

# Calculs d'unification sur les arbres de dérivation TAG

Sylvain Schmitz, Joseph Le Roux

► **To cite this version:**

Sylvain Schmitz, Joseph Le Roux. Calculs d'unification sur les arbres de dérivation TAG. TALN'08, 15ème Conférence sur le Traitement Automatique des Langues Naturelles, Jun 2008, Avignon, France. p. 320–329, 2008. <inria-00270922>

**HAL Id: inria-00270922**

**<https://hal.inria.fr/inria-00270922>**

Submitted on 7 Apr 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Calculs d'unification sur les arbres de dérivation TAG

Sylvain SCHMITZ   Joseph LE ROUX

LORIA, INRIA Nancy - Grand Est, Nancy, France

LORIA, Université Nancy 2, Nancy, France

Sylvain.Schmitz@loria.fr, Joseph.LeRoux@loria.fr

## Abstract

Nous définissons un formalisme, les grammaires rationnelles d'arbres avec traits, et une traduction des grammaires d'arbres adjoints avec traits vers ce nouveau formalisme. Cette traduction préserve les structures de dérivation de la grammaire d'origine en tenant compte de l'unification de traits. La construction peut être appliquée aux réalisateurs de surface qui se fondent sur les arbres de dérivation.

The derivation trees of a tree adjoining grammar provide a first insight into the sentence semantics, and are thus prime targets for generation systems. We define a formalism, feature based regular tree grammars, and a translation from feature based tree adjoining grammars into this new formalism. The translation preserves the derivation structures of the original grammar, and accounts for feature unification.

*Mots clefs* : Unification, grammaire d'arbres adjoints, arbre de dérivation, grammaire rationnelle d'arbres

Unification, tree adjoining grammar, derivation tree, regular tree grammar

## 1 Introduction

Le processus de dérivation dans les grammaires d'arbres adjoints (JOSHI et SCHABES, 1997, TAG) produit deux arbres : l'*arbre dérivé* qui correspond à un arbre syntagmatique classique (voir figure 1b), et l'*arbre de dérivation*, qui présente par quelles opérations les arbres élémentaires de la grammaire ont été combinés pour obtenir l'arbre dérivé (voir figure 1a). Selon la tâche de traitement de la langue, il sera plus adéquat de considérer l'un ou l'autre, l'arbre dérivé étant en correspondance avec les lexèmes d'une phrase, tandis que l'arbre de dérivation donne une vue sémantique primitive de la phrase, comme le montrent par exemple CANDITO et KAHANE (1998).

De fait, l'arbre de dérivation est privilégié dans plusieurs approches pour la réalisation de surface (KOLLER et STRIEGNITZ, 2002; KOLLER et STONE, 2007). Il sert aussi de pivot à partir duquel représentation sémantique et arbre dérivé peuvent être générés dans les approches de DE GROOTE (2002), POGODALLA (2004) et KANAZAWA (2007) à base de grammaires catégorielles abstraites.

Ces travaux ne sont cependant pas immédiatement applicables à des grammaires réalistes qui emploient une variante des TAG à base de structures de

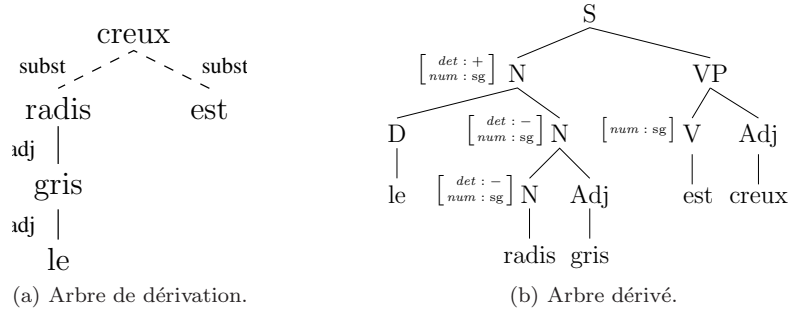


Figure 1: Dérivation de la phrase « Le radis gris est creux. » avec la grammaire de la figure 2.

traits (VIJAY-SHANKER, 1992, voir par exemple la figure 2). Cette variante munit les nœuds des arbres élémentaires de structures de traits, dont les unifications contraignent les opérations de substitution ou d’adjonction du nœud. Ces structures ne posent en théorie aucun problème, car les domaines de valeur des différents traits sont finis et il suffit de démultiplier le nombre de symboles non-terminaux pour émuler les différentes structures possibles. Mais le nombre de ces structures s’accommode mal de cette vision naïve : par exemple, les vingt-huit traits syntaxiques utilisés dans la grammaire SEMFRAG du français (GARDENT, 2006) décrivent un domaine, certes fini, mais comprenant plus de 214 milliards d’éléments. Enfin, l’argument du domaine fini ne tient tout simplement pas pour certains mécanismes de construction sémantique fondés sur l’unification de traits d’index sémantiques qui ont des domaines de valeur non finis (GARDENT et KALLMEYER, 2003; GARDENT, 2006).

Nous étudions dans cet article la traduction d’une grammaire d’arbres adjoints avec structures de traits en une grammaire rationnelle d’arbres de dérivation qui en préserve les mécanismes d’unification de traits. Plus en détail,

- nous rappelons comment traduire une grammaire TAG en une grammaire rationnelle d’arbres (RTG) qui en génère les arbres de dérivation (section 2.1),
- puis nous définissons un formalisme de grammaires rationnelles d’arbres enrichies par des structures de traits et montrons comment traduire une grammaire TAG dans ce nouveau formalisme (section 2.2) ;
- enfin, nous proposons une seconde traduction qui améliore l’efficacité de la génération des arbres de dérivation TAG (section 3).

