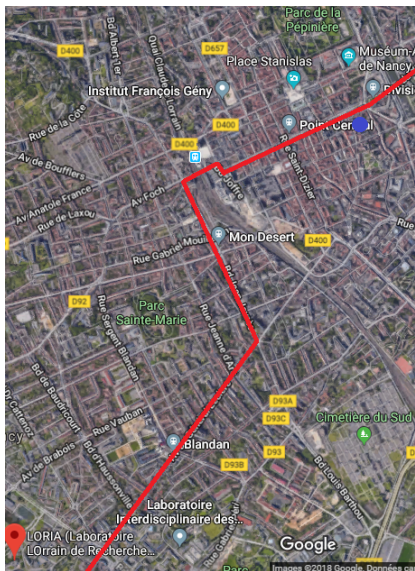


Traduction des grammaires catégorielles de Lambek dans les grammaires catégorielles abstraites

Valentin RICHARD
sous la direction de Philippe de Groote
Équipe SEMAGRAMME
Loria (Nancy)

du 11 juin au 27 juillet 2018

- 1 Environnement de travail
 - Situation géographique
 - Contexte scientifique
- 2 Sujet de recherche
- 3 Transformation des règles en coupure



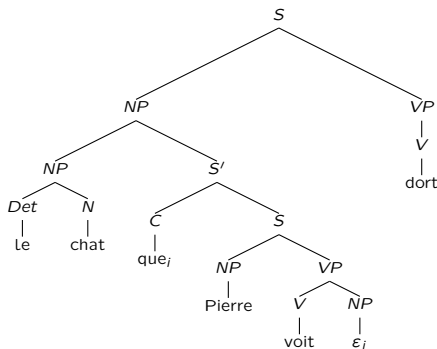
- Ligne 1 de tramway
- Mon lieu de résidence

Le Loria

- Unité Mixte de Recherche (UMR)
- 190 (enseignants-)chercheurs et 100 doctorants
- 27 équipes de recherche
- 3 grands bâtiments



La linguistique informatique et le traitement automatique des langues (TAL)



- Traduction automatique
- Recherche d'informations
- Résumé
- Modèles de la langue naturelle
- ...
- Outils de plus en plus demandés



L'équipe SEMMAGRAMME

- ACGtk : analyseur grammatical
- Outil d'aide au diagnostic de la schizophrénie
- Zombilingo : jeu pour l'annotation de ressources
- Sémantique à continuation et grammaires catégorielles



Mon maître de stage :
Philippe de Groote

1 Environnement de travail

2 Sujet de recherche

- Système logique et grammaire de Lambek
- λ -calcul typé
- Grammaire catégorielle abstraite

3 Transformation des règles en coupure

Types orientés et grammaire de Lambek

Pour Pr un ensemble de **types atomiques**, **types orientés** Tp :

$$\text{Tp} : \alpha, \beta, \gamma, \delta, \dots ::= s, r, p, q, \dots \in \text{Pr} \mid \beta/\alpha \mid \alpha \backslash \beta$$

Grammaire de Lambek [Lam58] (LG) $\mathcal{G} = (\text{Pr}, \mathcal{T}, \chi, s)$ où

- Pr types atomiques
- \mathcal{T} symboles terminaux, appelés **symboles lexicaux** ou lexèmes
- $\chi : \mathcal{T} \rightarrow \mathcal{P}(\text{Tp})$
- $s \in \text{Tp}$ un type distingué (dans la suite $s \in \text{Pr}$)

Système de dérivation et langage d'une LG

Système S_{IE} de dérivation formé sur les règles de séquents $\Gamma \vdash \beta$, avec Γ mot non vide sur $Tp \cup \mathcal{T}$ et $\beta \in Tp$:

$$\begin{array}{c}
 \frac{}{\alpha \vdash \alpha} \text{ (Ax.)} \\
 \frac{\Gamma, \alpha \vdash \beta}{\Gamma \vdash \beta/\alpha} \text{ (/I)} \\
 \frac{\alpha, \Gamma \vdash \beta}{\Gamma \vdash \alpha \backslash \beta} \text{ (\I)} \\
 \frac{\alpha \in \mathcal{X}(t)}{t \vdash \alpha} \text{ (Lex.)} \\
 \frac{\Gamma \vdash \beta/\alpha \quad \Delta \vdash \alpha}{\Gamma, \Delta \vdash \beta} \text{ (/E)} \\
 \frac{\Gamma \vdash \alpha \backslash \beta \quad \Delta \vdash \alpha}{\Delta, \Gamma \vdash \beta} \text{ (\E)} \\
 \frac{\Phi, \gamma, \Delta \vdash \beta \quad \Gamma \vdash \gamma}{\Phi, \Gamma, \Delta \vdash \beta} \text{ (CUT)}
 \end{array}$$

