

Génération d'un lexique syntaxique

Une application Perl pour le Traitement Automatique des Langues

Ingrid Falk

CNRS, Atilf, Nancy
en collaboration avec

Claire Gardent (CNRS), Guy Perrier (Université), Bruno Guillaume (INRIA)
travail financé dans le cadre du Projet « Ingénierie des Langues » CPER Lorrain
(Région Lorraine, Université, CNRS et INRIA)

Les Journées Perl 2006

Plan

Motivation

- Qu'est-ce qu'un lexique syntaxique ?
- Pour quoi faire ?
- ... et pour le Français ?

Objectifs

- Un lexique syntaxique pour le Français
- Les ressources initiales
- Le format cible
- Comment ?
- Résultats

La réalisation

- Le déroulement
- Le graphe d'une table
- Interprétation du graphe

Motivation

Objectifs

La réalisation

Au delà de la création du lexique

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

- ▶ un lexique est une liste qui associe des infos à chaque mot d'une langue :

Exemple : l'entrée d'un dictionnaire pour **bénéficier**.

BÉNÉFICIER, verbe.

A. Bénéficier de qqc. [Le suj. désigne un animé ou un inanimé]

2. Jouir, profiter de quelque chose. Bénéficier d'une aide; faire bénéficier qqc. ou qqn de :

- ▶ un lexique **syntactique** liste pour chaque usage d'un mot les constructions syntaxiques admises par cet usage.

Exemple : **bénéficier**

admets la construction : **suj v à obj**

mais n'admets pas : **suj v**

- ▶ un lexique est une liste qui associe des infos à chaque mot d'une langue :

Exemple : l'entrée d'un dictionnaire pour **bénéficier**.

BÉNÉFICIER, verbe.

A. Bénéficier de qqc. [Le suj. désigne un animé ou un inanimé]

2. Jouir, profiter de quelque chose. Bénéficier d'une aide; faire bénéficier qqc. ou qqn de :

- ▶ un lexique **syntaxique** liste pour chaque usage d'un mot les constructions syntaxiques admises par cet usage.

Exemple : **bénéficier**

admets la construction : **suj v à obj**

mais n'admets pas : **suj v**

Le cadre syntaxique (la construction syntaxique).

Définition

Une suite **a0**, **v**, **a1**, **a2**, etc. d'arguments, représentant chacun un ensemble de paires **trait-valeur**.

Un cadre syntaxique [a0 v a1] pour **bénéficiaire**

Que Max parte

a0[cat=p, fonc=sujet, mode=subj, comp='que']

bénéficie

v[cat=v, aux=avoir]

à Luc

a1[cat=n, fonc=datif, type-sem=humain,

prep='à']

Motivation

Objectifs

La réalisation

Au delà de la création du lexique

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Applications dans le Traitement Automatique des Langues :

- ▶ Analyse syntaxique et ses applications :
 - ▶ Traduction automatique.
 - ▶ Extraction d'information
 - ▶ ...

Motivation

Objectifs

La réalisation

Au delà de la création du lexique

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Pour le Français . . .

Leff production sémi-automatique,
disponible librement,
format directement compatible avec le TAL,

Dicovalence production manuelle,
librement utilisable depuis peu,
pas directement utilisable pour le TAL,

Ladl production manuelle,
disponible partiellement sous LGPL-LR,
format pas directement utilisable pour le TAL.

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Objectif : réaliser un lexique syntaxique pour le Français

- ▶ à partir des tables du Ladi (\approx 5000 verbes)
- ▶ dans un format compatible avec le TAL
- ▶ dans un format normalisé
- ▶ libre d'accès

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?
Pour quoi faire ?
... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible
Comment ?
Résultats

La réalisation

Le déroulement
Le graphe d'une table
Interprétation du graphe

Les tables du Ladl

- ▶ \approx 50 tables
- ▶ chaque table regroupe des verbes ayant un comportement commun.
- ▶ un verbe par ligne
- ▶ infos détaillée sur l'usage par les + et -
- ▶ conversion tables \rightsquigarrow lexique syntaxique n'est pas évidente !

