

# Parsing Beyond Context-Free Grammars: Linear Context-Free Rewriting Systems

Laura Kallmeyer & Tatiana Bladier  
Heinrich-Heine-Universität Düsseldorf

Sommersemester 2018

# Overview

- 1 Basic Ideas
- 2 LCFRS and CL
- 3 LCFRS and MCFG
- 4 LCFRS with Simple RCG syntax

[Kal10]

# Basic Ideas (1)

Linear Context-Free Rewriting Systems (LCFRS) can be conceived as a natural extension of CFG:

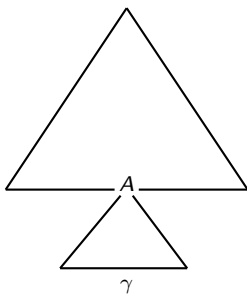
- In a CFG, non-terminal symbols  $A$  can span single strings, i.e., the language derivable from  $A$  is a subset of  $T^*$ .
- Extension to LCFRS: non-terminal symbols  $A$  can span tuples of (possibly non-adjacent) strings, i.e., the language derivable from  $A$  is a subset of  $(T^*)^k$

⇒ LCFRS displays an extended domain of locality

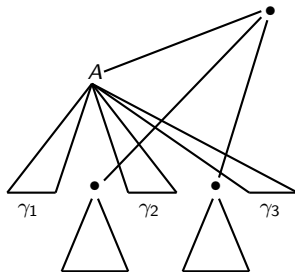
## Basic Ideas (2)

Different spans in CFG and LCFRS:

CFG:



LCFRS:



## Basic Ideas (3)

Example for a non-terminal with a yield consisting of 2 components:

$$\text{yield}(A) = \langle a^n b^n, c^n d^n \rangle, \text{ with } n \geq 1.$$

The rules in an LCFRS describe how to compute an element in the yield of the lefthand-side (lhs) non-terminal from elements in the yields of the right-hand side (rhs) non-terminals.

$$\text{Ex.: } A(ab, cd) \rightarrow \varepsilon \quad A(aXb, cYd) \rightarrow A(X, Y)$$

The start symbol  $S$  is of dimension 1, i.e., has single strings as yield elements.

$$\text{Ex.: } S(XY) \rightarrow A(X, Y)$$

Language generated by this grammar (yield of  $S$ ):

$$\{a^n b^n c^n d^n \mid n \geq 1\}.$$

## Basic Ideas (4)

- In a CFG derivation tree (parse tree), dominance is determined by the relations between lhs symbol and rhs symbols of a rule.
- Furthermore, there is a linear order on the terminals and on all rhs of rules.

In an LCFRS, we can also obtain a derivation tree from the rules that have been applied:

- Dominance is also determined by the relations between lhs symbol and rhs symbols of a rule.
- There is a linear order on the terminals. BUT: there is no linear order on all rhs of rules.

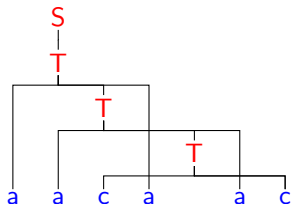
As a convention, we draw a non-terminal  $A$  left of a non-terminal  $B$  if the first terminal in the span of  $A$  precedes the first terminal in the span of  $B$ .

# Basic Ideas (5)

Ex.: LCFRS for  $\{wcwc \mid w \in \{a, b\}^*\}$ :

$$\begin{array}{ll}
 S(XY) \rightarrow T(X, Y) & T(aY, aU) \rightarrow T(Y, U) \\
 T(bY, bU) \rightarrow T(Y, U) & T(c, c) \rightarrow \varepsilon
 \end{array}$$

Derivation tree for *aacaac*:



# LCFRS and CL (1)

Interest of LCFRS for CL:

- ① Applications in CL (parsing, grammar engineering, etc.).
- ② Mild context-sensitivity.
- ③ Equivalence with several important CL formalisms.



# LCFRS and CL (2)

## Applications in CL

- **Grammar engineering and language modeling:** *Grammatical Framework* is a framework which is equivalent to LCFRS [Lju04]. It is actively used for multilingual grammar development and allows for an easy treatment of discontinuities [Ran11].
- **Grammar engineering and parsing:** In TuLiPA [KMPD10], a multi-formalism parser used in a development environment for variants of Tree Adjoining Grammar (TAG), LCFRS acts as a pivot formalism, i.e., instead of parsing directly with a TAG variant, TuLiPA parses with its equivalent LCFRS, obtained through a suitable grammar transformation [KP08].

