

# Iterated Dependencies and Kleene Iteration

Michael Dekhtyar, Alexander Dikovsky and Boris Karlov

LINA CNRS UMR 6241, Université de Nantes (France) and Tver State University (Russia)

# Plan

- 1. Dependency syntax and d-structures
- 2. Categorical dependency grammars
- 3. Kleene star closure problem
  - 3.1. A solution using multi-modal CDG
  - 3.2. Solution complexity
- 4. Conclusion

# 1. Dependency syntax and d-structures

# Dependency syntax

The traditional grammar is *word-oriented* [W.K.Percival'90]. It defines the syntax in terms of *dependencies*: relations between words.

First formalized approach to dependency structures is due to L.Tesnière [1934].

Syntactic structure is defined in terms of two kinds of binary relations on words:

- linear **precedence order** of words in the sentence, and
- **syntactic dependencies** whose graph is **acyclic** and **rooted**.



*c-copul*, *pred*, *det*, *prepos-in* are dependency names.

**Formal dependency grammars** have appeared in the early 60ies [D.Hays'61, E.Gaifman'62,65]

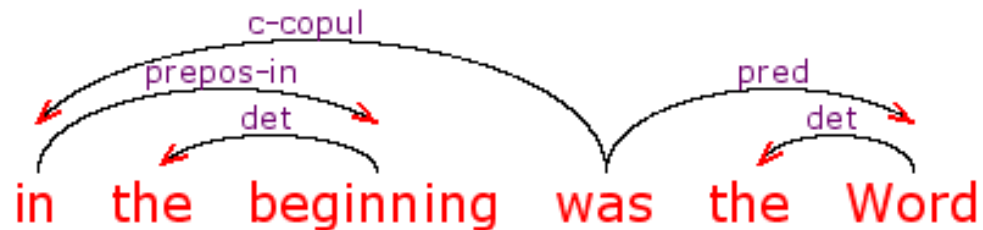
# Dependency structures

Dependency relations are asymmetric:  $W_1 \xrightarrow{d} W_2$ :

$W_1$  is a **governor (head)** of  $W_2$  through dependency  $d$

$W_2$  is a **subordinate** of  $W_1$  through  $d$ .

Dependency structure (DS) of a sentence  $x$  is linearly ordered by the precedence order in  $x$ .



$was \xrightarrow{pred} Word$ ;  $was \xrightarrow{c-copul} in$

Dominance  $W_1 \rightarrow^* W_2$ :

reflexive-transitive closure of dependency relations:  $\exists d(W_1 \xrightarrow{d} W_2)$

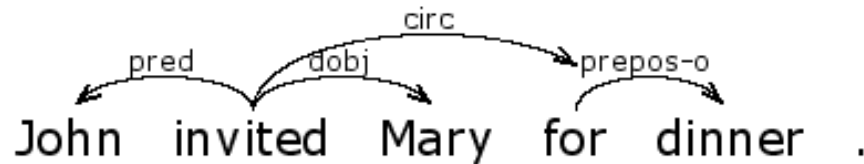
$was \rightarrow^* the$ ;  $Word \rightarrow^* Word$ .

Projection  $Proj(W) = \{W' \mid W \rightarrow^* W'\}$  ordered by precedence (a generalized *constituent* with the *head*  $W$ ).

# Projective DS

Projective DS: words' projections are **intervals**. So the dependencies don't cross

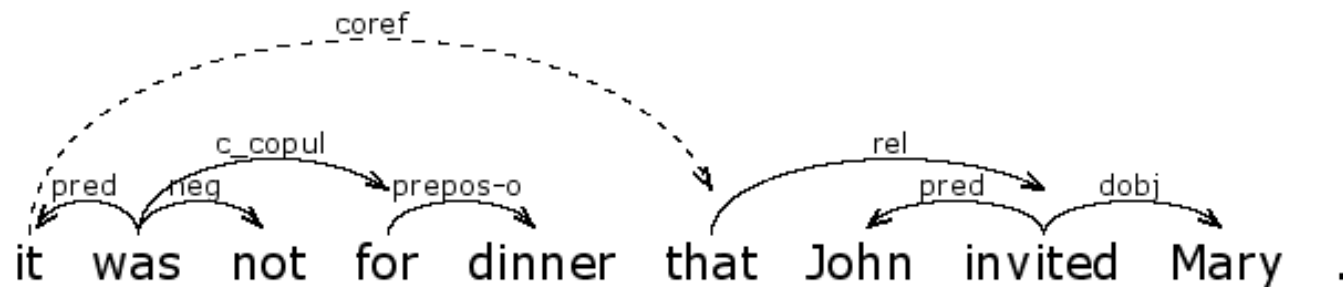
In all languages, the majority of DS are **projective**



Many authors of parsers assume the projectivity. Not the linguists. . .