Les tables du Ladi

| | N0 =: Nhum | N0 =: Ou P | N0 =: le fait Ou P | N0 = V1 W | N0 = V1c W | 5 | aux =: avoir | aux =: être | N0 V | | N1 =: Nhum | Prép Nhum = Ppv | N1 =: N-hum | N1 =: le fait Ou P | Loc N1 = à Nq | Loc N1 = dans Nq | Loc N1 = de Nq | N1 =: Ppv | N2 =: Ou P | N2 =: Ou Psubj | N2 = de V1 W | N2 = de V1c W | Tc =: passé | Tc =: futur | Vc =: devoir | Vc =: pouvoir | Vc =: savoir | N0 V de ce Ou P | |
|---|------------|------------|--------------------|-----------|------------|--------------|--------------|-------------|------|-----|------------|-----------------|-------------|--------------------|---------------|------------------|----------------|-----------|------------|----------------|--------------|---------------|-------------|-------------|--------------|---------------|--------------|-----------------|---------------------------------------|
| - | + | - | + | - | | advenir | - | + | | à | + | + | - | - | - | - | - | - | + | + | + | - | - | - | + | + | - | - | Il est advenu à Max que Léa a divorcé |
| + | + | + | + | + | | agir | + | - | + | sur | + | - | + | + | - | - | - | - | + | + | + | + | + | - | + | + | + | - | Que Max l'ait grondé a agi sur Luc |
| - | + | + | + | - | | bénéficiaire | + | - | - | à | + | + | - | - | - | - | - | - | + | + | - | + | + | + | + | + | - | - | Que Max parte bénéficie à Luc |

Figure: Échantillon de la table 5

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Le Format cible

- ▶ Pas de format normalisé pour les lexiques syntaxiques.
- ▶ \rightsquigarrow format de sortie **Perl**

Exemple : un cadre syntaxique en **Perl**

```
[  
  { arg => 'a0',  
    fonc => 'sujet',  
    cat => 'n' },  
  
  { arg => 'v',  
    lemme => "b\\x{e9}n\\x{e9}ficier",  
    cat => 'v',  
    aux => 'avoir' },  
  
  { arg => 'a1',  
    fonc => 'oblique',  
    cat => 'n',  
    prep => "\\x{e0}",  
    type_sem => 'humain',  
    cliticisable => 'vrai' }  
]
```

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Comment ?

En Perl à l'aide de :

- ▶ `XML::LibXML`
- ▶ `Graph::ReadWrite`
- ▶ `Graph`
- ▶ `Parse::RecDescent`
- ▶ ...

Pourquoi Perl ?

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Résultats

Un lexique syntaxique pour le Français

- ▶ ≈ 5000 verbes avec ≈ 28000 constructions.
- ▶ ≈ 550 cadres syntaxiques.
- ▶ nombre moyen de cadres par verbe : 6, 8.

En

- ▶ ≈ 10 mois/ingénieur.
- ▶ 373 de scripts et modules Perl.
- ▶ ≈ 90000 lignes de code.

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Le déroulement dans la production du lexique

1. tables en format excel \rightsquigarrow format xml (gnumeric)
2. Pour chaque table :
 - ▶ Modélisation : production manuelle du graphe de la table (graphviz dot)
 - ▶ Génération du lexique (Perl) :
 - ▶ Interpréter le graphe : quels cadres syntaxiques peut-il engendrer ?
 - ▶ Pour chaque verbe (\approx ligne de la table) générer ses cadres syntaxiques.