# LCFRS and CL (3)

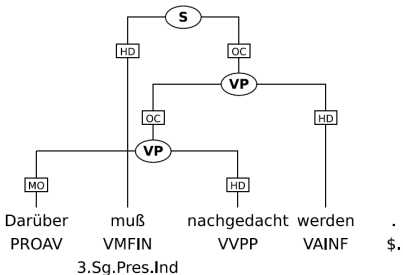
- Modeling of **non-concatenative morphology** [BB13], such as stem derivation in Semitic languages. In such languages, words are derived by combining a discontinuous root with a discontinuous template.

Ex. (Arabic):

*k i t a b* (“book”), *k a t i b* (“writer”), ma *k t a b* (“desk”)

# LCFRS and CL (4)

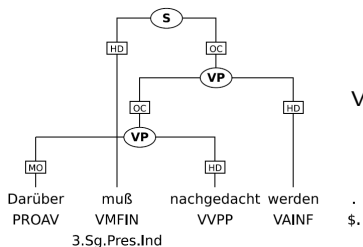
- **Syntax and data-driven parsing:**
  - Just like phrase structure trees (without crossing branches) can be described with CFG rules, trees with crossing branches can be described with LCFRS rules.
  - Trees with crossing branches allow to describe discontinuous constituents, as for example in the Negra and Tiger treebanks.



# LCFRS and CL (4)

Trees with crossing branches can be interpreted as LCFRS derivation trees.

⇒ an LCFRS can be straight-forwardly extracted from such treebanks. This makes LCFRS particularly interesting for data-driven parsing.



PROAV(Darüber)	→	$\epsilon$
VMFIN(muß)	→	$\epsilon$
VVPP(nachgedacht)	→	$\epsilon$
VAINF(werden)	→	$\epsilon$
$S(X_1, X_2, X_3)$	→	$VP(X_1, X_3) VMFIN(X_2)$
$VP(X_1, X_2, X_3)$	→	$VP(X_1, X_2) VAINF(X_3)$
$VP(X_1, X_2)$	→	$PROAV(X_1) VVPP(X_2)$

LCFRS has been successfully used for data-driven probabilistic syntactic parsing [KM13, vC12, AL14].

# LCFRS and CL (5)

- **Machine translation:** Synchronous LCFRS have been used for the modeling of translational equivalence [Kae15]. They can model certain alignment configurations that cannot be modeled with synchronous CFGs [Kae13].

(1)    *je ne veux plus jouer*  
          *I do not want to play anymore*

$\langle X(\text{jouer}) \rightarrow \varepsilon, X(\text{to play}) \rightarrow \varepsilon \rangle$

$\langle X(\text{veux}) \rightarrow \varepsilon, X(\text{do, want}) \rightarrow \varepsilon \rangle$

$\langle X(\text{ne } x_1 \text{ plus } x_2) \rightarrow X_{\boxed{1}}(x_1)X_{\boxed{2}}(x_2),$

$X(x_1 \text{ not } x_2 x_3 \text{ anymore}) \rightarrow X_{\boxed{1}}(x_1, x_2)X_{\boxed{2}}(x_3) \rangle$

...

# LCFRS and CL (6)

## Mild Context-Sensitivity:

- Natural languages are not context-free.
- Question: How complex are natural languages? In other words, what are the properties that a grammar formalism for natural languages should have?
- Goal: extend CFG only as far as necessary to deal with natural languages in order to capture the complexity of natural languages.

This effort has led to the definition of *mild context-sensitivity* (Aravind Joshi).

# LCFRS and CL (7)

A formalism is mildly context-sensitive if the following holds:

- ① It generates at least all context-free languages.
- ② It can describe a limited amount of crossing dependencies.
- ③ Its string languages are polynomial.
- ④ Its string languages are of constant growth.

# LCFRS and CL (8)

- LCFRS are mildly context-sensitive.
- We do not have any other formalism that is also mildly context-sensitive and whose set of string languages properly contains the string languages of LCFRS.
- Therefore, LCFRS are often taken to provide a grammar-formalism-based characterization of mild context-sensitivity.