Non-projective dependencies are caused by:

- **discontinuous constructions** and
- **fronting lexical units which change their communicative status** (e.g. are topicalized or defocalized).



# Non-projective DS

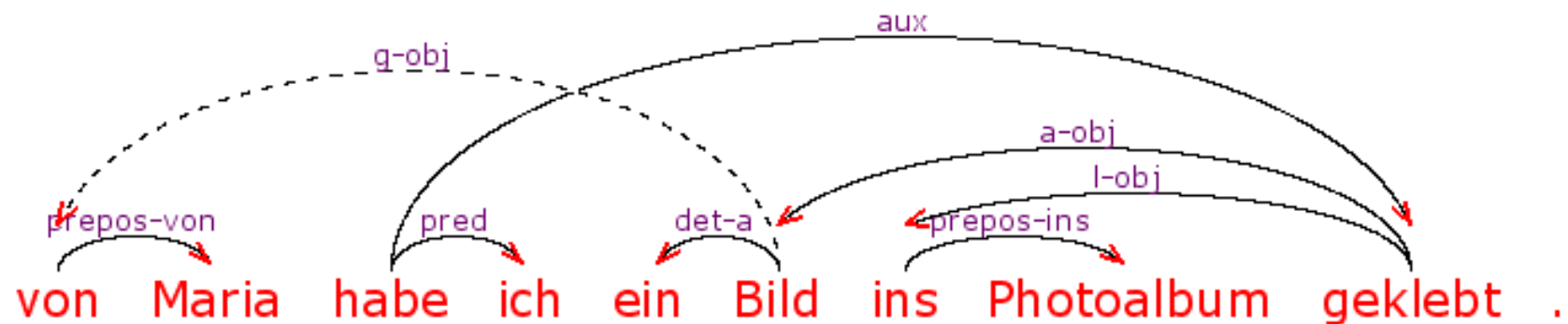
There are regular **non-projective** constructions in many languages.  
**Discontinuous comparatives in English:**



**Negation, comparatives, clitics, etc. in French:**



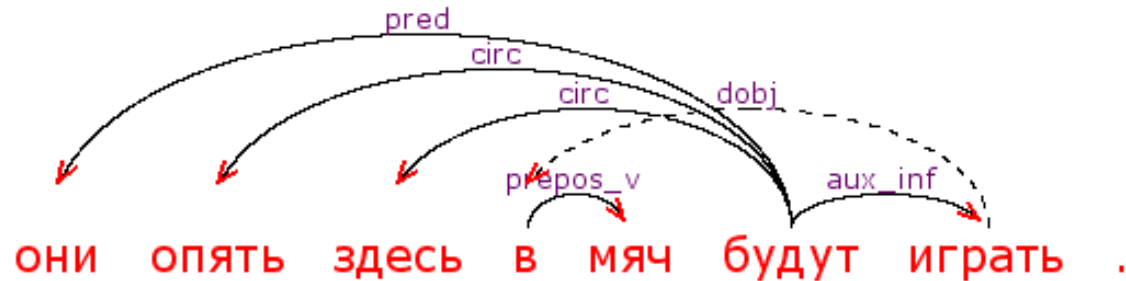
**Scrambling in German:**



(Gr. \**of, Maria have I a photo in the album stucked on*)

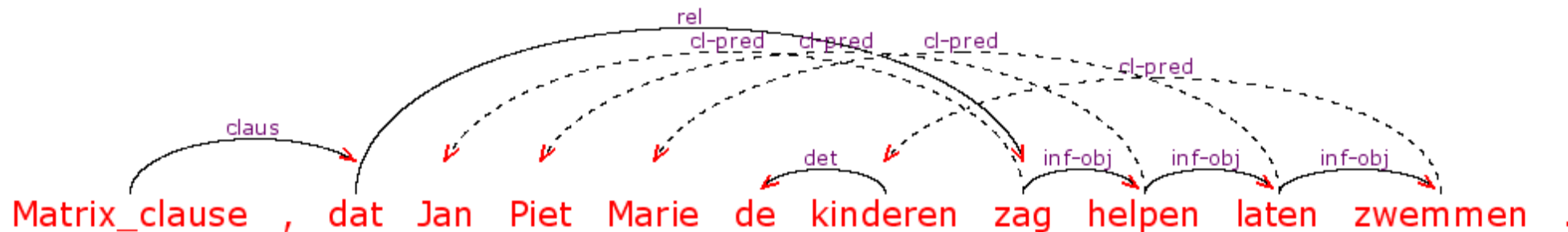
# Complex non-projective DS

Flexible word order in Slavonic languages:



