

# Parsing Beyond Context-Free Grammars: LCFRS Parsing

Laura Kallmeyer & Tatiana Bladier  
Heinrich-Heine-Universität Düsseldorf

Sommersemester 2018

# Overview

- 1 Ranges
- 2 CYK Parsing
- 3 Incremental Earley Parsing
  - Filters

[Kal10]

# Ranges (1)

- During parsing we have to link the terminals and variables in our LCFRS rules to portions of the input string.
- These can be characterized by their start and end positions.
- A *range* is an pair of indices  $\langle i, j \rangle$  that characterizes the span of a component within the input and a range vector characterizes a tuple in the yield of a non-terminal.
- The range instantiation of a rule specifies the computation of an element from the lefthand side yield from elements of in the yields of the right-hand side non-terminals based on the corresponding range vectors.

## Ranges (2)

Example: Rule  $A(aXa, bYb) \rightarrow B(X)C(Y)$  and input string  $abababcb$ . We assume without loss of generality that our LCFRSs are monotone and  $\varepsilon$ -free. Furthermore, because of the linearity, the components of a tuple in the yield of an LCFRS non-terminal are necessarily non-overlapping. Then, given this input, we have the following possible instantiations for this rule:

$$\begin{array}{l}
 A_{(0\text{aba}_3, 3\text{bab}_6)} \rightarrow B_{(1b_2)}C_{(4a_5)} \quad A_{(0\text{aba}_3, 3\text{babcb}_8)} \rightarrow B_{(1b_2)}C_{(4abc_7)} \\
 A_{(0\text{aba}_3, 5\text{bcb}_8)} \rightarrow B_{(1b_2)}C_{(6c_7)} \quad A_{(0\text{ababa}_5, 5\text{bcb}_8)} \rightarrow B_{(1bab_4)}C_{(6c_7)} \\
 A_{(2\text{aba}_5, 5\text{bcb}_8)} \rightarrow B_{(3b_4)}C_{(6c_7)}
 \end{array}$$

# Ranges (3)

## Definition 1 (Range instantiation, [Bou00])

Let  $G = (N, T, V, P, S)$  be a LCFRS,  $w = t_1 \dots t_n \in T^n$  ( $n \geq 0$ ) and  $r = A(\vec{\alpha}) \rightarrow A_1(\vec{\alpha}_1) \cdots A_m(\vec{\alpha}_m) \in P$  ( $0 \leq m$ ). A **range instantiation** of  $r$  wrt.  $w$  is a function  $f : V \cup \{Eps_i \mid \vec{\alpha}(i) = \varepsilon\} \cup \{t' \mid t' \text{ an occurrence of some } t \in T \text{ in } \vec{\alpha}\} \rightarrow \{\langle i, j \rangle \mid 0 \leq i \leq j \leq n\}$  such that

- for all occurrences  $t'$  of a  $t \in T$  in  $\vec{\alpha}$ ,  $f(t') = \langle i - 1, i \rangle$  for some  $i$  with  $t_i = t$ ,
- for all  $x, y$  adjacent in one of the  $\vec{\alpha}(i)$  there are  $i, j, k$  with  $f(x) = \langle i, j \rangle$ ,  $f(y) = \langle j, k \rangle$ ; we define then  $f(xy) = \langle i, k \rangle$ ,
- for all  $Eps \in \{Eps_i \mid \vec{\alpha}(i) = \varepsilon\}$ ,  $f(Eps) = \langle j, j \rangle$  for some  $j$ ; we define then for every  $\varepsilon$ -argument  $\vec{\alpha}(i)$  that  $f(\vec{\alpha}(i)) = f(Eps_i)$ .

$A(f(\vec{\alpha})) \rightarrow A_1(f(\vec{\alpha}_1)) \cdots A_m(f(\vec{\alpha}_m))$  with

$f(\langle x_1, \dots, x_k \rangle) = \langle f(x_1), \dots, f(x_k) \rangle$  is then called an **instantiated rule**.