BUT: There are polynomial languages of constant growth that cannot be generated by LCFRS.



# LCFRS and CL (8)

## Equivalence with CL formalisms:

LCFRS are weakly equivalent to

- *set-local Multicomponent Tree Adjoining Grammar*, an extension of TAG that has been motivated by linguistic considerations;
- *Minimalist Grammar*, a formalism that was developed in order to provide a formalization of a GB-style grammar with transformational operations such as movement;
- *finite-copying Lexical Functional Grammar*, a version of LFG where the number of nodes in the c-structure that a single f-structure can be related with is limited by a grammar constant.

# LCFRS and MCFG (1)

- *Multiple Context-Free Grammars (MCFG)* have been introduced by [SMFK91] while the equivalent *Linear Context-Free Rewriting Systems (LCFRS)* were independently proposed by [VSWJ87].
- The central idea is to extend CFGs such that non-terminal symbols can span a tuple of strings that need not be adjacent in the input string.
- The grammar contains productions of the form  $A_0 \rightarrow f[A_1, \dots, A_q]$  where  $A_0, \dots, A_q$  are non-terminals and  $f$  is a function describing how to compute the yield of  $A_0$  (a string tuple) from the yields of  $A_1, \dots, A_q$ .
- The definition of LCFRS is slightly more restrictive than the one of MCFG. However, [SMFK91] have shown that the two formalisms are equivalent.

# LCFRS and MCFG (2)

Example: MCFG/LCFRS for the double copy language.

Rewriting rules:

$$S \rightarrow f_1[A] \quad A \rightarrow f_2[A] \quad A \rightarrow f_3[A] \quad A \rightarrow f_4[ ] \quad A \rightarrow f_5[ ]$$

Operations:

$$f_1[\langle X, Y, Z \rangle] = \langle XYZ \rangle$$

$$f_4[ ] = \langle a, a, a \rangle$$

$$f_2[\langle X, Y, Z \rangle] = \langle aX, aY, aZ \rangle$$

$$f_5[ ] = \langle b, b, b \rangle$$

$$f_3[\langle X, Y, Z \rangle] = \langle bX, bY, bZ \rangle$$

# LCFRS and MCFG (3)

## Definition 1 (Multiple Context-Free Grammar)

A *multiple context-free grammar (MCFG)* is a 5-tuple  $\langle N, T, F, P, S \rangle$  where

- $N$  is a finite set of non-terminals, each  $A \in N$  has a *fan-out*  $\dim(A) \geq 1$ ,  $\dim(A) \in \mathbb{N}$ ;
- $T$  is a finite set of terminals;
- $F$  is a finite set of mcf-functions;
- $P$  is a finite set of rules of the form  $A_0 \rightarrow f[A_1, \dots, A_k]$  with  $k \geq 0$ ,  $f \in F$  such that  $f : (T^*)^{\dim(A_1)} \times \dots \times (T^*)^{\dim(A_k)} \rightarrow (T^*)^{\dim(A_0)}$ ;
- $S \in N$  is the start symbol with  $\dim(S) = 1$ .

A MCFG with maximal non-terminal fan-out  $k$  is called a  $k$ -MCFG.

# LCFRS and MCFG (4)

Mcf-functions are such that

- each component of the value of  $f$  is a concatenation of some constant strings and some components of its arguments.
- Furthermore, each component of the right-hand side arguments of a rule is not allowed to appear in the value of  $f$  more than once.

# LCFRS and MCFG (5)

## Definition 2 (mcf-function)

$f$  is an mcf-function if there is a  $k \geq 0$  and there are  $d_i > 0$  for  $0 \leq i \leq k$  such that  $f$  is a total function from  $(T^*)^{d_1} \times \dots \times (T^*)^{d_k}$  to  $(T^*)^{d_0}$  such that

- the components of  $f(\vec{x}_1, \dots, \vec{x}_k)$  are concatenations of a limited amount of terminal symbols and the components  $x_{ij}$  of the  $\vec{x}_i$  ( $1 \leq i \leq k, 1 \leq j \leq d_i$ ), and
- the components  $x_{ij}$  of the  $\vec{x}_i$  are used at most once in the components of  $f(\vec{x}_1, \dots, \vec{x}_k)$ .

A LCFRS is a MCFG where the mcf-functions  $f$  are such that the the components  $x_{ij}$  of the  $\vec{x}_i$  are used exactly once in the components of  $f(\vec{x}_1, \dots, \vec{x}_k)$ .