(Rus. \**they again here the ball will play*)

Cross-serial dependencies in Dutch:



(Dutch \*... *that Jan Piet Marie children saw help teach swim*)

**CONCLUSION:** Dependency grammars must define non-projective DS.



## 2. Categorical dependency grammars

- are type logical grammars (categorical grammars with FO types),
- define unlimited DS (including discontinuous and long distance dependencies),
- are lexicalized and completely local (grammars of words' valencies),
- are adapted to flexible order (extended by flat regular expressions),
- are polynomially parsed (efficient deterministic offline parser is implemented),
- are well suited for practical use (cf. a wide scope large scale CDG of French).

# Types of projective dependencies

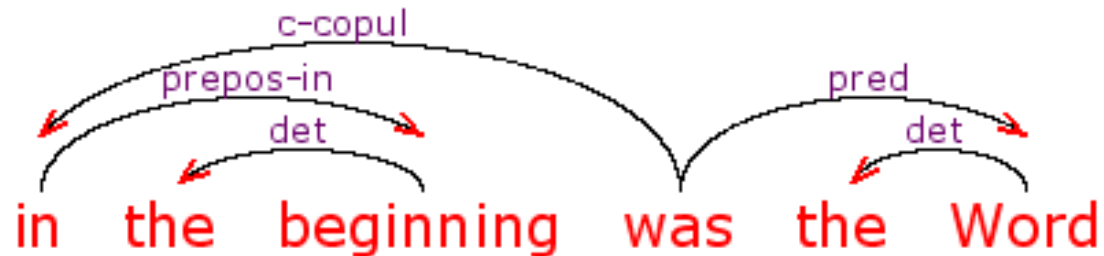
$Gov \xrightarrow{d} Sub$  is typed as :

$Gov \mapsto [..\backslash..\backslash..\backslash d/..]$  for the governor

and as

$Sub \mapsto [..\backslash d/..]$  for the subordinate.

**EX:**



$in \mapsto [c\_copul/prepos-in]$     $the \mapsto [det]$     $beginning \mapsto [det\backslash prepos-in]$   
 $was \mapsto [c\_copul\backslash S/pred]$     $Word \mapsto [det\backslash pred]$

# Necessity of iterated types

Specific interpretation of subtype elimination:

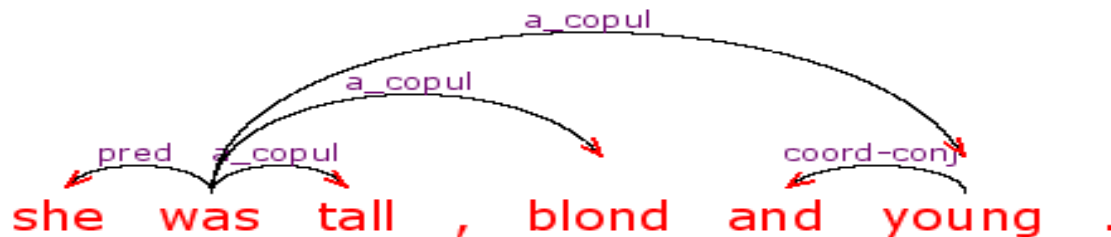
$a \mapsto [\alpha \setminus d], b \mapsto [d \setminus \beta]$  derives the dependency  $a \xleftarrow{d} b$  for  $ab$

So the recursive types derive sequenced dependencies:

$v \mapsto [c1 \setminus S], c \mapsto [c1 \setminus c1], [c1]$  derives for  $ccccc$  the DS: 

Principle of optional repeatable dependencies (see [Mel'cuk'88]):

- modifiers of a noun  $n$  share the governor  $n$ ,
- circumstantials of a verb  $v$  share the governor  $v$ .