Nous supposons que le lecteur est familier avec les aspects théoriques des grammaires d’arbres adjoints (JOSHI et SCHABES, 1997), des grammaires rationnelles d’arbres (COMON et al., 2007) et de l’unification (ROBINSON, 1965)<sup>1</sup>.

<sup>1</sup>Pour éviter toute confusion avec l’opération de substitution dans les TAG, la notion de substitution que l’on trouve associée à l’unification sera appelée u-substitution dans la suite.

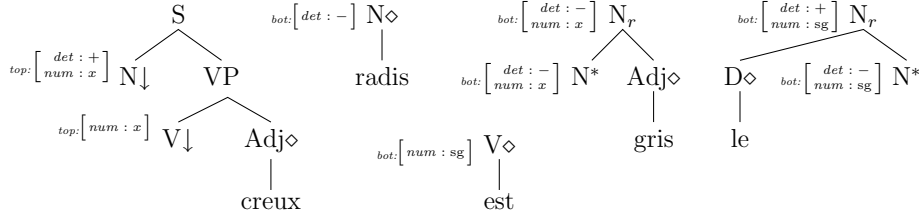


Figure 2: Exemple de grammaire d'arbres adjoints avec structures de traits.

## 2 Arbres de dérivation et unification

Un arbre de dérivation d'une grammaire d'arbres adjoints a des nœuds étiquetés par des arbres élémentaires de la grammaire et en guise d'arêtes les relations d'adjonctions et substitutions permises par la grammaire entre arbres élémentaires. Dans un premier temps, nous reformulons la description donnée par DE GROOTE (2002) des arbres de dérivation qu'une grammaire d'arbres adjoints peut engendrer, en utilisant explicitement une grammaire rationnelle d'arbres. Dans un second temps, nous montrons comment les calculs d'unification de l'arbre dérivé peuvent s'intégrer simplement dans cette grammaire rationnelle.

### 2.1 Grammaire rationnelle des arbres de dérivation

Formellement, une grammaire d'arbres adjoints  $\langle \Sigma, N, I, A, S \rangle$  est constituée d'un alphabet terminal  $\Sigma$ , d'un alphabet non-terminal  $N$ , d'un ensemble  $I$  d'arbres initiaux  $\alpha$ , d'un ensemble  $A$  d'arbres auxiliaires  $\beta$ , et d'un non-terminal distingué  $S$  de  $N$ . Nous désignons par  $\gamma_r$  le nœud racine de l'arbre élémentaire  $\gamma$  et par  $\beta_f$  le nœud pied de l'arbre auxiliaire  $\beta$ .

Les nœuds d'un arbre élémentaire  $\gamma$  de  $I \cup A$  qui nous intéressent sont étiquetés par des non-terminaux, et permettent une opération de substitution ou d'adjonction ; nous considérons en particulier que le pied d'un arbre auxiliaire ne permet pas d'adjonction<sup>2</sup>. Nous numérotons ces nœuds par un parcours arbitraire depuis la racine, de sorte que  $\gamma_1 = \gamma_r$ . Nous notons  $\mathbf{lab}(\gamma_i)$  l'étiquette dans  $N$  du nœud  $\gamma_i$ .

Pour construire la grammaire rationnelle  $\langle S, N \cup N_A, \mathcal{F}, R \rangle$  des arbres de dérivation, nous définissons :

- l'ensemble des arbres élémentaires comme notre alphabet ordonné  $\mathcal{F} = I \cup A \cup \{\varepsilon\}$ , où le rang  $n = \mathbf{rg}(\gamma)$  d'un arbre élémentaire  $\gamma$  est le nombre de ses nœuds où une substitution ou une adjonction est possible, et où  $\varepsilon$ , de rang 0, représente une feuille vide ;
- l'alphabet non-terminal  $N$  et un duplicata  $N_A = \{X_A \mid X \in N\}$  comme alphabet de la grammaire rationnelle ; à chaque nœud non terminal  $\gamma_i$  d'un arbre étiqueté par  $X = \mathbf{lab}(\gamma_i)$ , on associe un non terminal  $\mathbf{nt}(\gamma_i)$  de forme  $X$  dans  $N$  s'il permet une substitution ou  $X_A$  dans  $N_A$  s'il permet une adjonction ;

<sup>2</sup>Dans un souci de concision, nous ne traitons pas les contraintes d'adjonction sélective, qui ne posent aucune difficulté conceptuelle.

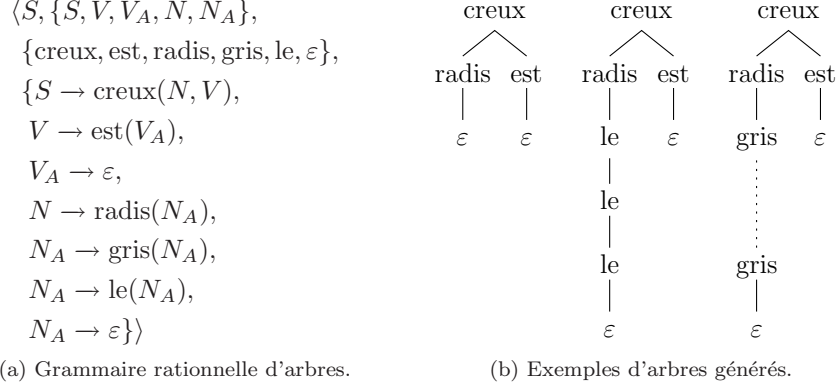


Figure 3: Grammaire rationnelle correspondant à la grammaire TAG de la figure 2.