# LCFRS and MCFG (6)

- We can understand a MCFG as a generative device that specifies the yields of the non-terminals.
- The language of a MCFG is then the yield of the start symbol  $S$ .

Ex.: LCFRS for the double copy language.

$$\text{yield}(A) = \{\langle w, w, w \rangle \mid w \in \{a, b\}^*\}$$

$$\text{yield}(S) = \{\langle www \rangle \mid w \in \{a, b\}^*\}$$

# LCFRS and MCFG (7)

## Definition 3 (String Language of an MCFG/LCFRS)

Let  $G = \langle N, T, F, P, S \rangle$  be a MCFG/LCFRS.

- ① For every  $A \in N$ :
  - For every  $A \rightarrow f[ ] \in P$ ,  $f( ) \in \text{yield}(A)$ .
  - For every  $A \rightarrow f[A_1, \dots, A_k] \in P$  with  $k \geq 1$  and all tuples  $\tau_1 \in \text{yield}(A_1), \dots, \tau_k \in \text{yield}(A_k)$ ,  $f(\tau_1, \dots, \tau_k) \in \text{yield}(A)$ .
  - Nothing else is in  $\text{yield}(A)$ .
- ② The string language of  $G$  is  $L(G) = \{w \mid \langle w \rangle \in \text{yield}(S)\}$ .



# LCFRS with Simple RCG syntax (1)

- *Range Concatentation Grammars (RCG)* and the restricted *simple RCG* have been introduced in [Bou00].
- Simple RCG are not only equivalent to MCFG and LCFRS but also represent a useful syntactic variant.

Example: Simple RCG for the double copy language.

$$S(XYZ) \rightarrow A(X, Y, Z)$$

$$A(aX, aY, aZ) \rightarrow A(X, Y, Z)$$

$$A(bX, bY, bZ) \rightarrow A(X, Y, Z)$$

$$A(a, a, a) \rightarrow \varepsilon$$

$$A(b, b, b) \rightarrow \varepsilon$$

## LCFRS with Simple RCG syntax (2)

We redefine LCFRS with the simple RCG syntax:

### Definition 4 (LCFRS)

A LCFRS is a tuple  $G = (N, T, V, P, S)$  where

- 1  $N$ ,  $T$  and  $V$  are disjoint alphabets of non-terminals, terminals and variables resp. with a fan-out function  $dim: N \rightarrow \mathbb{N}$ .  $S \in N$  is the start predicate;  $dim(S) = 1$ .
- 2  $P$  is a finite set of rewriting rules of the form

$$A_0(\vec{\alpha}_0) \rightarrow A_1(\vec{x}_1) \cdots A_m(\vec{x}_m)$$

with  $m \geq 0$ ,  $\vec{\alpha}_0 \in [(T \cup V)^*]^{dim(A_0)}$ ,  $\vec{x}_i \in V^{dim(A_i)}$  for  $1 \leq i \leq m$  and it holds that every variable  $X \in V$  occurring in the rule occurs exactly once in the left-hand side and exactly once in the right-hand side.

## LCFRS with Simple RCG syntax (3)

In order to apply a rule, we have to map variables to strings of terminals:

### Definition 5 (LCFRS rule instantiation)

Let  $G = \langle N, T, V, S, P \rangle$  be a LCFRS.

For a rule  $c = A(\vec{\alpha}) \rightarrow A_1(\vec{\alpha}_1) \dots A_m(\vec{\alpha}_m) \in P$ , every function  $f : \{x \mid x \in V, x \text{ occurs in } c\} \rightarrow T^*$  is an *instantiation* of  $c$ .

We call  $A(f(\vec{\alpha})) \rightarrow A_1(f(\vec{\alpha}_1)) \dots A_m(f(\vec{\alpha}_m))$  then an *instantiated clause* where  $f$  is extended as follows:

- ①  $f(\varepsilon) = \varepsilon$ ;
- ②  $f(t) = t$  for all  $t \in T$ ;
- ③  $f(xy) = f(x)f(y)$  for all  $x, y \in T^*$ ;
- ④  $f(\langle \alpha_1, \dots, \alpha_m \rangle) = (\langle f(\alpha_1), \dots, f(\alpha_m) \rangle)$  for all  $(\langle \alpha_1, \dots, \alpha_m \rangle) \in [(T \cup V)^*]^m, m \geq 1$ .

