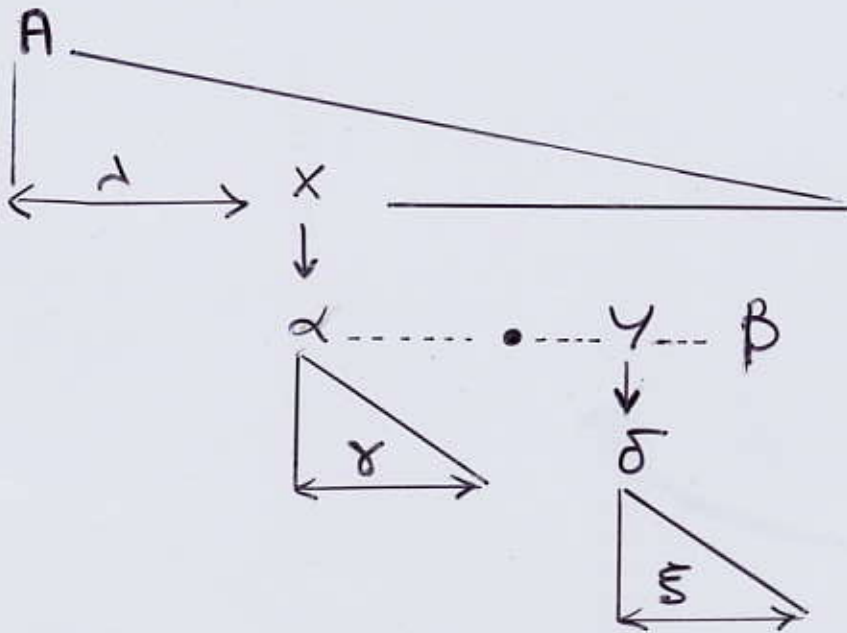
$$\frac{[\lambda, X \rightarrow \alpha \cdot \gamma \beta, \gamma]}{[\lambda \gamma, \gamma \rightarrow \cdot \delta, \varepsilon]}$$



- Réduction $(X \rightarrow \alpha \gamma \beta, \gamma \rightarrow \delta \in R) :$

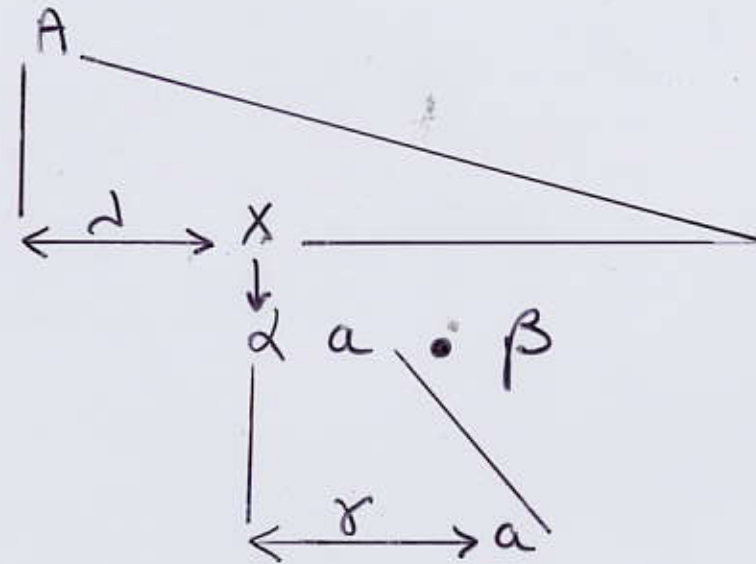
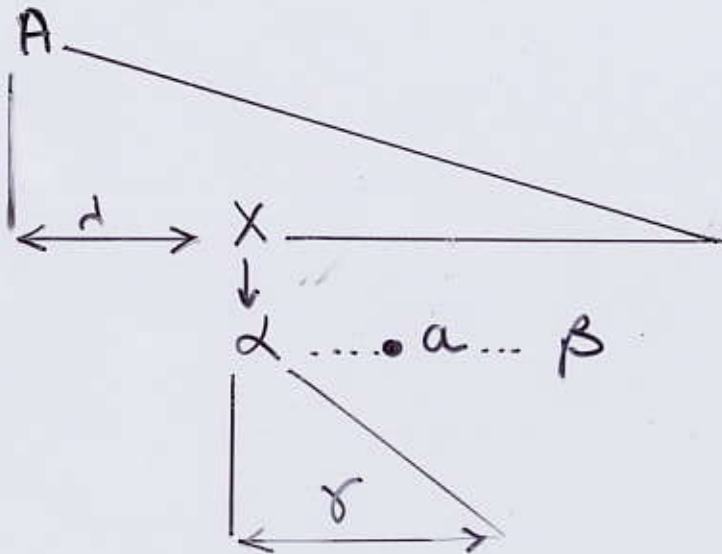
$$[\Delta, X \rightarrow \alpha \cdot \gamma \beta, \delta] \quad [\Delta \delta, \gamma \rightarrow \delta \cdot, \xi]$$