# CYK Parsing (1)

First introduced in [SMFK91]; deduction-based definition in, e.g., [KM10].

Idea: Once all elements in the RHS of an instantiated rule have been found, complete the LHS.

- We start with the terminal symbols: whenever we can find a range instantiation of a rule with rhs  $\varepsilon$ , we conclude that this rule can be applied (*scan*).
- We parse bottom-up: whenever, for an instantiated rule, all elements in the rhs have been found, we conclude that this rule can be applied and the lhs of the instantiated rule is deduced (*complete*).
- Our input  $w$  is in the language iff  $S$  with range vector  $\langle\langle 0, n \rangle\rangle$  is in the final set of results that we have deduced.

# CYK Parsing (2)

Deduction rules:

Items  $[A, \vec{\rho}]$  with  $A \in N$ ,  $\vec{\rho}$  is a  $\dim(A)$ -dimensional range vector in  $w$ .

Axioms (scan):  $\frac{}{[A, \vec{\rho}]} A(\vec{\rho}) \rightarrow \varepsilon$  a range instantiated rule

Complete:  $\frac{[A_1, \vec{\rho}_1], \dots, [A_m, \vec{\rho}_m]}{[A, \vec{\rho}]} A(\vec{\rho}) \rightarrow A_1(\vec{\rho}_1), \dots, A_m(\vec{\rho}_m)$   
a range instantiated rule

Goal item:  $[S, \langle\langle 0, n \rangle\rangle]$

# CYK Parsing: Example (1)

Example: MCFG/LCFRS for the double copy language, **input word:** *ababab*

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:

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

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

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

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

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



# CYK Parsing: Example (2)

	Item	Rule
1	$[A, \langle\langle 0, 1 \rangle, \langle 2, 3 \rangle, \langle 4, 5 \rangle\rangle]$	axiom with $A \rightarrow f_4[ ]$
2	$[A, \langle\langle 0, 1 \rangle, \langle 4, 5 \rangle, \langle 2, 3 \rangle\rangle]$	axiom with $A \rightarrow f_4[ ]$
3	$[A, \langle\langle 2, 3 \rangle, \langle 0, 1 \rangle, \langle 4, 5 \rangle\rangle]$	axiom with $A \rightarrow f_4[ ]$
		...
4	$[A, \langle\langle 1, 2 \rangle, \langle 3, 4 \rangle, \langle 5, 6 \rangle\rangle]$	axiom with $A \rightarrow f_5[ ]$
5	$[A, \langle\langle 1, 2 \rangle, \langle 5, 6 \rangle, \langle 3, 4 \rangle\rangle]$	axiom with $A \rightarrow f_5[ ]$
		...
6	$[A, \langle\langle 0, 2 \rangle, \langle 2, 4 \rangle, \langle 4, 6 \rangle\rangle]$	complete, with 4 and $A \rightarrow f_2[A]$
7	$[A, \langle\langle 0, 2 \rangle, \langle 4, 6 \rangle, \langle 2, 4 \rangle\rangle]$	complete, with 5 and $A \rightarrow f_2[A]$
		...
8	$[S, \langle\langle 0, 6 \rangle\rangle]$	complete, with 6 and $S \rightarrow f_1[A]$

# CYK Parsing with binarized LCFRS

Deduction rules for binarized  $\varepsilon$ -free grammars where, without loss of generality, either the lhs contains a single terminal and the rhs is  $\varepsilon$  or the rule contains only variables:

Items and goal as before.

$$\text{Scan: } \frac{}{[A, \langle\langle i, i+1 \rangle\rangle]} A(w_{i+1}) \rightarrow \varepsilon \in P$$

$$\text{Unary: } \frac{[B, \vec{\rho}]}{[A, \vec{\rho}]} A(\vec{\alpha}) \rightarrow B(\vec{\alpha}) \in P$$

$$\text{Binary: } \frac{[B, \vec{\rho}_B], [C, \vec{\rho}_C]}{[A, \vec{\rho}_A]} \quad \begin{array}{l} A(\vec{\rho}_A) \rightarrow B(\vec{\rho}_B)C(\vec{\rho}_C) \\ \text{is a range instantiated rule} \end{array}$$

# CYK Parsing: Another example (1)

LCFRS  $G$  in sRCG format and the [input word](#) *aabbacbbac*.