- l'ensemble de règles  $R$  défini comme l'union

$$\begin{aligned}
 & \{X \rightarrow \alpha(\text{nt}(\alpha_1), \dots, \text{nt}(\alpha_n)) \mid \alpha \in I, n = \text{rg}(\alpha), X = \text{lab}(\alpha_r)\} \\
 \cup & \{X_A \rightarrow \beta(\text{nt}(\beta_1), \dots, \text{nt}(\beta_n)) \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \\
 \cup & \{X_A \rightarrow \varepsilon \mid X_A \in N_A\}
 \end{aligned} \quad (1)$$

Les arbres initiaux de la grammaire TAG sont ainsi associés à des règles de la forme  $X \rightarrow \alpha(Y_1, \dots, Y_n)$  et les arbres auxiliaires à des règles  $X_A \rightarrow \beta(Y_1, \dots, Y_n)$ , où  $X$  est le non terminal qui étiquette la racine de l'arbre élémentaire TAG. Enfin, la possibilité d'une adjonction non réalisée est simulée par les règles  $X_A \rightarrow \varepsilon$ . On peut observer que la taille de la grammaire RTG obtenue est équivalente à la taille de la grammaire TAG d'origine. La traduction elle-même peut être calculée en temps linéaire.

Puisque la grammaire TAG de la figure 2 ne propose pas d'arbre auxiliaire enraciné par  $S$ ,  $VP$ ,  $Adj$  ou  $D$ , on peut simplifier les règles en ignorant ces nœuds d'adjonction. La figure 3a montre la grammaire simplifiée pour les arbres de la figure 2. Il est aisé de vérifier que cette grammaire rationnelle génère les arbres enracinés par « creux », avec « radis » et « est » pour deux fils, et une combinaison arbitraire de nœuds « le » et « gris » comme descendance de « radis » (voir figure 3b) : la grammaire rationnelle génère les arbres de dérivation d'une version sans structures de traits de la grammaire TAG d'origine.

## 2.2 Calculs d'unification

**Grammaire rationnelle d'arbres avec traits** Afin de traduire les restrictions imposées par les structures de traits de la grammaire TAG, nous considérons dans notre RTG non plus de simples réécritures entre termes, mais des *surréductions* (HANUS, 1994), c'est-à-dire des réécritures assorties d'unifications, avec des variables dans  $(N \cup N_A) \times \mathcal{D}$  où  $\mathcal{D}$  désigne l'ensemble des structures de traits possibles<sup>3</sup>.

<sup>3</sup>Étant données deux structures de traits  $d$  et  $d'$  de  $\mathcal{D}$ , on désigne par l'u-substitution  $\sigma = \text{mgu}(d, d')$  leur *unificateur le plus général* s'il existe. Nous notons  $\top$  l'élément le plus

**Definition 1.** Une *grammaire rationnelle d'arbres avec traits*  $\langle S, N, \mathcal{F}, \mathcal{D}, R \rangle$  est composée d'un axiome  $S$ , d'un ensemble de symboles non-terminaux  $N$  contenant  $S$ , d'un alphabet ordonné de terminaux  $\mathcal{F}$ , d'un ensemble de structures de traits  $\mathcal{D}$ , et d'un ensemble de règles de forme  $(A, d) \rightarrow a((B_1, d'_1), \dots, (B_n, d'_n))$  avec  $A, B_1, \dots, B_n$  des non-terminaux de  $N$ ,  $d, d'_1, \dots, d'_n$  des structures de traits de  $\mathcal{D}$ , et  $a$  un terminal d'arité  $n$  de  $\mathcal{F}$ .

La relation de dérivation  $\Rightarrow$  associée à  $G = \langle S, N, \mathcal{F}, \mathcal{D}, R \rangle$  met en relation des paires associant un terme<sup>4</sup> de  $T(\mathcal{F}, N \times \mathcal{D})$  et une u-substitution, de telle sorte que  $(s, e) \Rightarrow (t, e')$  si et seulement s'il existe un contexte<sup>5</sup>  $C$ , une règle  $(A, d) \rightarrow a((B_1, d'_1), \dots, (B_n, d'_n))$  dans  $R$  avec des variables fraîches dans les structures  $d, d'_1, \dots, d'_n$  et une u-substitution  $\sigma$  tels que

$$s = C[(A, d)], \sigma = \text{mgu}(d, e(d')), t = C[a((B_1, \sigma(d'_1)), \dots, (B_n, \sigma(d'_n)))] \text{ et } e' = \sigma \circ e$$

Le langage généré par  $G$  est  $L(G) = \{t \in T(\mathcal{F}) \mid \exists e, ((S, \top), id) \Rightarrow^* (t, e)\}$ .  $\square$

La propagation des unifications de traits se fait hiérarchiquement par la recherche de l'unificateur le plus général  $\text{mgu}$  à chaque étape de dérivation. L'u-substitution  $e$  globale associée comme environnement à notre terme sert à communiquer les résultats des unifications dans les différentes branches du terme.

**Traduction de TAG vers RTG avec traits** Munis de cette définition opérationnelle d'une RTG avec unification, nous enrichissons notre traduction de TAG vers RTG pour tenir compte des structures de traits des nœuds des arbres TAG. Nous définissons  $\text{feats}(\gamma_i)$  comme la structure de traits de  $\mathcal{D}$  associée au nœud  $\gamma_i$  de l'arbre élémentaire  $\gamma$ . Cette structure est composée de deux ensembles hauts et bas de traits atomiques  $\text{top}(\gamma_i)$  et  $\text{bot}(\gamma_i)$ .