$$[\Delta, X \rightarrow \alpha \gamma \cdot \beta, \delta \xi]$$



- Avancement ($X \rightarrow \alpha a \beta \in \mathcal{R}, a \in T$):

$$\frac{[\lambda, X \rightarrow \alpha \cdot a \beta, \gamma]}{[\lambda, X \rightarrow \alpha a \cdot \beta, \gamma] \text{ (I)}}$$



Ensemble défini par des règles :

- \mathcal{U} univers

- Règles $R = \left\{ \frac{P_i}{C_i} \mid i \in \Delta \right\}$ Prémisses $P_i \subseteq \mathcal{U}$ (\emptyset pour axiome)
Conclusion $C_i \subseteq \mathcal{U}$

- $F_R(X) = \{ C \mid \exists \frac{P}{C} \in R : P \subseteq X \}$

(conclusions que l'on peut tirer de X par R)

- F_R est croissante sur le treillis complet $(\mathcal{F}(\mathcal{U}), \subseteq)$

- $\text{lfp } F_R$ existe (Tarski)

$$= \bigcup_{n \geq 0} F_R^n(\emptyset) \quad f^0(x) = x \quad f^{n+1}(x) = f(F^n(x))$$

(conclusions successives que l'on peut tirer de R sans hypothèse initiale).

Sémantique de dérivation d'une grammaire algébrique

$$G = \langle N, T, R, A \rangle$$

lfp F_G où

$$F_G(X) = \{ [\varepsilon, A \rightarrow \cdot \beta, \varepsilon] \mid A \rightarrow \beta \in R \}$$

$$\cup \{ [\alpha \gamma, Y \rightarrow \cdot \delta, \varepsilon] \mid [\alpha, X \rightarrow \alpha \cdot Y \beta, \gamma] \in X \wedge Y \rightarrow \delta \in R \}$$

$$\cup \{ [\alpha, X \rightarrow \alpha Y \cdot \beta, \gamma \delta] \mid [\alpha, X \rightarrow \alpha \cdot Y \beta, \gamma] \in X \wedge [\alpha \gamma, Y \rightarrow \delta \cdot, \delta] \in X \}$$

$$\cup \{ [\alpha, X \rightarrow \alpha a \cdot \beta, \gamma a] \mid [\alpha, X \rightarrow \alpha \cdot a \beta, \gamma] \in X \}$$

Abstraction par une correspondance de Galois

- (P, \leq) ensemble partiellement ordonné
- (Q, \sqsubseteq) ensemble partiellement ordonné
- $\alpha \in P \rightarrow Q$, $\gamma \in Q \rightarrow P$ tels que:
 $\forall x \in P: \forall y \in Q: \alpha(x) \sqsubseteq y \iff x \leq \gamma(y)$

$\iff ((P, \leq), (Q, \sqsubseteq), \alpha, \gamma)$ est une correspondance de Galois