Le déroulement dans la production du lexique

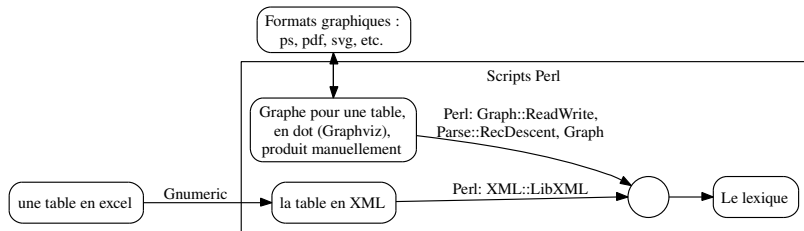


Figure: Le déroulement

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

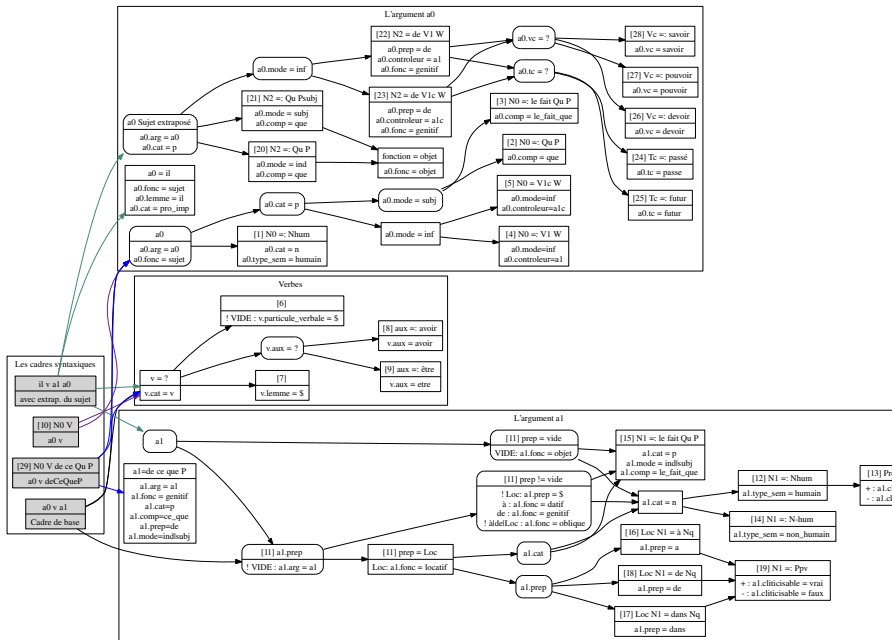
Le déroulement

Le graphe d'une table

Interprétation du graphe

Le graphe d'une table

- ▶ **Vise à rendre explicite**
 - ▶ dépendences implicites entre les colonnes de la table.
 - ▶ connaissances provenant de la description de la table, d'autres publications, de discussions.
- ▶ produit manuellement, format **dot**



Description du graphe

Graphe orienté, acyclique, *ET/OU*

Les nœuds

OU ovales, **disjonction** d'informations entre leurs nœuds fils.

ET carrée, **conjonction** d'informations

cadre racines, boîtes grisés, cadres syntaxiques associées à une table

▶ exemple

arguments les nœuds fils des nœuds *cadre*

▶ exemple

Les nœuds

Les informations associés aux nœuds sont codées dans leurs étiquettes :

Les étiquettes des nœuds

[▶ exemple](#)

[c] une colonne de la table.

conditions sur les valeurs d'une ligne de la table pour la colonne [c].

aires traits - valeurs à rajouter à l'argument si la *condition* est satisfaite pour une ligne de la table.

Motivation

Objectifs

La réalisation

Au delà de la création du lexique

Le déroulement

Le graphe d'une table

Interprétation du graphe

Rallier le graphe et la table

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Calcul des cadres syntaxiques

Cadres syntaxiques

Nœuds racine du graphe.

▶ exemple

Calculer les arguments

Calculer les chemins partant des nœuds arguments.

▶ exemple

Calculer les paires **trait-valeur** des arguments

Décoder les étiquettes des nœuds sur le chemin.

▶ exemple

En Perl

`dot` → Graph

```
Graph::ReadWrite ~> Graph
```

Calculer les chemins

```
Graph
```

Décoder les étiquettes des nœuds

```
Parse::RecDescent
```


En Perl

Décoder les étiquettes des nœuds

[▶ exemple](#)

pour chaque ligne dans la partie inférieure générer un objet **Traits** avec une fonction :

```
sub {  
    @_ = valeurs d'une ligne de la table  
    return (trait => 'valeur');  
}
```

Après le parcours du graphe

Résultat

- ▶ les cadres engendrés par la table
- ▶ chaque cadre est une suite d'arguments
- ▶ chaque argument est une liste d'objets **Traits**

Plan

Motivation

Qu'est-ce qu'un lexique syntaxique ?

Pour quoi faire ?

... et pour le Français ?