Figure – Règles du calcul de Lambek (déduction naturelle) sans produit lexicalisé S_{IE}

- **Dérivation** (ou preuve) : arbre dont les nœuds sont des instances des règles
- **Langage** de la grammaire de Lambek : $\mathcal{L}(\mathcal{G}) \triangleq \{\Gamma \in \mathcal{T}^+ \mid \Gamma \vdash_{IE} s\}$
- Analyse grammaticale (parsing) = recherche de preuve

Exemple

Grammaire \mathcal{G}_0 définie sur $\text{Pr} \triangleq \{s, n, np\}$

LE :	np/n	CHAT :	n	DORT :	$np \backslash s$
QUE :	$(n \backslash n)/(s/np)$	PIERRE :	np	VOIT :	$(np \backslash s)/np$

Table – Exemple 0 : Grammaire de Lambek \mathcal{G}_0

Exemple

Grammaire \mathcal{G}_0 définie sur $\text{Pr} \triangleq \{s, n, np\}$

LE :	np/n	CHAT :	n	DORT :	$np \backslash s$
QUE :	$(n \backslash n)/(s/np)$	PIERRE :	np	VOIT :	$(np \backslash s)/np$

Table – Exemple 0 : Grammaire de Lambek \mathcal{G}_0

$$\begin{array}{c}
 \frac{V \vdash (np \backslash s)/np \quad np \vdash np}{V, np \vdash np \backslash s} (/E) \quad \frac{P \vdash np}{P, V, np \vdash s} (\backslash E) \\
 \frac{Q \vdash (n \backslash n)/(s/np) \quad \frac{P, V, np \vdash s}{P, V \vdash s/np} (/I)}{Q, P, V \vdash n \backslash n} (/E) \\
 \frac{Q, P, V \vdash n \backslash n \quad C \vdash n}{C, Q, P, V \vdash n} (\backslash E)
 \end{array}$$

Figure – Exemple 0 : Haut de la dérivation de LE CHAT QUE PIERRE VOIT DORT de \mathcal{G}_0 dans S_{IE}

λ-calcul linéaire avec constantes

Pour \mathcal{T} considéré comme un ensemble de constantes et \mathcal{X} un ensemble infini dénombrable de variables, **λ-termes** linéaire $\Lambda(\mathcal{T})$:

$u, v, w \dots ::= x, y, z \dots \in \mathcal{X} \mid t, \dots \in \mathcal{T} \mid \lambda x. u$ (x au plus une fois libre dans u) $\mid u v$

- Convention de parenthésage $u v w \triangleq (u v) w$
- Typage simple (non orienté) sur un ensemble de types de base \mathcal{B} :

$\text{Tp}_{no}(\mathcal{B}) : \alpha, \beta, \gamma, \delta, \dots ::= s, r, p, q, \dots \in \mathcal{B} \mid \alpha \rightarrow \beta$

$$\frac{}{x : \alpha \vdash x : \alpha} \text{ (Ax.no)}$$

$$\frac{\alpha \in \mathcal{X}_{no}(t)}{\vdash t : \alpha} \text{ (Lex.no)}$$

$$\frac{\Gamma, x : \alpha \vdash u : \beta}{\Gamma \vdash \lambda x. u : \alpha \rightarrow \beta} (\rightarrow I_{no})$$

$$\frac{\Gamma \vdash u : \alpha \rightarrow \beta \quad \Delta \vdash v : \alpha}{\Gamma, \Delta \vdash u v : \beta} (\rightarrow E_{no})$$

Figure – Règles de typage simple (non orienté) \mathcal{S}_{no} du λ-calcul linéaire avec constantes

Exemple : $D(L(Q(\lambda x. V x P) C)) \vdash s$

Signature d'ordre supérieure et lexique

Signature d'ordre supérieur $\Sigma = (\mathcal{B}, C, \tau)$ où

- \mathcal{B} ensemble fini de types de base
- C ensemble fini de constantes
- $\tau : C \rightarrow \text{Tp}_{no}(\mathcal{B})$

Signature d'ordre supérieure et lexique

Signature d'ordre supérieur $\Sigma = (\mathcal{B}, C, \tau)$ où

- \mathcal{B} ensemble fini de types de base
- C ensemble fini de constantes
- $\tau : C \rightarrow \text{Tp}_{no}(\mathcal{B})$