$G = \langle \{S, A, B, C\}, \{a, b, c\}, \{U, V, W, X, Y, Z\}, P, S \rangle$ , where

$$P = \left\{ \begin{array}{l} S(VYWZX) \rightarrow A(V, W, X)B(Y, Z), \\ A(a, a, a) \rightarrow \epsilon, \\ A(XU, YV, ZW) \rightarrow A(X, Y, Z)C(U, V, W), \\ B(b, b) \rightarrow \epsilon \\ B(XV, WY) \rightarrow B(X, Y)B(V, W) \\ C(a, c, c) \rightarrow \epsilon \end{array} \right\}$$

# CYK Parsing: Another example (2)

	Item	Rule
1	$[A, \langle\langle 0, 1 \rangle, \langle 4, 5 \rangle, \langle 8, 9 \rangle\rangle]$	scan
2	$[A, \langle\langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 4, 5 \rangle\rangle]$	scan
	...	
3	$[B, \langle\langle 2, 3 \rangle, \langle 7, 8 \rangle\rangle]$	scan
4	$[B, \langle\langle 3, 4 \rangle, \langle 6, 7 \rangle\rangle]$	scan
5	$[B, \langle\langle 2, 3 \rangle, \langle 3, 4 \rangle\rangle]$	scan
	...	
6	$[C, \langle\langle 1, 2 \rangle, \langle 5, 6 \rangle, \langle 9, 10 \rangle\rangle]$	scan
7	$[C, \langle\langle 0, 1 \rangle, \langle 5, 6 \rangle, \langle 9, 10 \rangle\rangle]$	scan
8	$[A, \langle\langle 0, 2 \rangle, \langle 4, 6 \rangle, \langle 8, 10 \rangle\rangle]$	complete, with 1 and 6
9	$[B, \langle\langle 2, 4 \rangle, \langle 6, 8 \rangle\rangle]$	complete with 3 and 4
10	$[S, \langle\langle 0, 10 \rangle\rangle]$	complete with 8 and 9

# CYK Parsing: Complexity

Complexity of CYK parsing with binarized LCFRSs:

We have to consider the maximal number of possible applications of the complete rule.

$$\text{Binary: } \frac{[B, \vec{\rho}_B], [C, \vec{\rho}_C]}{[A, \vec{\rho}_A]} \quad \begin{array}{l} A(\vec{\rho}_A) \rightarrow B(\vec{\rho}_B)C(\vec{\rho}_C) \\ \text{is a range instantiated rule} \end{array}$$

If  $k$  is the maximal fan-out in the LCFRS, we have maximal  $2k$  range boundaries in each of the antecedent items of this rule. For variables  $X_1, X_2$  being in the same lhs side argument of the rule,  $X_1$  left of  $X_2$  and no other variables in between, the right boundary of  $X_1$  is the left boundary of  $X_2$ . In the worst case,  $A, B, C$  all have fan-out  $k$  and each lhs argument contains two variables. This gives  $3k$  independent range boundaries and consequently a time complexity of  $\mathcal{O}(n^{3k})$  for the entire algorithm.

# Incremental Earley Parsing

Strategy:

- Process LHS arguments incrementally, starting from an  $S$ -rule
- Whenever we reach a variable, move into rule of corresponding rhs non-terminal (**predict** or **resume**).
- Whenever we reach the end of an argument, **suspend** the rule and move into calling parent rule.
- Whenever we reach the end of the last argument **convert** item into a passive one and **complete** parent item.