Objectifs

Un lexique syntaxique pour le Français

Les ressources initiales

Le format cible

Comment ?

Résultats

La réalisation

Le déroulement

Le graphe d'une table

Interprétation du graphe

Rallier le graphe et la table

On parcourt les lignes de la table

En Perl

- ▶ `XML::LibXML` avec `XPath`
- ▶ pour chaque **ligne** de la table (= verbe) et pour chaque **Trait** on instancie la paire **trait-valeur** correspondante, si approprié.

Et...

on y est !

+++++ bénéficiaire +++++

```

5 v[pred='bénéficiaire<subj:scompl,obja:à-sn>', @avoir, cat=v]
5 v[pred='bénéficiaire<subj:sinf,obja:à-sn>', @avoir, cat=v]
5 v[pred='bénéficiaire<subj:sn,obja:à-sn,objj:scompl>', @avoir, @il]
5 v[pred='bénéficiaire<subj:sn,obja:à-sn,objjde:de-sinf>', @avoir,
8 v[pred='bénéficiaire<subj:scompl,objjde:de-scompl>', cat=v]
8 v[pred='bénéficiaire<subj:scompl,objjde:de-sinf>', cat=v]
8 v[pred='bénéficiaire<subj:scompl,objjde:de-sn>', cat=v]
8 v[pred='bénéficiaire<subj:sinf,objjde:de-scompl>', cat=v]
8 v[pred='bénéficiaire<subj:sinf,objjde:de-sinf>', cat=v]
8 v[pred='bénéficiaire<subj:sinf,objjde:de-sn>', cat=v]
8 v[pred='bénéficiaire<subj:sn,objjde:de-scompl>', cat=v]
8 v[pred='bénéficiaire<subj:sn,objjde:de-sinf>', cat=v]
8 v[pred='bénéficiaire<subj:sn,objjde:de-sn>', cat=v]

```

Beaucoup de travaux après la création du lexique :

en cours

- ▶ Analyse et validation des résultats
- ▶ Comptages, statistiques, représentations graphiques
- ▶ Comparaisons avec d'autres lexiques

... et pour l'avenir

- ▶ mieux définir le format de sortie
- ▶ des interfaces pour gérer le(s) graphe(s)
- ▶ des interfaces pour modifier et maintenir le lexique

Beaucoup de travaux après la création du lexique :

en cours

- ▶ Analyse et validation des résultats
- ▶ Comptages, statistiques, représentations graphiques
- ▶ Comparaisons avec d'autres lexiques

... et pour l'avenir

- ▶ mieux définir le format de sortie
- ▶ des interfaces pour gérer le(s) graphe(s)
- ▶ des interfaces pour modifier et maintenir le lexique

Pourquoi Perl s'est avéré un bon choix

- ▶ Souplesse : les méthodes n'étaient pas claires dès le début
- ▶ La richesse et puissance des modules, en particulier
 - ▶ `Graph`,
 - ▶ `XML::LibXML` et
 - ▶ `Parse::RecDescent`

Conclusion

- ▶ Grace à Perl nous avons rendu possible une tache assez complexe et difficile !

Références I

Théorie

SynLex <http://www.loria.fr/~gardent/ladl/index.php>

Ladl <http://infolingu.univ-mlv.fr/DonneesLinguistiques/Lexiques-Grammaires/Presentation.html>

Implémentation : modules et scripts Perl spécifiques (non CPAN)

- ▶ <http://www.loria.fr/~falk/Ladl.html>
- ▶ <http://www.loria.fr/~falk/Ladl.html#tools>

Implémentation : modules du CPAN utilisés

Références II

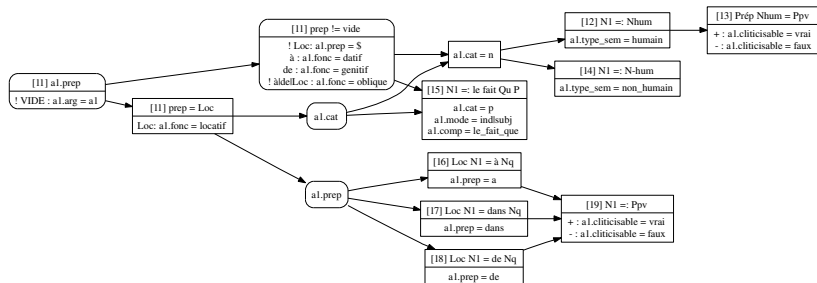
manipuler des graphes Graph, <http://search.cpan.org/author/JHI/Graph-0.80/lib/Graph.pod>

dot \rightsquigarrow Graph Graph::ReadWrite, <http://search.cpan.org/~neilb/Graph-ReadWrite-2.00/>