La traduction est établie sur la notion d'*interface*  $\text{in}(\gamma)$  offerte par chaque arbre élémentaire TAG  $\gamma$ , qui servira de structure de traits de la partie gauche des règles de la grammaire rationnelle d'arbres avec traits. Dans le cas d'un arbre initial  $\alpha$ , la structure  $[\text{top} : \text{top}(\alpha_r)]$  doit s'unifier avec celle du nœud de substitution. Dans le cas d'un arbre auxiliaire  $\beta$ , la structure  $[\text{top} : \text{top}(\beta_r), \text{bot} : \text{bot}(\beta_f)]$  doit s'unifier avec celle du nœud d'adjonction. Il reste à coindexer ces interfaces avec les structures de la partie droite de chaque règle ; le seul cas à traiter est celui de la racine de l'arbre élémentaire, pour laquelle nous définissons une fonction  $\text{feats}_r$ . Nous définissons ainsi pour tout  $\alpha$  dans  $I$ ,  $\beta$  dans  $A$  et  $\gamma$  dans  $I \cup A$ , à l'aide d'une variable fraîche  $t$

$$\text{in}(\alpha) = \begin{bmatrix} \text{top} : t \\ \text{top} : \text{top}(\alpha_r) \end{bmatrix} \quad (2)$$

$$\text{in}(\beta) = \begin{bmatrix} \text{top} : t \\ \text{top} : \text{top}(\beta_r) \\ \text{bot} : \text{bot}(\beta_f) \end{bmatrix} \quad (3)$$

$$\text{feats}_r(\gamma_1) = \begin{bmatrix} \text{top} : t \\ \text{bot} : \text{bot}(\gamma_1) \end{bmatrix} \quad (4)$$

Pour un nœud  $\gamma_i$ , nous définissons  $\text{tr}(\gamma_i) = (\text{nt}(\gamma_i), \text{feats}(\gamma_i))$  et  $\text{tr}_r(\gamma_1) = (\text{nt}(\gamma_1), \text{feats}_r(\gamma_1))$ . L'ensemble de règles de notre grammaire rationnelle d'arbres

général de  $\mathcal{D}$ , et  $id$  l'u-substitution identité.

<sup>4</sup>L'ensemble des termes sur l'alphabet  $\mathcal{F}$  et l'ensemble de variables  $\mathcal{X}$  est noté  $T(\mathcal{F}, \mathcal{X})$  ; en particulier  $T(\mathcal{F}, \emptyset) = T(\mathcal{F})$  est l'ensemble des arbres sur  $\mathcal{F}$ .

<sup>5</sup>Un contexte  $C$  de  $T(\mathcal{F}, \mathcal{X})$  est un terme de  $T(\mathcal{F}, \mathcal{X} \cup \{x\})$ ,  $x \notin \mathcal{X}$ , qui ne contient qu'une seule occurrence de  $x$ , et le terme  $C[t]$  pour un terme  $t$  de  $T(\mathcal{F}, \mathcal{X})$  est obtenu en remplaçant cette occurrence par  $t$  dans  $C$ .

avec traits pour une grammaire TAG  $\langle \Sigma, N, I, A, S \rangle$  est alors

$$\begin{aligned} & \{(X, \text{in}(\alpha)) \rightarrow \alpha(\text{tr}_r(\alpha_1), \text{tr}(\alpha_2), \dots, \text{tr}(\alpha_n)) \mid \alpha \in I, n = \text{rg}(\alpha), X = \text{lab}(\alpha_r)\} \\ \cup & \{(X_A, \text{in}(\beta)) \rightarrow \beta(\text{tr}_r(\beta_1), \text{tr}(\beta_2), \dots, \text{tr}(\beta_n)) \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \\ \cup & \{X_A \left[ \begin{smallmatrix} \text{top} : x \\ \text{bot} : x \end{smallmatrix} \right] \rightarrow \varepsilon \mid X_A \in N_A, x \text{ variable de } \mathcal{D}\} \end{aligned} \quad (5)$$

Les règles dérivant la feuille vide  $\varepsilon$  effectuent l'unification finale entre traits hauts et bas des nœuds de la grammaire TAG.

Nous obtenons alors l'ensemble de règles suivant pour la grammaire rationnelle enrichie de structures de traits correspondant à la grammaire TAG de la figure 2 :

$$\begin{aligned} (S, \top) & \rightarrow \text{creux} \left( N \left[ \begin{smallmatrix} \text{top} : \left[ \begin{smallmatrix} \text{det} : + \\ \text{num} : x \end{smallmatrix} \right] \\ \text{bot} : \left[ \begin{smallmatrix} \text{top} : x \end{smallmatrix} \right] \end{smallmatrix} \right], V \left[ \begin{smallmatrix} \text{top} : \left[ \begin{smallmatrix} \text{num} : x \end{smallmatrix} \right] \end{smallmatrix} \right] \right) \\ V \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{num} : \text{sg} \end{smallmatrix} \right] \end{smallmatrix} \right] & \rightarrow \text{est} \left( V_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{num} : \text{sg} \end{smallmatrix} \right] \end{smallmatrix} \right] \right) \\ V_A \left[ \begin{smallmatrix} \text{top} : x \\ \text{bot} : x \end{smallmatrix} \right] & \rightarrow \varepsilon \\ N \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : - \end{smallmatrix} \right] \end{smallmatrix} \right] & \rightarrow \text{radis} \left( N_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : - \end{smallmatrix} \right] \end{smallmatrix} \right] \right) \\ N_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : - \\ \text{num} : x \end{smallmatrix} \right] \end{smallmatrix} \right] & \rightarrow \text{gris} \left( N_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : - \\ \text{num} : x \end{smallmatrix} \right] \end{smallmatrix} \right] \right) \\ N_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : - \\ \text{num} : \text{sg} \end{smallmatrix} \right] \end{smallmatrix} \right] & \rightarrow \text{le} \left( N_A \left[ \begin{smallmatrix} \text{top} : t \\ \text{bot} : \left[ \begin{smallmatrix} \text{det} : + \\ \text{num} : \text{sg} \end{smallmatrix} \right] \end{smallmatrix} \right] \right) \\ N_A \left[ \begin{smallmatrix} \text{top} : x \\ \text{bot} : x \end{smallmatrix} \right] & \rightarrow \varepsilon \end{aligned} \quad (6)$$