This parser is described in [KM09] and inspired by the Thread Automata in [Vil02]

# Incremental Earley Parsing: Items

**Passive items:**  $[A, \vec{\rho}]$  where  $A$  is a non-terminal of fan-out  $k$  and  $\vec{\rho}$  is a range vector of fan-out  $k$

**Active items:**

$$[A(\vec{\phi}) \rightarrow A_1(\vec{\phi}_1) \dots A_m(\vec{\phi}_m), pos, \langle i, j \rangle, \vec{\rho}]$$

where

- $A(\vec{\phi}) \rightarrow A_1(\vec{\phi}_1) \dots A_m(\vec{\phi}_m) \in P$ ;
- $pos \in \{0, \dots, n\}$ : We have reached input position  $pos$ ;
- $\langle i, j \rangle \in \mathbb{N}^2$ : We have reached the  $j$ th element of  $i$ th argument (dot position);
- $\vec{\rho}$  is a range vector containing variable and terminal bindings. All elements are initialized to “?”, an initialized vector is called  $\vec{\rho}_{init}$ .

# Incremental Earley Parsing: Example (1)

$$S(X_1X_2) \longrightarrow A(X_1, X_2) \quad A(aX_1, bX_2) \longrightarrow A(X_1, X_2) \quad A(a, b) \longrightarrow \varepsilon$$

Parsing trace for input  $w = aabb$ :

	pos.	item and dot position	$\vec{\rho}$ (bindings)	rule
1	0	$S(\bullet X_1 X_2) \longrightarrow A(X_1, X_2)$	$(?, ?)$	axiom
2	0	$A(\bullet a X_1, b X_2) \longrightarrow A(X_1, X_2)$	$(?, ?, ?, ?)$	predict, 1
3	0	$A(\bullet a, b) \longrightarrow \varepsilon$	$(?, ?)$	predict, 1
4	1	$A(a \bullet X_1, b X_2) \longrightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, ?, ?, ?)$	scan, 2
5	1	$A(a \bullet, b) \longrightarrow \varepsilon$	$(\langle 0, 1 \rangle, ?)$	scan, 3
6	1	$A(\bullet a X_1, b X_2) \longrightarrow A(X_1, X_2)$	$(?, ?, ?, ?)$	predict, 4
7	1	$A(\bullet a, b) \longrightarrow \varepsilon$	$(?, ?)$	predict 4
8	1	$S(X_1 \bullet X_2) \longrightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, ?)$	susp. 5, 1
9	1	$A(a, \bullet b) \longrightarrow \varepsilon$	$(\langle 0, 1 \rangle, ?)$	resume 5, 8



# Incremental Earley Parsing: Example (2)

	pos.	item and dot position	$\vec{p}$ (bindings)	rule
10	2	$A(a \bullet X_1, bX_2) \rightarrow A(X_1, X_2)$	$(\langle 1, 2 \rangle, ?, ?, ?)$	scan 6
11	2	$A(a \bullet, b) \rightarrow \varepsilon$	$(\langle 1, 2 \rangle, ?)$	scan 7
12	2	$A(\bullet aX_1, bX_2) \rightarrow A(X_1, X_2)$	$(?, ?, ?, ?)$	predict 10
13	2	$A(\bullet a, b) \rightarrow \varepsilon$	$(?, ?)$	predict 10
14	2	$A(aX_1 \bullet, bX_2) \rightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, \langle 1, 2 \rangle, ?, ?)$	susp. 11, 4
15	2	$S(X_1 \bullet X_2) \rightarrow A(X_1, X_2)$	$(\langle 0, 2 \rangle, ?)$	susp. 14, 1
16	2	$A(aX_1, \bullet bX_2) \rightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, \langle 1, 2 \rangle, ?, ?)$	resume 14, 15
17	3	$A(aX_1, b \bullet X_2) \rightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, ?)$	scan 16
18	3	$A(a, \bullet b) \rightarrow \varepsilon$	$(\langle 1, 2 \rangle, ?)$	resume 11, 17