Lexique de Σ_1 vers Σ_2 : un morphisme de signature d'ordre supérieur $\mathcal{L} = (F, G)$ où

- $F : \mathcal{B}_1 \rightarrow \text{Tp}_{no}(\mathcal{B}_2)$
- $G : C_1 \rightarrow \Lambda(C_2)$
- pour tout $c \in C_1$, on a $\vdash_{no} G(c) : F(\tau_1(c))$

Grammaire catégorielle abstraite

Grammaire catégorielle abstraite [dG01] $\mathfrak{G} = (\Sigma_1, \Sigma_2, \mathfrak{L}, s)$ où

- Σ_1 (resp. Σ_2) une signature d'ordre supérieur appelé le **vocabulaire abstrait** (resp. **objet**)
- $\mathfrak{L} : \Sigma_1 \rightarrow \Sigma_2$ un lexique de Σ_1 vers Σ_2 .
- $s \in \text{Tp}_{no}(\mathcal{B}_1)$ un type distingué

Langage objet : $\mathfrak{O}(\mathfrak{G}) \triangleq \{v \in \Lambda(C_2) \mid \exists u \in \Lambda(C_1), \vdash_{no} u : s \wedge v = \hat{G}(u)\}$

$L : [n] \rightarrow [np]$	$\mapsto \lambda x. LE + x$	$C : [n]$	$\mapsto CHAT$
$D : [np] \rightarrow [s]$	$\mapsto \lambda x. x + DORT$	$P : [np]$	$\mapsto PIERRE$

Table – Exemple d'ACG sous la forme $c_1 : \tau(c_1) \mapsto \mathfrak{L}(c_1)$

- 1 Environnement de travail
- 2 Sujet de recherche
- 3 Transformation des règles en coupure
 - Principe
 - Transformation des éliminations
 - Transformation des introductions
 - Des coupures aux ACG
 - Algorithme de mise à niveau

Grammaire algébrique

Grammaire algébrique (CFG) $\mathcal{G}_A = (N, A, P, s)$ où

- N alphabet fini de symboles non terminaux (K, L, \dots)
- A alphabet fini de symboles terminaux (a, b, \dots), ($N \cap A = \emptyset$)
- $P \in \mathcal{P}_f(N \times (N \cup A)^*)$ ensemble fini de productions notées $K \rightarrow w$
- $s \in N$ un symbole distingué

Grammaire algébrique

Grammaire algébrique (CFG) $\mathcal{G}_A = (N, A, P, s)$ où

- N alphabet fini de symboles non terminaux (K, L, \dots)
- A alphabet fini de symboles terminaux (a, b, \dots), ($N \cap A = \emptyset$)
- $P \in \mathcal{P}_f(N \times (N \cup A)^*)$ ensemble fini de productions notées $K \rightarrow w$
- $s \in N$ un symbole distingué

munie d'un système de réécriture de règle le passage au contexte de P :

$$\frac{K \rightarrow w \in P}{w_1 K w_2 \Rightarrow w_1 w w_2} .$$

- \Rightarrow^* : la clôture réflexive transitive de \Rightarrow
- **langage** de la grammaire algébrique : $\mathcal{L}(\mathcal{G}_A) \triangleq \{w \in A^* \mid s \Rightarrow^* w\}$
- un symbole non terminal K est **accessible** s'il existe une dérivation $s \Rightarrow^* w_1 K w_2$

Exemple pour $P : s \rightarrow np, DORT ; np \rightarrow LE, n$ et $n \rightarrow CHAT$:

$$s \Rightarrow np, DORT \Rightarrow LE, n, DORT \Rightarrow LE, CHAT, DORT$$

Une première traduction insatisfaisante

- Traduction naïve des types orientés en types non orientés :

$$\begin{aligned} \tau_0(r) &= r && \text{si } r \in \text{Pr} \\ \tau_0(\beta/\alpha) &= \tau_0(\alpha) \rightarrow \tau_0(\beta) \\ \tau_0(\alpha \setminus \beta) &= \tau_0(\alpha) \rightarrow \tau_0(\beta) \end{aligned}$$

et λ -termes associés à la preuve. Mais **surgénère** ! ex :

QUE: $(np \rightarrow s) \rightarrow n \rightarrow n$ indifférentiable de QUI: $(np \rightarrow s) \rightarrow n \rightarrow n$

Une première traduction insatisfaisante

- Traduction naïve des types orientés en types non orientés :

$$\begin{aligned}\tau_0(r) &= r && \text{si } r \in \text{Pr} \\ \tau_0(\beta/\alpha) &= \tau_0(\alpha) \rightarrow \tau_0(\beta) \\ \tau_0(\alpha \setminus \beta) &= \tau_0(\alpha) \rightarrow \tau_0(\beta)\end{aligned}$$

et λ -termes associés à la preuve. Mais **surgénère** ! ex :