**Exemple de dérivation** Nous reprenons dans la figure 4 le cas de la phrase « Le radis gris est creux. » en employant les règles enrichies de structures de traits de l'équation (6). Chaque nœud de l'arbre de la figure est constitué d'une étiquette et d'un couple formé d'un terme de  $T(\mathcal{F}, (N \cup N_A) \times \mathcal{D})$  et d'un environnement<sup>6</sup>. La création de variables fraîches utilise l'adresse de Gorn du nœud où la réécriture a lieu. Les étiquettes de chaque nœud indiquent l'ordre dans lequel s'effectuent les surréductions. Enfin, l'on remplace les variables par leur valeur associée dans l'environnement dès que possible.

On ne peut pas dériver l'arbre correspondant à « \* Le radis gris sont creux. ». La partie gauche de la tentative de dérivation aurait été similaire. En revanche, dans la partie droite, le trait *bot* associé au nœud  $V_A$  de *sont* aurait eu pour valeur *num* : pl (pluriel). L'analyse aurait donc échoué à l'étape suivante, puisqu'en atteignant la feuille  $\varepsilon$  il aurait fallu unifier des traits *top* et *bot* avec respectivement *sg* et *pl* comme valeurs de *num*.

Bien sûr, pour une analyse qui visite d'abord le sous-arbre droit avant le sous-arbre gauche, le résultat serait le même, avec encore pour étape décisive du point de vue de l'unification la réécriture finale à  $\varepsilon$ .

### 3 Transformation par coin gauche

Comme nous venons de le voir, la génération d'un arbre de dérivation TAG à l'aide d'une grammaire RTG avec traits n'est pas très prédictive, dans le sens où il est nécessaire de patienter jusqu'à la réécriture à  $\varepsilon$  pour vérifier si une substitution réussit. Dans l'exemple de la figure 4, la substitution de « radis » dans « creux » n'est véritablement entérinée qu'au moment de la réécriture à  $\varepsilon$ ,

<sup>6</sup>Nous ne faisons apparaître que les parties de l'environnement calculées à ce point de la dérivation.

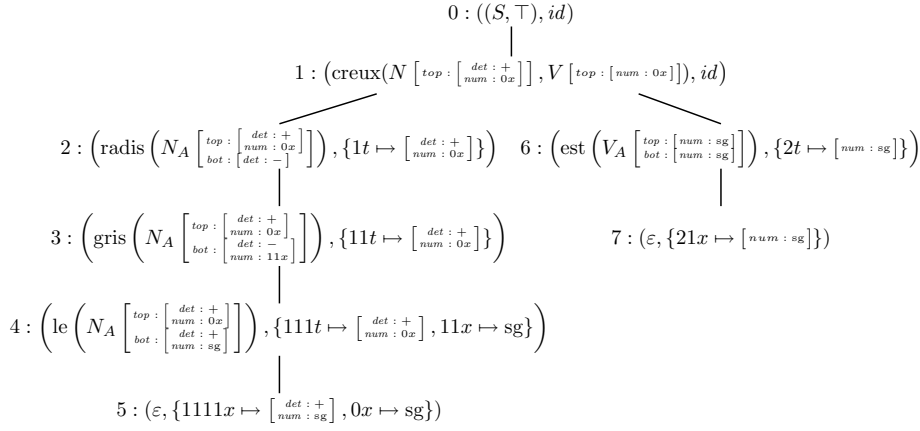


Figure 4: Une dérivation dans la RTG enrichie pour la phrase « Le radis gris est creux. »

et potentiellement toutes les opérations intermédiaires seraient à défaire si cette réécriture n'avait pas été possible.

Le seul filtrage immédiatement exercé par l'arbre « radis » lors de sa substitution au nœud  $N$  de « creux » est l'unification de sa structure  $top$  avec la structure  $top$  de  $N$ . Or, l'arbre « radis » suit l'usage dans les grammaires TAG, qui est que sa structure  $top$  est vide, et il n'y a en fait aucun filtrage par ce biais.

Nous présentons dans cette section une transformation du langage d'arbres de dérivation qui permet d'inverser l'ordre des réécritures, en commençant par  $\varepsilon$ , en opérant à toutes les adjonctions à la racine, et en finissant par l'arbre initial. Comme nous avons convenu que la racine d'un arbre élémentaire TAG apparaissait en fils gauche dans nos arbres de dérivation, cette transformation revient à une transformation par coin gauche (ROSENKRANTZ et LEWIS II, 1970) appliquée à nos grammaires rationnelles d'arbres de dérivation. Cette transformation est simple, et nous semble plus naturelle que la transformation correspondante sur les arbres dérivés.

### 3.1 Grammaire rationnelle transformée

Les règles que nous souhaitons transformer sont de la forme  $X \rightarrow \alpha(X_A, \dots)$ ,  $X_A \rightarrow \beta(X_A, \dots)$  ou  $X_A \rightarrow \varepsilon$ . À l'issue de la transformation, un appel à  $X$  devra commencer par invoquer  $\varepsilon$ , puis les adjonctions  $\beta$  en ordre inverse, et enfin  $\alpha$  en dernier lieu, dont l'arité est décrémentée. Pour notre grammaire (figure 3a), cela revient simplement à utiliser de nouveaux non-terminaux  $N_S$  et  $V_S$  et les règles

$$\begin{array}{l}
S \rightarrow \text{creux}(N_S, V_S) \\
N_S \rightarrow \varepsilon(N) \\
N \rightarrow \text{radis} \mid \text{gris}(N) \mid \text{le}(N) \\
V_S \rightarrow \varepsilon(V) \\
V \rightarrow \text{est}
\end{array} \tag{7}$$

Il manque à ces règles la possibilité d'une adjonction ailleurs qu'à la racine d'un arbre initial ; il suffit alors de conserver les règles  $X_A \rightarrow \beta(X_A, \dots)$  et



$X_A \rightarrow \varepsilon$  qui s'appliqueront comme auparavant.