# Incremental Earley Parsing: Example (3)

	pos.	item and dot position	$\vec{\rho}$ (bindings)	rule
19	4	$A(a, b\bullet) \rightarrow \varepsilon$	$(\langle 1, 2 \rangle, \langle 3, 4 \rangle)$	scan 18
20	4	$A(\langle 1, 2 \rangle, \langle 3, 4 \rangle)$		convert 19
21	4	$A(aX_1, bX_2\bullet) \rightarrow A(X_1, X_2)$	$(\langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, \langle 3, 4 \rangle)$	compl. 17, 20
22	4	$A(\langle 0, 2 \rangle, \langle 2, 4 \rangle)$		convert 21
23	4	$S(X_1X_2\bullet) \rightarrow A(X_1, X_2)$	$(\langle 0, 2 \rangle, \langle 2, 4 \rangle)$	compl. 15, 22
24	4	$S(\langle 0, 4 \rangle)$		convert 23

# Incremental Earley Parsing: Deduction Rules

- Notation:
  - $\vec{\rho}(X)$ : range bound to variable  $X$ .
  - $\vec{\rho}(\langle i, j \rangle)$ : range bound to  $j$ th element of  $i$ th argument on LHS.
- Applying a range vector  $\vec{\rho}$  containing variable bindings for given rule  $c$  to the argument vector of the lefthand side of  $c$  means mapping the  $i$ th element in the arguments to  $\vec{\rho}(i)$  and concatenating adjacent ranges. The result is defined iff every argument is thereby mapped to a range.

# Incremental Earley Parsing: Initialize, Goal item

**Initialize:**  $\frac{S(\vec{\phi}) \rightarrow \vec{\Phi} \in P}{[S(\vec{\phi}) \rightarrow \vec{\Phi}, 0, \langle 1, 0 \rangle, \vec{\rho}_{init}]}$

**Goal Item:**  $[S(\vec{\phi}) \rightarrow \vec{\Phi}, n, \langle 1, j \rangle, \psi]$  with  $|\vec{\phi}(1)| = j$  (i.e., dot at the end of lhs argument).

# Incremental Earley Parsing: Scan

If next symbol after dot is next terminal in input, scan it.

$$\text{Scan: } \frac{[A(\vec{\phi}) \rightarrow \vec{\Phi}, pos, \langle i, j \rangle, \vec{\rho}]}{[A(\vec{\phi}) \rightarrow \vec{\Phi}, pos + 1, \langle i, j + 1 \rangle, \vec{\rho}^*]} \vec{\phi}(i, j + 1) = w_{pos+1}$$

where  $\vec{\rho}^*$  is  $\vec{\rho}$  updated with  $\vec{\rho}(\langle i, j + 1 \rangle) = \langle pos, pos + 1 \rangle$ .

# Incremental Earley Parsing: Predict

Whenever our dot is left of a variable that is the first argument of some rhs non-terminal  $B$ , we predict new  $B$ -rules:

$$\text{Predict: } \frac{[A(\vec{\phi}) \rightarrow \dots B(X, \dots) \dots, pos, \langle i, j \rangle, \vec{\rho}_A]}{[B(\vec{\psi}) \rightarrow \vec{\Psi}, pos, \langle 1, 0 \rangle, \vec{\rho}_{init}]}$$

where  $\vec{\phi}(i, j + 1) = X, B(\vec{\psi}) \rightarrow \vec{\Psi} \in P$

# Incremental Earley Parsing: Suspend

## Suspend:

$$\frac{[B(\vec{\psi}) \rightarrow \vec{\Psi}, pos', \langle i, j \rangle, \vec{\rho}_B], [A(\vec{\phi}) \rightarrow \dots B(\vec{\xi}) \dots, pos, \langle k, l \rangle, \vec{\rho}_A]}{[A(\vec{\phi}) \rightarrow \dots B(\vec{\xi}) \dots, pos', \langle k, l + 1 \rangle, \vec{\rho}]}$$

where

- the dot in the antecedent  $A$ -item precedes the variable  $\vec{\xi}(i)$ ,
- $|\vec{\psi}(i)| = j$  ( $i$ th argument has length  $j$ , i.e., is completely processed),
- $|\vec{\psi}| < i$  ( $i$ th argument is not the last argument of  $B$ ),
- $\vec{\rho}_B(\vec{\psi}(i)) = \langle pos, pos' \rangle$
- and for all  $1 \leq m < i$ :  $\vec{\rho}_B(\vec{\psi}(m)) = \vec{\rho}_A(\vec{\xi}(m))$ .