This is why the explicit iterated types are needed:

$v \mapsto [c^* \setminus S], c \mapsto [c]$  derives for  $ccccc$  the DS:



# Types of non-projective dependencies

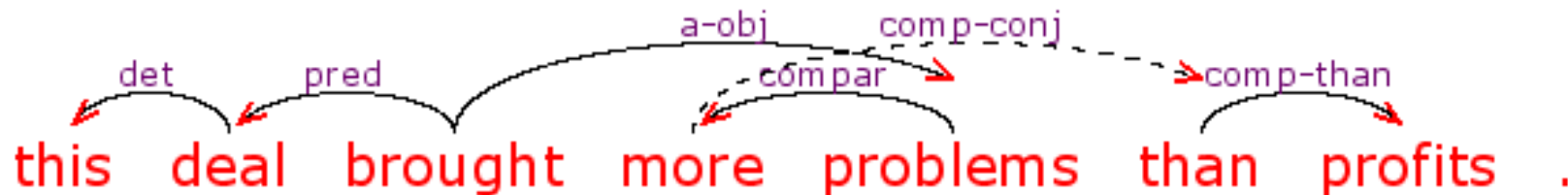
are defined using dual polarized valencies:

$Gov \overset{d}{-} > Sub$ :

$Gov \mapsto [..] \nearrow^d$  for the governor

and

$Sub \mapsto [..] \searrow^d$  for the subordinate.



$this \mapsto [det]$     $deal \mapsto [det \setminus pred]$     $brought \mapsto [pred \setminus S / a - obj]$   
 $problems \mapsto [compar \setminus a - obj]$     $profits \mapsto [comp - than]$   
 $more \mapsto [compar] \nearrow^{comp-conj}$     $than \mapsto [ / comp - than ] \searrow^{comp-conj}$

# Anchoring distant subordinates

Host words may **anchor** distant subordinates using **anchor types** :

*Host word*  $\mapsto [..\backslash..\backslash..\backslash\#(\searrow d)/..]$

*Anchored word*  $\mapsto [..\backslash\#(\searrow d)/..]\searrow d$



*elle*  $\mapsto [pred]$

*la*  $\mapsto [\#(\swarrow clit-a-obj)]\swarrow clit-a-obj$

*lui*  $\mapsto [\#(\swarrow clit-3d-obj)]\swarrow clit-3d-obj$

*a*  $\mapsto [\#(\swarrow clit-3d-obj)\backslash\#(\swarrow clit-a-obj)\backslash pred\backslash S/aux]$

*donnée*  $\mapsto [aux]\nwarrow clit-3d-obj\nwarrow clit-a-obj$

# CDG calculus

$$\mathbf{L}^1. C^{P_1} [C \setminus \beta]^{P_2} \vdash [\beta]^{P_1 P_2}$$

$$\mathbf{I}^1. C^{P_1} [C^* \setminus \beta]^{P_2} \vdash [C^* \setminus \beta]^{P_1 P_2}$$

$$\mathbf{\Omega}^1. [C^* \setminus \beta]^P \vdash [\beta]^P$$

$$\mathbf{D}^1. \alpha^{P_1} (\swarrow C) P (\searrow C) P_2 \vdash \alpha^{P_1 P P_2},$$

if  $(\swarrow C) P (\searrow C)$  satisfies the following valency pairing principle:

**FA:** in  $(\swarrow C) P (\searrow C)$ , the valency  $\swarrow C$  is the **first available** for the dual valency  $\searrow C$ , i.e.  $P$  has no occurrences of  $\swarrow C, \searrow C$ .

Similar for right valencies.

Elimination of **dual** valencies  $v = (\swarrow C), \check{v} = (\searrow C)$  by the rule  $\mathbf{D}^1$  creates the discontinuous dependency  $C$ .

Elimination of **anchored** valency  $\#(v)$  by the rules  $\mathbf{L}^1, \mathbf{I}^1$  creates not a dependency, but an anchor link.

# Proofs of typing correctness

EX: *Elle la lui a donnée*

**Without anchoring** (clitics may permute):

$[pred] [\varepsilon] \swarrow_{clit-dobj} [\varepsilon] \swarrow_{clit-iobj} [pred \setminus S/aux] [aux] \swarrow_{clit-iobj} \swarrow_{clit-dobj}$