Exemple :

$$f: A \rightarrow B$$

$$\alpha \in \mathcal{F}(A) \rightarrow \mathcal{F}(B)$$

$$\alpha(X) = \{f(x) \mid x \in X\}$$

$$\gamma \in \mathcal{F}(B) \rightarrow \mathcal{F}(A)$$

$$\gamma(Y) = \{x \in A \mid f(x) \in Y\}$$

Abstraction de point fixe par une correspondance de Galois

- $\langle P, \leq \rangle$ et $\langle Q, \sqsubseteq \rangle$, treillis complets
- $\langle P, \leq \rangle \xrightleftharpoons[\alpha]{\gamma} \langle Q, \sqsubseteq \rangle$ correspondance de Galois
- $f: P \rightarrow P$, croissante
- $g: Q \rightarrow Q$, croissante
- $\alpha \circ f = g \circ \alpha$

$$\Rightarrow \alpha(\text{lfp } f) = \text{lfp } g$$

preuve

$$- \alpha(\text{lfp } f)$$

$$= \alpha(\bigvee \{x \mid f(x) \leq x\}) \quad \text{-- Tarshi}$$

$$= \bigwedge \{ \alpha(x) \mid f(x) \leq x \} \quad \text{-- } \alpha \text{ préserve les bornes sup. existantes}$$

$$f(x) \leq x \Rightarrow \alpha(f(x)) \sqsubseteq \alpha(x) \Rightarrow g(\alpha(x)) \sqsubseteq \alpha(x) \Rightarrow \alpha(x) \in \{y \mid g(y) \sqsubseteq y\}$$

$$\sqsubseteq \bigwedge \{y \mid g(y) \sqsubseteq y\}$$

$$= \text{lfp } g$$

-- Tarshi

$$- g(\alpha(\text{lfp } f))$$

$$= \alpha(f(\text{lfp } f))$$

$$= \alpha(\text{lfp } f)$$

$$\Rightarrow \text{lfp } g \sqsubseteq \alpha(\text{lfp } f)$$

$$\text{-- } \alpha \circ f = g \circ \alpha$$

-- $\text{lfp } f$ est un point fixe (le plus petit)

-- $\text{lfp } g$ est le plus petit point fixe de g

$$- \text{lfp } g = \alpha(\text{lfp } f)$$

-- antisymétrie

□

Abstraction de la sémantique de dérivation de
 $G = \langle N, T, R, A \rangle$ en la sémantique de Schützenberger
 (langage fini engendré pour chaque non-terminal)

$$\alpha_s(I) = \{ \langle X, \gamma \rangle \mid \exists \alpha \in T^* : [\alpha, X \rightarrow \alpha., \gamma] \in I \}$$

$$\alpha_s(F_G(I)) =$$

$$\dots \quad \text{-- 2 pages de calcul}$$

$$S_G(\alpha_s(I))$$

$$S_G(L) = \{ (X, \alpha) \mid X \rightarrow \alpha \in R \wedge \alpha \in T^* \}$$

$$\cup \{ (X, \alpha_1 \dots \alpha_n) \mid X \rightarrow X_1 \dots X_n \in R \wedge \forall i \in [1, n] : \\ (X_i = \alpha_i \in T) \vee ((X, \alpha_i) \in L) \}$$

$$\alpha_s(\text{lfp } F_G) = \text{lfp } S_G \quad \text{sémantique de Schützenberger}$$

Analyse syntaxique

Etant donnée une grammaire $G = (N, T, R, A)$ et une phrase terminale $\sigma \in T^*$, répondre à la question :

– Est-ce que σ est engendrée par G ?