QUE: $(np \rightarrow s) \rightarrow n \rightarrow n$ indifférentiable de QUI: $(np \rightarrow s) \rightarrow n \rightarrow n$

- En utilisant le résultat de Pentus [Pen97] : Les langages des grammaires de Lambek sont exactement les **langages algébriques**

Problème : Beaucoup de séquents inutiles ! exemple :

VOIT: $(n \rightarrow np) \rightarrow np \rightarrow n \rightarrow s$ dans (PIERRE VOIT LE) CHAT

Une première traduction insatisfaisante

- Traduction naïve des types orientés en types non orientés :

$$\begin{aligned} \tau_0(r) &= r & \text{si } r \in \text{Pr} \\ \tau_0(\beta/\alpha) &= \tau_0(\alpha) \rightarrow \tau_0(\beta) \\ \tau_0(\alpha \setminus \beta) &= \tau_0(\alpha) \rightarrow \tau_0(\beta) \end{aligned}$$

et λ -termes associés à la preuve. Mais **surgénère** ! ex :

QUE: $(np \rightarrow s) \rightarrow n \rightarrow n$ indifférentiable de QUI: $(np \rightarrow s) \rightarrow n \rightarrow n$

- En utilisant le résultat de Pentus [Pen97] : Les langages des grammaires de Lambek sont exactement les **langages algébriques**

Problème : Beaucoup de séquents inutiles ! exemple :

VOIT: $(n \rightarrow np) \rightarrow np \rightarrow n \rightarrow s$ dans (PIERRE VOIT LE) CHAT

Idée :

- Garder le principe des grammaires algébriques : par les **coupures**
- Transformation par réécriture de preuve

Suppression des règles d'élimination

- **Aplatissement** du type orienté comme mot autour du symbole lexical : variables pour coupures. ex : $n, Q, s/np \vdash n$
- de même pour certaines variables

$$\begin{array}{c}
 \frac{n, Q, s/np \vdash n}{\frac{\frac{\frac{np, V, np \vdash s \quad np \vdash np}{np, V, np \vdash s} \text{ (CUT)} \quad P \vdash np}{P, V, np \vdash s} \text{ (CUT)} \quad \frac{P, V, np \vdash s}{P, V \vdash s/np} \text{ (/I)}}{n, Q, P, V \vdash n} \text{ (CUT)} \quad C \vdash n}{C, Q, P, V \vdash n} \text{ (CUT)}
 \end{array}$$

- Ensemble des axiomes créés : \mathcal{A}_0

Suppression des règles d'introduction

- Réécriture \rightsquigarrow en intégrant les introductions (abstraction) directement dans les axiomes :

$$\frac{\begin{array}{c} | \\ \Gamma, \alpha \vdash \beta \end{array}}{\Gamma \vdash \beta/\alpha} \text{ (I)} \rightsquigarrow \frac{\begin{array}{c} | \\ \Gamma \vdash \beta/\alpha \end{array}}{\Gamma \vdash \beta/\alpha} \text{ (idem pour } \backslash \text{)}$$

- Axiomes construits par induction $\mathcal{A}_2 ::= \mathcal{A}_0 \mid \mathcal{Q}_1(\mathcal{A}_2) \mid \mathcal{Q}_2(\mathcal{A}_2)$, avec les **opérateurs** \mathcal{Q}_1 et \mathcal{Q}_2

$$\mathcal{Q}_1(\mathcal{A}) : \quad \text{si } \begin{array}{l} \Gamma, \gamma \vdash \beta \in \mathcal{A} \\ \text{et } \alpha \vdash_{IE} \gamma \end{array} \quad \text{alors ajouter } \Gamma \vdash \beta/\alpha$$

$$\mathcal{Q}_2(\mathcal{A}) : \quad \text{si } \Gamma, \alpha \vdash \beta \in \mathcal{A} \quad \text{alors ajouter } \Gamma, \alpha/\delta \vdash \beta/\delta \\ \text{(idem pour } \backslash \text{)}$$

Suppression des règles d'introduction

- Réécriture \rightsquigarrow en intégrant les introductions (abstraction) directement dans les axiomes :

$$\frac{\begin{array}{c} | \\ \Gamma, \alpha \vdash \beta \end{array}}{\Gamma \vdash \beta/\alpha} \text{ (I)} \rightsquigarrow \frac{\begin{array}{c} | \\ \Gamma \vdash \beta/\alpha \end{array}}{\Gamma \vdash \beta/\alpha} \text{ (idem pour } \backslash \text{)}$$