$$\begin{array}{c}
 \frac{[pred] \quad \frac{[\varepsilon] \swarrow_{clit-dobj} \quad \frac{[\varepsilon] \swarrow_{clit-iobj} [pred \setminus S/aux]}{[pred \setminus S/aux] \swarrow_{clit-iobj}} (\mathbf{L}^l)}}{[pred \setminus S/aux] \swarrow_{clit-dobj} \swarrow_{clit-iobj}} (\mathbf{L}^l)}{[pred \setminus S/aux] \swarrow_{clit-dobj} \swarrow_{clit-iobj}} (\mathbf{L}^1)} \\
 \frac{[S/aux] \swarrow_{clit-dobj} \swarrow_{clit-iobj} \quad [aux] \swarrow_{clit-iobj} \swarrow_{clit-dobj}}{[S] \swarrow_{clit-dobj} \swarrow_{clit-iobj} \swarrow_{clit-iobj} \swarrow_{clit-dobj}} (\mathbf{L}^r)} \\
 \frac{\quad}{S} (\mathbf{D}^l \times 2)
 \end{array}$$

**With anchoring** (clitics cannot permute):

$[pred] [\#^l(\swarrow_{clit-dobj})] \swarrow_{clit-dobj} [\#^l(\swarrow_{clit-iobj})] \swarrow_{clit-iobj} [\#^l(\swarrow_{clit-iobj}) \setminus \#^l(\swarrow_{clit-dobj}) \setminus pred \setminus S/aux] [aux] \swarrow_{clit-iobj} \swarrow_{clit-dobj}$

$$\begin{array}{c}
 \frac{[\#^l(\swarrow_{clit-dobj})] \swarrow_{clit-dobj} \quad \frac{[\#^l(\swarrow_{clit-iobj})] \swarrow_{clit-iobj} [\#^l(\swarrow_{clit-iobj}) \setminus \#^l(\swarrow_{clit-dobj}) \setminus pred \setminus S/aux]}{[\#^l(\swarrow_{clit-dobj}) \setminus pred \setminus S/aux] \swarrow_{clit-iobj}} (\mathbf{L}^l)}}{[\#^l(\swarrow_{clit-dobj}) \setminus pred \setminus S/aux] \swarrow_{clit-iobj}} (\mathbf{L}^1)} \\
 \frac{[pred] \quad [S/aux] \swarrow_{clit-dobj} \swarrow_{clit-iobj}}{[pred \setminus S/aux] \swarrow_{clit-dobj} \swarrow_{clit-iobj}} (\mathbf{L}^1)} \\
 \frac{\quad \quad [aux] \swarrow_{clit-iobj} \swarrow_{clit-dobj}}{[S] \swarrow_{clit-dobj} \swarrow_{clit-iobj} \swarrow_{clit-iobj} \swarrow_{clit-dobj}} (\mathbf{L}^r)} \\
 \frac{\quad}{S} (\mathbf{D}^l \times 2)
 \end{array}$$

# CDG Definition

**CDG**  $G = (W, \mathbf{Cat}, S, \lambda)$ :

**Dictionary:**  $W$

**Types:** the set of (polarized) categories  $\mathbf{Cat}$ ,  $S \in \mathbf{Cat}$  (axiom)

**Lexicon (type assignment):**  $\lambda : W \rightarrow 2^{\mathbf{Cat}}$  (i.e.  $w \mapsto \{c_1, \dots, c_k\} \subset \mathbf{Cat}$ ).

**Calculus (one for all CDG):**  $\alpha_1 \dots \alpha_n \vdash \alpha$ ,  $\alpha, \alpha_1, \dots, \alpha_n \in \mathbf{Cat}$

**Language:**  $L(G) = \{x \in W^+ \mid (\exists \Gamma) (\Gamma \in \lambda(x) \wedge \Gamma \vdash^* S)\}$

**DS-language:**  $D(G) =$   
 $\{D(\pi) \mid (\exists x)(\exists \Gamma) (x \in L(G) \wedge \Gamma \in \lambda(x) \wedge \pi = \Gamma \vdash^* S)\}$

**EX:**  $L(G_{abc}) = \{a^n b^n c^n \mid n > 0\}$  for the CDG  $G_{abc}$  :

$$\begin{aligned} a &\mapsto A \swarrow^A, [A \setminus A] \swarrow^A, \\ b &\mapsto [B/C] \nwarrow^A, [A \setminus S/C] \nwarrow^A, \\ c &\mapsto C, [B \setminus C] \end{aligned}$$