Nous pouvons ensuite éliminer les  $\varepsilon$ -termes ; la grammaire de la figure 3a transformée est alors :

$$\begin{aligned} S &\rightarrow \text{creux}(N, V) \\ N &\rightarrow \text{radis} \mid \text{gris}(N) \mid \text{le}(N) \\ N_A &\rightarrow \text{gris}(N_A) \mid \text{le}(N_A) \mid \varepsilon \\ V &\rightarrow \text{est} \end{aligned} \quad (8)$$

Les règles de dérivation de  $N_A$  sont cependant inutiles puisqu'il n'y a jamais d'adjonction sur un nœud de catégorie  $N$  qui n'est pas une racine dans notre grammaire d'arbres adjoints.

Formellement, étant donnée une grammaire d'arbres adjoints  $\langle \Sigma, N, I, A, S \rangle$ , sa grammaire rationnelle d'arbres de dérivation transformée par coin gauche  $G_{\text{lc}} = \langle S, N \cup N_A, \mathcal{F}_{\text{lc}}, R_{\text{lc}} \rangle$  utilise un alphabet terminal  $\mathcal{F}_{\text{lc}} = I \cup A \cup \{\varepsilon\}$  mais où l'arité d'un arbre initial  $\alpha$  est  $\text{rg}(\alpha) - 1$ , et un ensemble de règles  $R_{\text{lc}}$  défini comme l'union

$$\begin{aligned} &\{X \rightarrow \alpha(\text{nt}(\alpha_2), \dots, \text{nt}(\alpha_n)) \mid \alpha \in I, n = \text{rg}(\alpha), X = \text{lab}(\alpha_r)\} \\ \cup &\{X \rightarrow \beta(X, \text{nt}(\beta_2), \dots, \text{nt}(\beta_n)) \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \\ \cup &\{X_A \rightarrow \beta(\text{nt}(\beta_1), \dots, \text{nt}(\beta_n)) \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \end{aligned} \quad (9)$$

La taille de cette grammaire est au pire doublée par rapport à la grammaire rationnelle d'arbres de dérivation puisque chaque arbre auxiliaire apparaît maintenant deux fois. En pratique, les règles utiles dans la grammaire obtenue sont probablement moins nombreuses. Par exemple, dans la grammaire SEMFRAG et en se basant sur l'existence de nœuds d'adjonction ailleurs qu'à la racine pour chaque catégorie syntaxique, seuls un tiers des arbres auxiliaires, soit encore un dixième des arbres élémentaires, est concerné par cette duplication.

Notons enfin que la transformation est aisément réversible. Nous définissons pour cela la fonction  $\text{lc}^{-1}$  de  $T(\mathcal{F}_{\text{lc}})$  dans  $T(\mathcal{F})$  par

$$\text{lc}^{-1}(t) = \text{revlc}(t, \varepsilon) \quad (10)$$

$$\text{revlc}(\beta(t_1, t_2, \dots, t_n), t) = \text{revlc}(t_1, \beta(t, f_{\beta_2}(t_2), \dots, f_{\beta_n}(t_n))) \quad (11)$$

$$\text{revlc}(\alpha(t_1, \dots, t_n), t) = \alpha(t, f_{\alpha_2}(t_1), \dots, f_{\alpha_{n+1}}(t_n)) \quad (12)$$

$$f_{\gamma_i}(t) = \begin{cases} \text{recur}(t) & \text{si } \gamma_i \text{ est un nœud d'adjonction} \\ \text{lc}^{-1}(t) & \text{si } \gamma_i \text{ est un nœud de substitution} \end{cases} \quad (13)$$

$$\text{recur}(\gamma(t_1, \dots, t_n)) = \gamma(f_{\gamma_1}(t_1), \dots, f_{\gamma_n}(t_n)) \quad (14)$$

On peut ainsi procéder à la génération d'un arbre dérivé dans  $L(G_{\text{lc}})$  et retrouver l'arbre correspondant de  $L(G)$  en lui appliquant  $\text{lc}^{-1}$ .

### 3.2 Unification dans la grammaire transformée

Nous procédons maintenant à la définition d'une grammaire rationnelle d'arbres de dérivation transformée par coin gauche avec structures de traits. En reprenant les règles transformées (7) de la section 3.1, nous obtenons dans un premier

temps les règles transformées avec structures de traits

$$\begin{aligned}
(S, \top) &\rightarrow \text{creux} \left( N_S \left[ \begin{array}{l} \text{top} : [ \begin{array}{l} \text{det} : + \\ \text{num} : x \end{array} ] \\ \text{bot} : [ \text{num} : x ] \end{array} \right], V \left[ \begin{array}{l} \text{top} : [ \text{num} : x ] \\ \text{bot} : [ \text{num} : x ] \end{array} \right] \right) \\
N_S \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : t \end{array} \right] &\rightarrow \varepsilon \left( N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : t \end{array} \right] \right) \\
N \left[ \begin{array}{l} \text{bot} : [ \text{det} : - ] \end{array} \right] &\rightarrow \text{radis} \\
N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : x \end{array} ] \end{array} \right] &\rightarrow \text{gris} \left( N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : x \end{array} ] \end{array} \right] \right) \\
N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : + \\ \text{num} : \text{sg} \end{array} ] \end{array} \right] &\rightarrow \text{le} \left( N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : \text{sg} \end{array} ] \end{array} \right] \right) \\
V_S \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : t \end{array} \right] &\rightarrow \varepsilon \left( V \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : t \end{array} \right] \right) \\
V \left[ \begin{array}{l} \text{bot} : [ \text{num} : \text{sg} ] \end{array} \right] &\rightarrow \text{est}
\end{aligned} \tag{15}$$