$\vec{\rho}$  is  $\vec{\rho}_A$  updated with  $\vec{\rho}_A(\vec{\xi}(i)) = \langle pos, pos' \rangle$ .

# Incremental Earley Parsing: Convert

Whenever we arrive at the end of the last argument, we convert the item into a passive one:

**Convert:**

$$\frac{[B(\vec{\psi}) \rightarrow \vec{\Psi}, pos, \langle i, j \rangle, \vec{\rho}_B]}{[B, \rho]} \quad |\vec{\psi}(i)| = j, |\vec{\psi}| = i, \vec{\rho}_B(\vec{\psi}) = \rho$$



# Incremental Earley Parsing: Complete

Whenever we have a passive  $B$  item we can use it to move the dot over the variable of the last argument of  $B$  in a parent  $A$ -rule:

**Complete:** 
$$\frac{[B, \vec{\rho}_B], [A(\vec{\phi}) \rightarrow \dots B(\vec{\xi}) \dots, pos, \langle k, l \rangle, \vec{\rho}_A]}{[A(\vec{\phi}) \rightarrow \dots B(\vec{\xi}) \dots, pos', \langle k, l + 1 \rangle, \vec{\rho}]}$$
 where

- the dot in the antecedent  $A$ -item precedes the variable  $\vec{\xi}(|\vec{\rho}_B|)$ ,
- the last range in  $\vec{\rho}_B$  is  $\langle pos, pos' \rangle$ ,
- and for all  $1 \leq m < |\vec{\rho}_B|$ :  $\vec{\rho}_B(m) = \vec{\rho}_A(\vec{\xi}(m))$ .

$\vec{\rho}$  is  $\vec{\rho}_A$  updated with  $\vec{\rho}_A(\vec{\xi}(|\vec{\rho}_B|)) = \langle pos, pos' \rangle$ .

# Incremental Earley Parsing: Resume

Whenever we are left of a variable that is not the first argument of one of the rhs non-terminals, we resume the rule of the rhs non-terminal.

$$\text{Resume: } \frac{[A(\vec{\phi}) \rightarrow \dots B(\vec{\xi}) \dots, pos, \langle i, j \rangle, \vec{\rho}_A], [B(\vec{\psi}) \rightarrow \vec{\Psi}, pos', \langle k-1, l \rangle, \vec{\rho}_B]}{[B(\vec{\psi}) \rightarrow \vec{\Psi}, pos, \langle k, 0 \rangle, \vec{\rho}_B]}$$

where

- $\vec{\phi}(i, j+1) = \vec{\xi}(k), k > 1$  (the next element is a variable that is the  $k$ th element in  $\vec{\xi}$ , i.e., the  $k$ th argument of  $B$ ),
- $|\vec{\psi}(k-1)| = l$ , and
- $\vec{\rho}_A(\vec{\xi}(m)) = \vec{\rho}_B(\vec{\psi}(m))$  for all  $1 \leq m \leq k-1$ .

# Incremental Earley Parsing: Filters

- Filters can be applied to decrease the number of items in the chart
- A filter is an additional condition on the form of items.
- E.g., in a  $\varepsilon$ -free grammar, the number of variables in the part of the lefthand side arguments of a rule that has not been processed yet must be lower or equal to the length of the remaining input.

We will discuss in the following some filters that are particularly useful when dealing with natural languages.

# Incr. Earley Parsing: Remaining Input Length Filter

- In  $\varepsilon$ -free grammars each variable must cover at least one input symbol.
- $i$  input symbols left implies no prediction of a clause with more than  $i$  variables or terminals on LHS since no instantiation is possible
- Condition on active items, can be applied with predict, resume, suspend and complete

The length of the remaining input must be  $\geq$  the number of variables and terminal occurrences to the right of the dot in the lefthand side of the clause, i.e.