# Properties of CDG

Projections. *local* :  $\|C^P\|_l = C$  and *valency* :  $\|C^P\|_v = P$

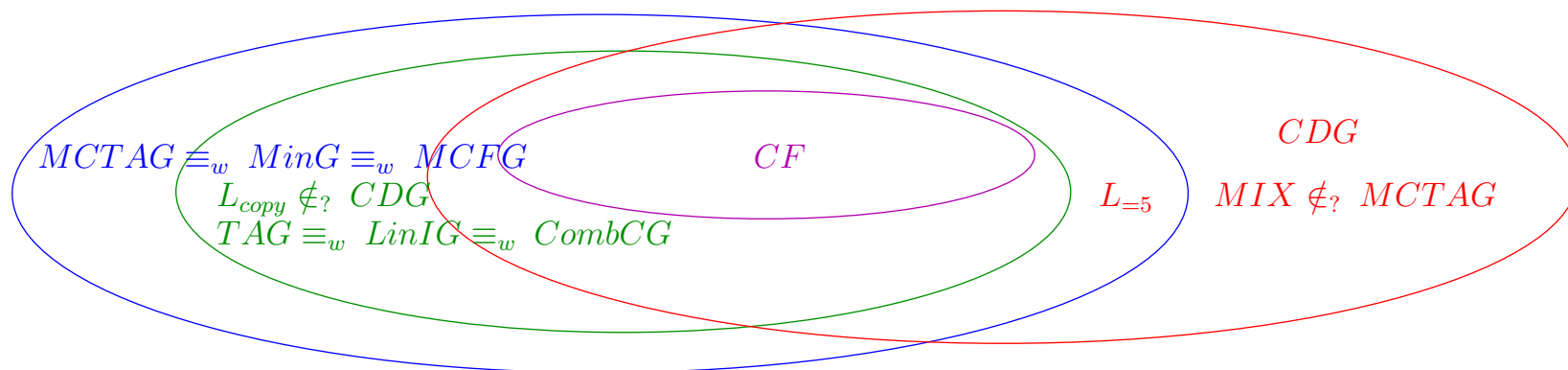
Projective dependency calculus  $\vdash_c$  : rules  $L^1, I^1, \Omega^1$ .

**TH:** For a CDG  $G = (W, C, S, \lambda)$  and  $x \in W^+$ ,  
 $x \in L(G)$  iff there is  $\Gamma \in \lambda(x)$  such that:

1.  $\|\Gamma\|_l \vdash_c S$ ,
2.  $\|\Gamma\|_v$  is well-bracketed for every valency.

Parsing complexity (Dekhtyar, Dikovsky'2004,2008) :

1. **theoretical** :  $O(n^{3+2p})$  (for  $p$  polarized valencies)
2. **in practice**:  $O(n^4)$  ( $O(n^3)$  for projective dependencies).



$$L_{copy} = \{ww \mid w \in W^+\}, \quad L_{=i} = \{a_1^n \dots a_i^n \mid n > 0\}, \quad MIX = \{w \in \{a, b, c\}^+ \mid |w|_a = |w|_b = |w|_c\}$$

# 3. Kleene star closure problem

# Closure under operations

$\mathcal{L}(CDG)$  is closed under  $HOM$ ,  $HOM^{-1}$ ,  $\cap REG$ ,  $\cup$ ,  $\circ$   
 ([Dekhtyar,Dikovsky'2008])

Explicit iterated types may suggest the idea that  $\mathcal{L}(CDG)$  is also closed under iteration  $*$  ([Dekhtyar,Dikovsky'2008] without proof). BUT:

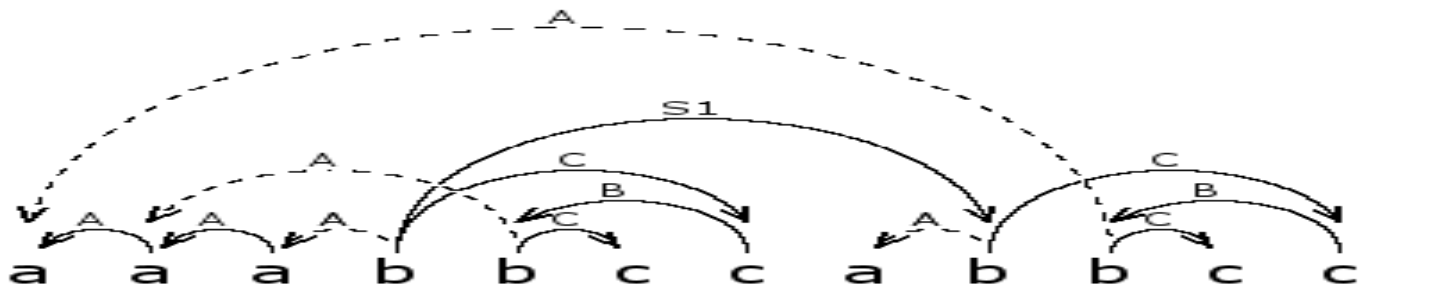
**EX:** let us consider the CDG  $G$ :