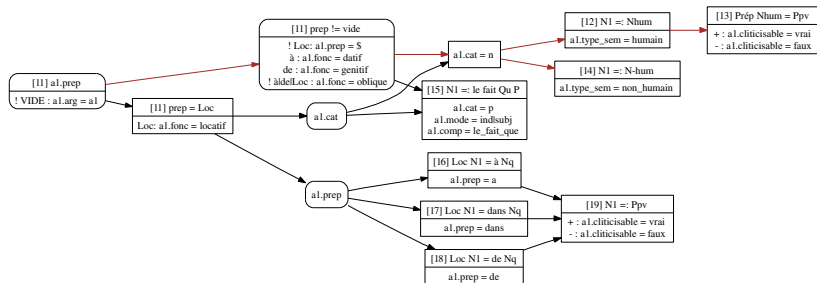
analyse Parse::RecDescent, <http://search.cpan.org/~dconway/Parse-RecDescent-1.94/lib/Parse/RecDescent.pod>

XML XML::LibXML,
<http://search.cpan.org/author/PAJAS/XML-LibXML-1.62001/LibXML.pod>

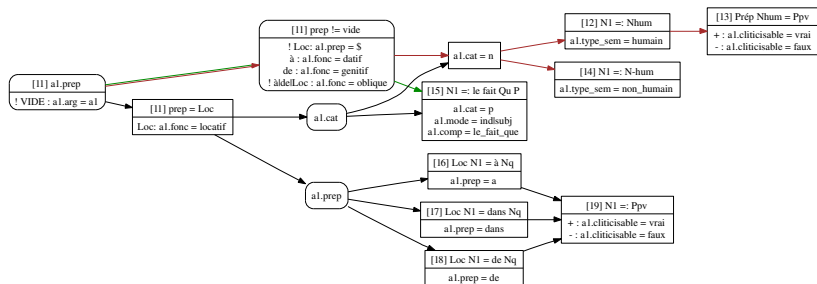
Exemple : chemins dans un graphe *ET/OU*



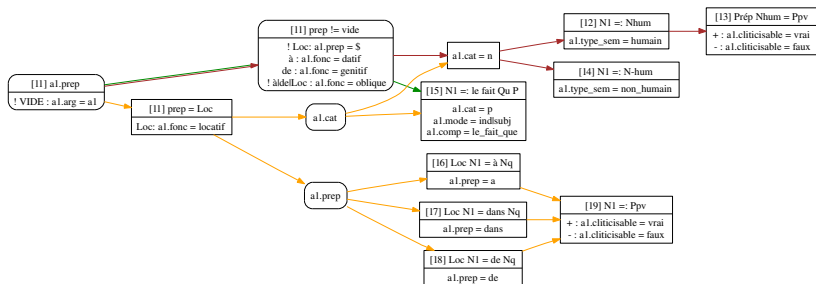
Exemple : chemins dans un graphe *ET/OU*