Comme la récursion au sein des arbres auxiliaires est inversée, les structures de traits de la partie gauche de chaque règle sont les structures de son nœud racine dans la grammaire TAG, et inversement (on observe ce changement pour la règle qui dérive « le »).

Nous pouvons comme auparavant éliminer les règles dérivant  $\varepsilon$ , ce qui a pour effet de copier la structure de traits *top* des nœuds de substitution dans la structure *bot*. Nous obtenons l'ensemble de règles suivant pour la grammaire TAG de la figure 2 :

$$\begin{aligned}
(S, \top) &\rightarrow \text{creux} \left( N \left[ \begin{array}{l} \text{top} : [ \begin{array}{l} \text{det} : + \\ \text{num} : x \end{array} ] \\ \text{bot} : [ \begin{array}{l} \text{det} : + \\ \text{num} : x \end{array} ] \end{array} \right], V \left[ \begin{array}{l} \text{top} : [ \text{num} : x ] \\ \text{bot} : [ \text{num} : x ] \end{array} \right] \right) \\
N \left[ \begin{array}{l} \text{bot} : [ \text{det} : - ] \end{array} \right] &\rightarrow \text{radis} \\
N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : x \end{array} ] \end{array} \right] &\rightarrow \text{gris} \left( N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : x \end{array} ] \end{array} \right] \right) \\
N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : + \\ \text{num} : \text{sg} \end{array} ] \end{array} \right] &\rightarrow \text{le} \left( N \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : [ \begin{array}{l} \text{det} : - \\ \text{num} : \text{sg} \end{array} ] \end{array} \right] \right) \\
V \left[ \begin{array}{l} \text{bot} : [ \text{num} : \text{sg} ] \end{array} \right] &\rightarrow \text{est}
\end{aligned} \tag{16}$$

Cette grammaire d'arbres avec traits est bien plus lisible que celle décrite dans l'équation (6) : le premier fils de « creux » ne peut être que « le » de par la présence du trait  $\text{det} = +$  dans la structure *bot* associée à *N*. Les seuls fils de « le » possibles sont « gris » et « radis », seuls compatibles avec le trait  $\text{det} = -$ . Le filtrage dû aux unifications est maintenant immédiat.

**Construction de la grammaire rationnelle transformée** Nous définissons les variantes suivantes des fonctions de calcul de structures de traits, pour tout arbre auxiliaire  $\beta$  de  $A$  et pour tout nœud  $\gamma_i$  d'un arbre élémentaire  $\gamma$  de  $I \cup A$  :

$$\text{in}_{1c}(\beta) = \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : \text{bot}(\beta_f) \end{array} \right] \tag{17}$$

$$\text{feats}_{1c}(\gamma_i) = \begin{cases} \left[ \begin{array}{l} \text{top} : \text{top}(\gamma_i) \\ \text{bot} : \text{top}(\gamma_i) \end{array} \right] & \text{si } \gamma_i \text{ est un nœud de substitution,} \\ \left[ \begin{array}{l} \text{top} : t \\ \text{bot} : \text{top}(\gamma_r) \end{array} \right] & \text{si } \gamma_i = \gamma_r, \\ \text{feats}(\gamma_i) & \text{sinon.} \end{cases} \tag{18}$$

Pour un nœud  $\gamma_i$ , nous notons  $\text{tr}_{1c}(\gamma_i)$  la paire  $(\text{nt}(\gamma_i), \text{feats}_{1c}(\gamma_i))$ .

Formellement, l'ensemble de règles de notre grammaire rationnelle d'arbres avec traits transformée pour une grammaire TAG  $\langle \Sigma, N, I, A, S \rangle$  est alors

$$\begin{aligned} & \{(X, \text{feats}(\alpha_1)) \rightarrow \alpha(\text{tr}_{1c}(\alpha_2), \dots, \text{tr}_{1c}(\alpha_n)) \mid \alpha \in I, n = \text{rg}(\alpha), X = \text{lab}(\alpha_r)\} \\ & \cup \{(X, \text{feats}_{1c}(\beta_1)) \rightarrow \beta((X, \text{in}_{1c}(\beta)), \text{tr}_{1c}(\beta_2), \dots, \text{tr}_{1c}(\beta_n)) \\ & \quad \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \\ & \cup \{(X_A, \text{in}(\beta)) \rightarrow \beta(\text{tr}_r(\beta_1), \text{tr}_{1c}(\beta_2), \dots, \text{tr}_{1c}(\beta_n)) \mid \beta \in A, n = \text{rg}(\beta), X = \text{lab}(\beta_r)\} \end{aligned} \tag{19}$$

## 4 Conclusion

Les grammaires rationnelles d'arbres avec structures de traits permettent de générer aisément les arbres de dérivation d'une grammaire TAG avec structures de traits. Les grammaires transformées par coin gauche permettent de plus de filtrer plus efficacement les opérations d'adjonction et de substitution possibles à partir d'un arbre élémentaire.