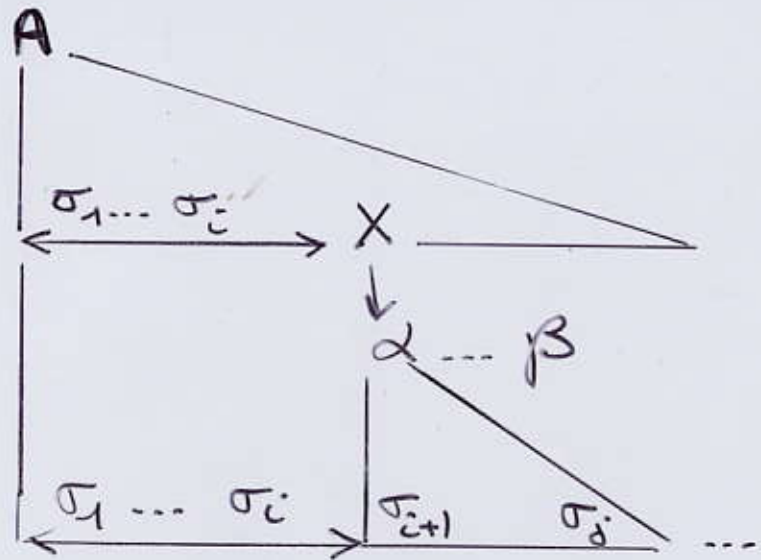
– i.e. $\sigma \in \text{LPP } S_G$ (sémantique de Schützenberger)

(et déterminer sa structure syntaxique)

Abstraction de la sémantique de dérivation en l'algorithme d'analyse syntaxique d'Earley.

- $\sigma = \sigma_1 \dots \sigma_n \in T^*$, phrase à analyser

- $\alpha_E(I) = \{ (X \rightarrow \alpha \cdot \beta, i, j) \mid 0 \leq i \leq j \leq n \wedge [\sigma_1 \dots \sigma_i, X \rightarrow \alpha \cdot \beta, \sigma_{i+1} \dots \sigma_j] \in I \}$



Dérivation de l'algorithme de Earley

$$\alpha_E(F_G(I)) =$$

$$\overline{\overline{E_G(\alpha_E(I))}} \quad \text{--- 2 pages de calculs}$$

$$\begin{aligned} E_G(I) = & \{ (A \rightarrow \bullet \gamma, 0, 0) \mid A \rightarrow \gamma \in R \} \\ & \cup \{ (\gamma \rightarrow \bullet \delta, i, j) \mid (x \rightarrow \alpha \bullet \gamma \beta, i, j) \in I \wedge \gamma \rightarrow \delta \in R \} \\ & \cup \{ (x \rightarrow \alpha \gamma \bullet \beta, k, j) \mid (x \rightarrow \alpha \bullet \gamma \beta, k, i) \in I \wedge \\ & \quad (\gamma \rightarrow \delta \bullet, i, j) \in I \} \\ & \cup \{ (x \rightarrow \alpha \sigma_j \bullet \beta, i, j) \mid (x \rightarrow \alpha \bullet \sigma_j \beta, i, j-1) \in I \} \end{aligned}$$

$\text{lfp } E_G$ est fini, donc calculable itérativement.

$\langle A \rightarrow \delta \bullet, 0, n \rangle \in \text{lfp } E_G$ est la réponse.

Correction de l'algorithme d'Earley

$$\begin{aligned} & (A \rightarrow \gamma. , 0, n) \in \text{eff} E_G \\ \Leftrightarrow & (A \rightarrow \gamma. , 0, n) \in \alpha_E(\text{eff} F_G) \\ \Leftrightarrow & (\varepsilon, A \rightarrow \gamma. , \sigma) \in \text{eff} F_G \\ \Leftrightarrow & (A, \sigma) \in \alpha_S(\text{eff} F_G) \\ \Leftrightarrow & (A, \sigma) \in \text{eff} S_G \quad \text{-- sémantique de Schützenberger} \end{aligned}$$

i.e. la phrase σ est engendrée par la grammaire G .

Pour connaître tous les détails :

P. Cousot & R. Cousot

Parsing as abstract interpretation of grammars
Theoret. Comput. Sci. 290 : 531 — 544, 2003

<http://www.di.ens.fr/~cousot/papers/TCS03-parsing.shtml>

— MERCI DE VOTRE ATTENTION —