- Axiomes construits par induction $\mathcal{A}_2 ::= \mathcal{A}_0 \mid \mathcal{Q}_1(\mathcal{A}_2) \mid \mathcal{Q}_2(\mathcal{A}_2)$, avec les **opérateurs** \mathcal{Q}_1 et \mathcal{Q}_2

$$\mathcal{Q}_1(\mathcal{A}) : \quad \text{si } \begin{array}{l} \Gamma, \gamma \vdash \beta \in \mathcal{A} \\ \text{et } \alpha \vdash_{IE} \gamma \end{array} \quad \text{alors ajouter } \Gamma \vdash \beta/\alpha$$

$$\mathcal{Q}_2(\mathcal{A}) : \quad \text{si } \Gamma, \alpha \vdash \beta \in \mathcal{A} \quad \text{alors ajouter } \Gamma, \alpha/\delta \vdash \beta/\delta \\ \text{(idem pour } \backslash \text{)}$$

$$\frac{\frac{n, Q, s/np \vdash n \quad \frac{\frac{np, V \vdash s/np \quad P \vdash np}{P, V \vdash s/np} \text{ (CUT)}}{n, Q, V, P \vdash n} \text{ (CUT)}}{C, Q, P, V \vdash n} \text{ (CUT)} \quad C \vdash n \text{ (CUT)}$$

Des coupures aux ACG

- Grammaire algébrique $\mathcal{G}_A \triangleq (N, \mathcal{T}, P, s)$ où N est l'ensemble des symboles lexicaux et types de \mathcal{A}_2 , et

$$P \triangleq \{\beta \rightarrow \Gamma \mid \Gamma \vdash \beta \in \mathcal{A}_2\}$$

- Traduction d'une grammaire algébrique en ACG par [dGPog04] avec vocabulaire objet des chaînes de caractères
- Exemple :

$L : [n] \rightarrow [np]$	$\mapsto \lambda x. LE + x$	$C : [n]$	$\mapsto CHAT$
$D : [np] \rightarrow [s]$	$\mapsto \lambda x. x + DORT$	$Q : [s/np] \rightarrow [n] \rightarrow [n]$	$\mapsto \lambda x. \lambda y. y + QUE + x$
$P : [np]$	$\mapsto PIERRE$	$V : [np] \rightarrow [s/np]$	$\mapsto \lambda y. y + VOIT$

Table – Exemple 0 : Version finale de \mathcal{G}_0 en ACG

Algorithme de mise à niveau

- **Problème** : \mathcal{A}_2 potentiellement infini
- **Algorithme** : création d'axiomes arrêtée à un certain niveau
 - polynomial en la grammaire initial
- En pratique, quelques itérations des opérateurs suffisent
 - introduction = extraction linguistique : PIERRE QUI MANGE LA POMME
- Amélioration possible : travailler en accessibilité (au type s)

Conclusion

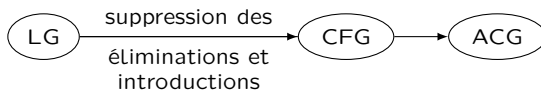


Figure – Schéma de la traduction d'une LG en ACG

Mon apport :

- Écriture et **démonstration** du procédé
- Ébauche de l'algorithme

Conclusion

Merci pour votre attention

- 1 Environnement de travail
 - Situation géographique
 - Contexte scientifique
- 2 Sujet de recherche
 - Système logique et grammaire de Lambek
 - λ -calcul typé
 - Grammaire catégorielle abstraite
- 3 Transformation des règles en coupure
 - Principe
 - Transformation des éliminations
 - Transformation des introductions
 - Des coupures aux ACG
 - Algorithme de mise à niveau

Annexes

Types de coupures :

$$\frac{\begin{array}{c} | \\ \Phi, \gamma, \Delta, \alpha \vdash \beta \end{array} \quad \begin{array}{c} | \\ \Gamma \vdash \gamma \end{array}}{\Phi, \Gamma, \Delta, \alpha \vdash \beta} \text{ (CUT) } \\ \text{Cas 1}$$

$$\frac{\begin{array}{c} | \\ \Phi, \gamma \vdash \beta \end{array} \quad \begin{array}{c} | \\ \Gamma, \alpha \vdash \gamma \end{array}}{\Phi, \Gamma, \alpha \vdash \beta} \text{ (CUT) } \\ \text{Cas 2}$$

(idem pour α à gauche)

Figure – Cas de la forme d'une coupure pour la suppression d'une introduction