Si des calculs d'unification sur arbres de dérivation ont déjà été considérés par le passé de manière spécialisée (KALLMEYER et ROMERO, 2004), les mécanismes que nous avons définis sont suffisamment généraux pour traduire fidèlement l'unification dans les grammaires d'arbres adjoints.

Parmi les perspectives ouvertes par ce traitement des structures de traits dans les arbres de dérivation, on pourra mentionner des calculs d'accessibilité plus fins entre les arbres élémentaires, utiles par exemple pour vérifier qu'une TAG est dans la classe restreinte des grammaires d'arbres par insertion (SCHABES et WATERS, 1995, TIG) ou sous forme rationnelle (ROGERS, 1994, RFTAG). On pourrait par ailleurs imaginer étendre notre approche à l'analyse syntaxique, pour peu que les informations topologiques d'ordre entre les ancres soient calculées dans nos arbres de dérivation (KUHLMANN, 2007).

## References

- Marie-Hélène CANDITO et Sylvain KAHANE, 1998. Une grammaire TAG vue comme une grammaire Sens-Texte précompilée. Dans Pierre ZWEIGENBAUM, éditeur, *TALN'98*, pages 102–111. ATALA. URL [http://www.kahane.fr/?u\\_act=download&dfile=TAG-MTT-TALN1998.pdf](http://www.kahane.fr/?u_act=download&dfile=TAG-MTT-TALN1998.pdf).
- Hubert COMON, Max DAUCHET, Rémi GILLERON, Christof LÖDING, Florent JACQUEMARD, Denis LUGIEZ, Sophie TISON, et Marc TOMMASI, 2007. *Tree Automata Techniques and Applications*. URL <http://www.grappa.univ-lille3.fr/tata>.
- Philippe DE GROOTE, 2002. Tree-adjointing grammars as abstract categorial grammars. Dans Robert FRANK, éditeur, *TAG+6*, pages 145–150. URL <http://www.loria.fr/~degroote/papers/tag02.pdf>.
- Claire GARDENT, 2006. Intégration d'une dimension sémantique dans les grammaires d'arbres adjoints. Dans Piet MERTENS, Cédric FAIRON, Anne DISTER, et Patrick WATRIN, éditeurs,

- TALN'06*, pages 149–158. Presses universitaires de Louvain. URL <http://www.loria.fr/~gardent/publis/taln06-semfrag.pdf>.
- Claire GARDENT et Laura KALLMEYER, 2003. Semantic construction in feature-based TAG. Dans *EACL'03*, pages 123–130. ACL Press. ISBN 1-333-56789-0. doi: 10.3115/1067807.1067825.
- Michael HANUS, 1994. The integration of functions into logic programming: From theory to practice. *Journal of Logic Programming*, 19–20:583–628. URL <http://citeseer.ist.psu.edu/hanus94integration.html>.
- Aravind K. JOSHI et Yves SCHABES, 1997. Tree-adjointing grammars. Dans Grzegorz ROZENBERG et Arto SALOMAA, éditeurs, *Handbook of Formal Languages*, volume 3: Beyond Words, chapitre 2, pages 69–124. Springer. ISBN 3-540-60649-1. URL <http://citeseer.ist.psu.edu/joshi97treeadjoining.html>.
- Laura KALLMEYER et Maribel ROMERO, 2004. LTAG semantics with semantic unification. Dans Owen RAMBOW et Matthew STONE, éditeurs, *TAG+7*, pages 155–162. URL <http://www.cs.rutgers.edu/TAG+7/papers/kallmeyer-c.pdf>.
- Makoto KANAZAWA, 2007. Parsing and generation as Datalog queries. Dans *ACL'07*, pages 176–183. ACL Press. URL <http://www.aclweb.org/anthology/P07-1023>.
- Alexander KOLLER et Matthew STONE, 2007. Sentence generation as a planning problem. Dans *ACL'07*, pages 336–343. ACL Press. URL <http://www.aclweb.org/anthology/P07-1043>.
- Alexander KOLLER et Kristina STRIEGNITZ, 2002. Generation as dependency parsing. Dans *ACL'02*, pages 17–24. ACL Press. doi: 10.3115/1073083.1073088.
- Marco KUHLMANN, 2007. *Dependency Structures and Lexicalized Grammars*. Doctoral dissertation, Saarland University, Saarbrücken, Germany. URL <http://www.ps.uni-sb.de/Papers/abstracts/kuhlmann2007dependency.pdf>.
- Sylvain POGODALLA, 2004. Vers un statut de l'arbre de dérivation : exemples de construction de représentations sémantiques pour les grammaires d'arbres adjoints. Dans Philippe BLACHE, éditeur, *TALN'04*, pages 377–386. LPL. URL <http://www.lpl.univ-aix.fr/jep-taln04/proceed/actes/taln2004-Fez/Pogodalla.pdf>.
- J. Alan ROBINSON, 1965. A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12(1):23–41. doi: 10.1145/321250.321253.
- James ROGERS, 1994. Capturing CFLs with tree adjoining grammars. Dans *ACL'94*, pages 155–162. ACL Press. doi: 10.3115/981732.981754.
- Daniel J. ROSENKRANTZ et Philip M. LEWIS II, 1970. Deterministic left corner parsing. Dans *11th Annual Symposium on Switching and Automata Theory*, pages 139–152. IEEE Computer Society.

- Yves SCHABES et Richard C. WATERS, 1995. Tree insertion grammar: a cubic-time parsable formalism that lexicalizes context-free grammar without changing the trees produced. *Computational Linguistics*, 21(4):479–513. doi: 10.1016/0165-0114(94)00364-7.
- K. VIJAY-SHANKER, 1992. Using descriptions of trees in a tree adjoining grammar. *Computational Linguistics*, 18(4):481–517. URL <http://www.aclweb.org/anthology/J92-4004>.