$a \mapsto A \swarrow^A, [A \setminus A] \swarrow^A,$

$b \mapsto [B/C] \swarrow^A, [A \setminus S1/C] \swarrow^A, [A \setminus S/S1/C] \swarrow^A,$

$c \mapsto C, [B \setminus C]$

It may seem that  $L(G) = L(G_{abc})L(G_{abc})$ , but it is not so because it contains, for example, the string **aaabbccabbcc** which has the DS:



So the straightforward construction does not work!

This is the effect of the completely local definition of discontinuous dependencies.

# Main Question

How to limit the scope of a polarized valency?

In particular,  
how to fix limits impenetrable for discontinuous dependencies?

E.g., in the preceding example, how to limit the scope of the valency  $\swarrow A$  in  $b \mapsto [B/C] \swarrow A, [A \setminus S_1 / C] \swarrow A$  to the substring  $aabcc$  aabbcc subordinate through  $S_1$  to the first  $b$  ?

# A multimodal solution

Multimodal CDG ([Dikovsky'07]):

Every polarized valency  $v$  may have **its own pairing rule**  $D_{R_v}$ .

For our problem, we propose the following **negative list modality**:

Let  $C$  be the name of a polarized valency and  $\pi(C)$  be a list of polarized valencies.

$$D_{FA_{C:\pi(C)}^l} \quad \alpha^{P_1(\swarrow C)P(\nwarrow C)P_2} \vdash \alpha^{P_1PP_2},$$

$(\swarrow C)P(\nwarrow C)$  satisfies the pairing rule

$FA_{C:\pi(C)}$ :  $P$  has **no occurrences** of  $\swarrow A, \nwarrow A, \nearrow A, \searrow A$  for  $A \in \pi(C)$  and also of  $\swarrow C, \nwarrow C$ .

**Intuitively**: the discontinuous dependency  $C$  **cannot cross** discontinuous dependencies in the list  $\pi(C)$ .

$\mathcal{L}(mmCDG^{-FA})$ : class of languages generated by CDG with negative list modalities.

**Theorem 1**  $\mathcal{L}(mmCDG^{-FA})$  is **closed under Kleene iteration** (so it is an **AFL**).

# Solution complexity

$mmCDG^{-FA}$  turn out to be very expressive:

$L_{exp} = \{1010^21 \dots 10^{2^n}1 \mid n > 1\}$  is generated by  $mmCDG^{-FA} G_{exp}$ :

$1 \mapsto [S/A_1] \nearrow^X, [D_1/C] \nearrow^Y, [B/A] \searrow^X \nearrow^X, [D/C] \searrow^Y \nearrow^Y, [B/A_2] \searrow^X \nearrow^X,$   
 $[D/C_2] \searrow^Y \nearrow^Y, [A_2] \searrow^X \searrow^Y, [C_2] \searrow^X \searrow^Y$

$0 \mapsto [A_1/D_1] \nearrow^B \nearrow^B, [A/A] \searrow^A \nearrow^B \nearrow^B, [A/D] \searrow^A \nearrow^B \nearrow^B, [C/C] \searrow^B \nearrow^A \nearrow^A,$   
 $[C/B] \searrow^B \nearrow^A \nearrow^A, [A_2/A_2] \searrow^A, [C_2/C_2] \searrow^B$

with  $\pi(X) = \{A\}, \pi(Y) = \{B\}$ .

**Corollary 1** *Languages in  $\mathcal{L}(mmCDG^{-FA})$  may be not semilinear.*

**Theorem 2** *Membership problem for  $mmCDG^{-FA}$  is NP-complete.*

(By reduction of 3-SAT; four discontinuous non-crossing dependencies are used:  $\pi(0) = \{1\}, \pi(1) = \{0\}, \pi(0') = \{1'\}, \pi(1') = \{0'\}$ .)

$\{ww^Rw \mid w \in \{a, b\}^+\}$  and  $L_{copy}$  also belong to  $\mathcal{L}(mmCDG^{-FA})$ .