# LCFRS with Simple RCG syntax (4)

## Definition 6 (LCFRS string language)

Let  $G = \langle N, T, V, S, P \rangle$  be a LCFRS.

- ① The set  $L_{pred}(G)$  of instantiated predicates  $A(\vec{\tau})$  where  $A \in N$  and  $\vec{\tau} \in (T^*)^k$  for some  $k \geq 1$  is defined by the following deduction rules:

$$\text{a) } \frac{}{A(\vec{\tau})} \quad A(\vec{\tau}) \rightarrow \varepsilon \text{ is an instantiated clause}$$

$$\text{b) } \frac{A_1(\vec{\tau}_1) \dots A_m(\vec{\tau}_m)}{A(\vec{\tau})} \quad \begin{array}{l} A(\vec{\tau}) \rightarrow A_1(\vec{\tau}_1) \dots A_m(\vec{\tau}_m) \\ \text{is an instantiated clause} \end{array}$$

- ② The string language of  $G$  is

$$\{w \in T^* \mid S(w) \in L_{pred}(G)\}.$$

# References I

- [AL14] Krasimir Angelov and Peter Ljunglöf.  
Fast statistical parsing with parallel multiple context-free grammars.  
*In Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 368–376, Gothenburg, Sweden, 2014.
- [BB13] Jan A. Botha and Phil Blunsom.  
Adaptor grammars for learning non-concatenative morphology.  
*In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 345–356, Seattle, WA, 2013.
- [Bou00] Pierre Boullier.  
Range Concatenation Grammars.  
*In Proceedings of the Sixth International Workshop on Parsing Technologies (IWPT2000)*, pages 53–64, Trento, Italy, February 2000.
- [Kae13] Miriam Kaeshammer.  
Synchronous linear context-free rewriting systems for machine translation.  
*In Proceedings of the Seventh Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 68–77, Atlanta, GA, 2013.
- [Kae15] Miriam Kaeshammer.  
Hierarchical machine translation with discontinuous phrases.  
*In Proceedings of the Tenth Workshop on Statistical Machine Translation*, Lisbon, Portugal, 2015.
- [Kal10] Laura Kallmeyer.  
*Parsing Beyond Context-Free Grammars*.  
Cognitive Technologies. Springer, Heidelberg, 2010.
- [KM13] Laura Kallmeyer and Wolfgang Maier.  
Data-driven parsing using probabilistic linear context-free rewriting systems.  
*Computational Linguistics*, 39(1):87–119, 2013.

# References II

- [KMPD10] Laura Kallmeyer, Wolfgang Maier, Yannick Parmentier, and Johannes Dellert.  
Tulipa-parsing extensions of tag with range concatenation grammars.  
*Bulletin of the Polish Academy of Sciences: Technical Sciences*, 58(3):377–391, 2010.
- [KP08] Laura Kallmeyer and Yannick Parmentier.  
On the relation between Multicomponent Tree Adjoining Grammars with Tree Tuples (TT-MCTAG) and Range Concatenation Grammars (RCG).  
In *Second International Conference on Language and Automata Theory and Applications (LATA 2008), Revised Papers*, Lecture Notes in Computer Science, pages 263–274. Springer, Tarragona, Spain, 2008.
- [Lju04] Peter Ljunglöf.  
*Expressivity and Complexity of the Grammatical Framework*.  
PhD thesis, Department of Computer Science, Gothenburg University and Chalmers University of Technology, November 2004.
- [Ran11] Aarne Ranta.  
*Grammatical Framework: Programming with Multilingual Grammars*.  
CSLI Publications, Stanford, 2011.
- [SMFK91] Hiroyuki Seki, Takahashi Matsumura, Mamoru Fujii, and Tadao Kasami.  
On multiple context-free grammars.  
*Theoretical Computer Science*, 88(2):191–229, 1991.
- [vC12] Andreas van Cranenburgh.  
Efficient parsing with linear context-free rewriting systems.  
In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 460–470, Avignon, France, 2012.
- [VSWJ87] K. Vijay-Shanker, David J. Weir, and Aravind K. Joshi.  
Characterizing structural descriptions produced by various grammatical formalisms.  
In *Proceedings of ACL*, Stanford, 1987.