An item  $[A(\vec{\phi}) \rightarrow A_1(\vec{\phi}_1) \dots A_m(\vec{\phi}_m), pos, \langle i, j \rangle, \vec{\rho}]$  satisfies the **length filter** iff

$$(n - pos) \geq (|\vec{\phi}(i)| - j) + \sum_{k=i+1}^{dim(A)} |\vec{\phi}(k)|$$

# Incremental Earley Parsing: Preterminal Filter (1)

- Check for the presence of (pre)terminals in the predicted part of a clause (the part to the right of the dot) in the remaining input, and
- check that terminals appear in the predicted order and that distance between two of them is at least the number of variables/terminals in between.

continued...

## Incremental Earley Parsing: Preterminal Filter (2)

In other words, an active item  $[A(\vec{\phi}) \rightarrow A_1(\vec{\phi}_1) \dots A_m(\vec{\phi}_m), pos, \langle i, j \rangle, \vec{p}]$  satisfies the **preterminal filter** iff we can find an injective mapping  $f_T : \text{Term} = \{\langle k, l \rangle \mid \vec{\phi}(k, l) \in T \text{ and either } k > i \text{ or } (k = i \text{ and } l > j)\} \rightarrow \{pos + 1, \dots, n\}$  such that

- ①  $w_{f_T(\langle k, l \rangle)} = \vec{\phi}(k, l)$  for all  $\langle k, l \rangle \in \text{Term}$ ;
- ② for all  $\langle k_1, l_1 \rangle, \langle k_2, l_2 \rangle \in \text{Term}$  with  $k_1 = k_2$  and  $l_1 < l_2$ :  
 $f_T(\langle k_2, l_2 \rangle) \geq f_T(\langle k_1, l_1 \rangle) + (l_2 - l_1)$ ;
- ③ for all  $\langle k_1, l_1 \rangle, \langle k_2, l_2 \rangle \in \text{Term}$  with  $k_1 < k_2$ :  
 $f_T(\langle k_2, l_2 \rangle) \geq f_T(\langle k_1, l_1 \rangle) + (|\vec{\phi}(k_1)| - l_1) + \sum_{k=k_1+1}^{k_2-1} |\vec{\phi}(k)| + l_2$ .

Checking this filtering condition amounts to a linear traversal of the part of the lefthand side of the clause that is to the right of the dot. We start with index  $i = pos + 1$ , for every variable or gap we increment  $i$  by 1. For every terminal  $a$ , we search the next  $a$  in the input, starting at position  $i$ . If it occurs at position  $j$ , then we set  $i = j$  and continue our traversal of the remaining parts of the lefthand side of the clause.

# References I

- [Bou00] Pierre Boullier.  
Range Concatenation Grammars.  
In *Proceedings of the Sixth International Workshop on Parsing Technologies (IWPT2000)*, pages 53–64, Trento, Italy, February 2000.
- [Kal10] Laura Kallmeyer.  
*Parsing Beyond Context-Free Grammars*.  
Cognitive Technologies. Springer, Heidelberg, 2010.
- [KM09] Laura Kallmeyer and Wolfgang Maier.  
An incremental Earley parser for simple Range Concatenation Grammar.  
In *Proceedings of IWPT 2009*, 2009.
- [KM10] Laura Kallmeyer and Wolfgang Maier.  
Data-driven parsing with probabilistic Linear Context-Free Rewriting Systems.  
In *Proceedings of the 23rd International Conference on Computational Linguistics (COLING 2010)*, Beijing, China, 2010.
- [SMFK91] Hiroyuki Seki, Takahashi Matsumura, Mamoru Fujii, and Tadao Kasami.  
On multiple context-free grammars.  
*Theoretical Computer Science*, 88(2):191–229, 1991.
- [Vil02] Éric Villemonte de La Clergerie.  
Parsing mildly context-sensitive languages with thread automata.  
In *Proc. of COLING'02*, August 2002.