Exemple : chemins dans un graphe *ET/OU*

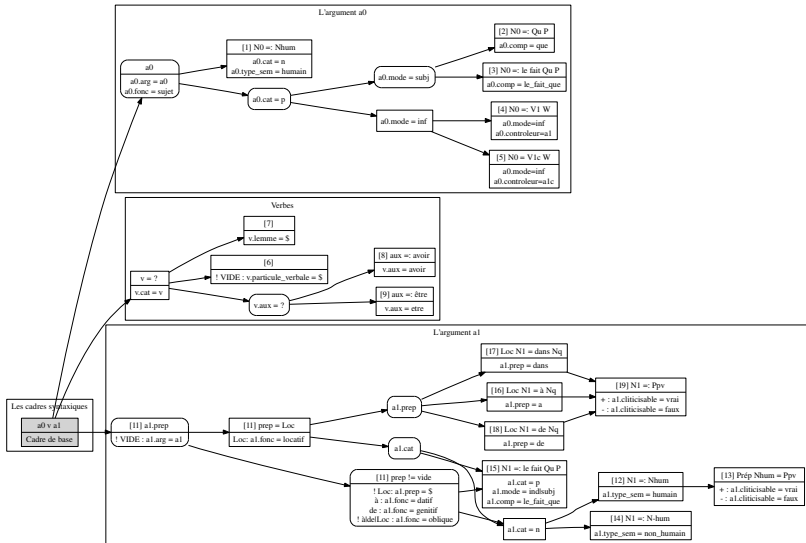


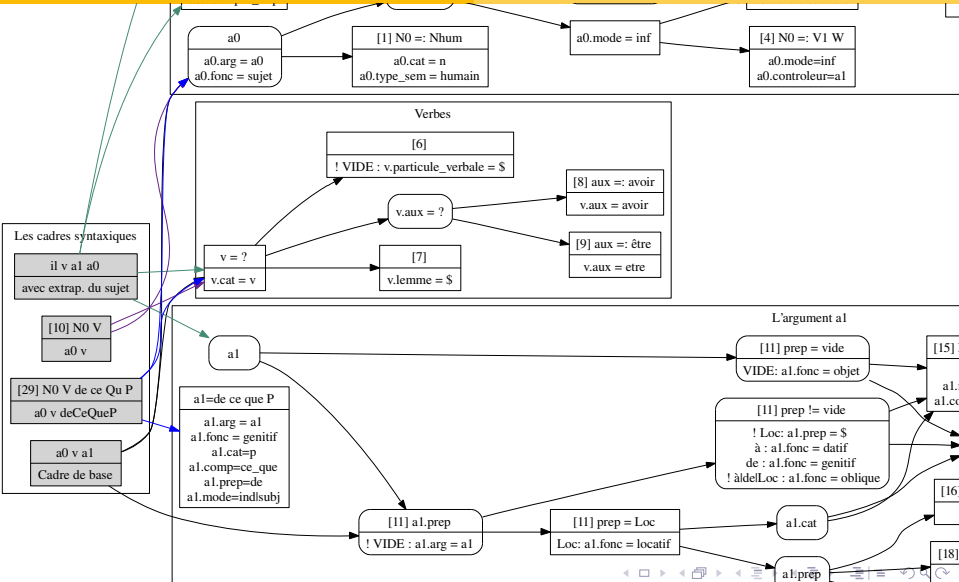
Exemple : chemins dans un graphe *ET/OU*

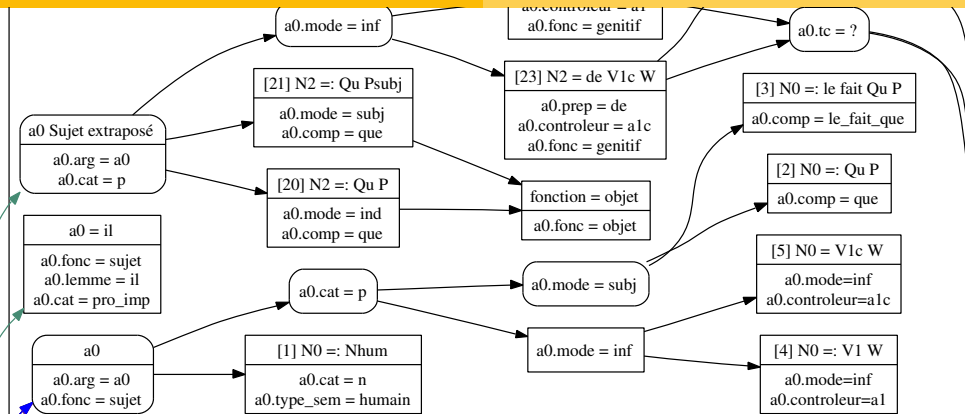


← retour

LADL table 5


[retour](#)





Verbes

[6]

! VIDE : v.particule_verbale = \$

[8] aux =: avoir

v.aux = avoir

v.aux = ?