# 4. Conclusion

- 1. Categorical Dependency Grammars
  - define unlimited dependency structures (projective, discontinuous, non-tree)
  - use explicit iterated types
  - do not use the non-cross modalities
  - are perfectly adapted for real applications
  - are polynomially analyzed
- 2. One non-crossed discontinuous dependency is sufficient for Kleene iteration closure
- 3. Four non-crossed discontinuous dependencies are sufficient for NP-hardness
- 4. Still remain open the problems:
  - $L_{copy} \notin \mathcal{L}(CDG)$
  - CDG-languages are semilinear
  - closure/non-closure of  $\mathcal{L}(CDG)$  under Kleene iteration

$$\{a^n b^n c^n \mid n > 0\}$$

EX: The CDG  $G_{abc}$  :

$$\begin{aligned}
 a &\mapsto A \swarrow^A, [A \setminus A] \swarrow^A, \\
 b &\mapsto [B/C] \nwarrow^A, [A \setminus S/C] \nwarrow^A, \\
 c &\mapsto C, [B \setminus C]
 \end{aligned}$$

generates the language  $\{a^n b^n c^n \mid n > 0\}$ .

E.g.,  $a^3 b^3 c^3$  obtains the DS:



due to the proof:

$$\begin{array}{c}
 \frac{[A] \swarrow^A \quad [A \setminus A] \swarrow^A}{[A] \swarrow^A \swarrow^A} (\mathbf{L}^l) \quad \frac{[A \setminus A] \swarrow^A \quad [A \setminus S/C] \nwarrow^A}{[A \setminus S] \nwarrow^A \swarrow^A} (\mathbf{L}^l) \quad \frac{[B/C] \nwarrow^A \quad [B \setminus C]}{[B \setminus C] \nwarrow^A} (\mathbf{L}^l) \\
 \frac{[A] \swarrow^A \swarrow^A \swarrow^A}{[A] \swarrow^A \swarrow^A \swarrow^A} (\mathbf{L}^l) \quad \frac{[A \setminus S] \nwarrow^A \swarrow^A \swarrow^A \swarrow^A}{[A \setminus S] \nwarrow^A \swarrow^A \swarrow^A \swarrow^A} (\mathbf{L}^l) \quad \frac{[B \setminus C] \nwarrow^A \quad [B \setminus C]}{[B \setminus C] \nwarrow^A} (\mathbf{L}^l) \\
 \frac{[S] \swarrow^A \swarrow^A \swarrow^A \swarrow^A \swarrow^A \swarrow^A}{S} (\mathbf{D}^l \times 3)
 \end{array}$$



# CDG of French

Parameters of CDG of French constructed by bootstrapping from examples:

Examples	Classes			Regular Expressions	Dependencies	
	total	verbal	nominal		projective	distant
500	190	46	6	~ 2500	85(9 <i>par</i> )	22(3 <i>par</i> )

(where  $n(m \text{ par})$  means  $n$  of which  $m$  are parametrized)

The grammar is developed by A. Dikovsky during 9 months (beginning of June 2009 - end of March 2010).

It explains :

- the kernel syntax (copulas/auxiliary verbs, verb complements (a-obj, d-obj, g-obj, l-obj, o-obj, Agent), noun modifiers, verb circumstantials, pronouns),
- genitive/partitive,
- negation (*ne..pas* (etc.), *ne..que*),
- clitics and reflexives,
- light verbs and infinitives,
- comparative constructions,
- clauses, coordination,
- word order inversion : interrogative sentences, topicalized complements,
- vocative constructions, co-reference, expletives (les incises),
- aggregation,
- complex lexical units and numbers, unknown terms.

# Integration of a morphological dictionary

**PROBLEM :** Lexicon of the bootstrapped CDG  $G_{approx}[T]$  is limited to the examples from which it was constructed ( $\sim 2000$  words).

To extend it to a wide scope French lexicon, was used the open resource morphological dictionary of French [Lefff 3.0](#) [B.Sagot'08].

Lefff 3.0 has more than 500000 forms annotated with morpho-syntactic feature values:

**EX :** *octroierions* 100 *v* [*pred = ' octroyer*\_\_\_\_\_1  
< *Suj : cln|scompl|sinf|sn, Obj : (cla|sn)* > ,  
*@pers, cat = v, @C1p] octroyer*\_\_\_\_\_1 *Default C1p*

Integration of Lefff with CDG :  $G_{approx}[T]$  :

- $G_{approx}[T]$  and Lefff are stored in the RDB *PostgreSQL* (R.Alfared)
- correspondence between CDG classes and the corresponding Lefff forms is realized through SQL queries or manually (R.Alfared, D.Béchet, A.Dikovsky, during 5 months).