

---

# Linear Context-Free Rewriting Systems: Thread Automata for Simple RCG

Laura Kallmeyer  
University of Tübingen  
Summer term 2010

---

Thread Automata 1 25 June 2010

Kallmeyer LCFRS

---

## Overview

1. Idea of Thread Automata
2. TA for Simple RCG
3. Example
4. General Definition of TA

---

Thread Automata 2 25 June 2010

---

## Idea of Thread Automata (1)

- Thread automata (TA) have been proposed in [Villemonde de La Clergerie, 2002].
- TA accept (at least) the class of all LCFRLs, maybe even a proper superset.
- TA are non-deterministic and, if all possibilities are pursued independently from each other, they are of exponential complexity.
- However, in combination with a compact representation of sub-derivations as items and with tabulation techniques, they become polynomial.

---

Thread Automata 3 25 June 2010

Kallmeyer LCFRS

---

## Idea of Thread Automata (2)

Overall idea of TA:

- We have a set of threads, one of which is the active thread.
- Each thread has a unique path that situates it within the tree-shaped thread structure.
- Whenever a new thread is started, its path is a concatenation of the parent thread path and a new symbol. This way, from a given active thread, we can always find its parent thread (the one that started it) and its daughter threads.
- The moves of the automaton are the following: We can change the content of the active thread, start a new daughter thread or move into an existing daughter thread, or move into the parent thread while, eventually, terminating the active thread.

---

Thread Automata 4 25 June 2010

---

### Idea of Thread Automata (3)

Example: Consider the following simple RCG for  $\{a^n b^n c^n d^n \mid n \geq 0\}$ :

$$\begin{aligned} S(XeY) &\rightarrow A(X, Y) & S(e) &\rightarrow \varepsilon \\ A(aXb, cYd) &\rightarrow A(X, Y) & A(ab, cd) &\rightarrow \varepsilon \end{aligned}$$

In the corresponding TA,

- we would start with an  $S$  thread.
- From here, we can read an  $e$  and then terminate.
- Or we start a daughter  $A$  thread that is suspended once the first component is finished.
- Then we continue the mother  $S$  thread, read the  $e$  and then resume the daughter.

### Idea of Thread Automata (4)

$[1 : S(\bullet X e Y) \rightarrow A(X, Y)]$	
$[1 : \dots], [11 : A(\bullet a X b, c Y d) \rightarrow A(X, Y)]$	predict
$[1 : \dots], [11 : A(a \bullet X b, c Y d) \rightarrow A(X, Y)]$	scan
$[1 : \dots], [11 : \dots], [111 : A(\bullet ab, cd) \rightarrow \varepsilon]$	predict
$[1 : \dots], [11 : \dots], [111 : A(ab \bullet, cd) \rightarrow \varepsilon]$	scan (twice)
$[1 : \dots], [11 : A(a X \bullet b, c Y d) \rightarrow A(X, Y)], [111 : \dots]$	suspend
$[1 : \dots], [11 : A(a X b \bullet, c Y d) \rightarrow A(X, Y)], [111 : \dots]$	scan
$[1 : S(X e \bullet Y) \rightarrow A(X, Y)], [11 : \dots], [111 : \dots]$	suspend, scan
$[1 : \dots], [11 : A(a X b, \bullet c Y d) \rightarrow A(X, Y)], [111 : \dots]$	resume
$[1 : \dots], [11 : A(a X b, c \bullet Y d) \rightarrow A(X, Y)], [111 : \dots]$	scan
$[1 : \dots], [11 : \dots], [111 : A(ab, cd \bullet) \rightarrow \varepsilon]$	resume, scan (twice)
$[1 : \dots], [11 : A(a X b, c Y \bullet d) \rightarrow A(X, Y)]$	suspend, scan
$[1 : S(X e Y \bullet) \rightarrow A(X, Y)]$	suspend

---

### Idea of Thread Automata (5)

- TA for simple RCGs perform a top-down recognition.
- If an additional tabulation is included, the automaton amounts to an incremental Earley recognition.
- [Villemonde de La Clergerie, 2002] has implemented TA in general.
- [Kallmeyer and Maier, 2009] implements an incremental Earley chart parsing of ordered simple RCG following the same strategy.

### TA for Simple RCG (1)

[Villemonde de La Clergerie, 2002] gives a general definition of TA and shows then how to construct equivalent TA for Tree Adjoining Grammars and for simple RCGs.

In the following, we first introduce the specific TA for ordered simple RCGs that we call *SRCG-TA*. Only later, general TA are defined.

We call a simple RCG rule with a dot in the lhs a dotted rule.

---

### TA for Simple RCG (2)

**Definition 1 (Thread Automaton for Simple RCG)** A thread automaton for an ordered simple RCG

$G = \langle N_G, T_G, S_G, P_G \rangle$  (SRCG-TA) is a tuple  $\langle N, T, S', ret, \mathcal{U}, \Theta \rangle$  where

- $N = N_G \cup \{S', ret\} \cup \{\gamma \mid \gamma \text{ is a dotted rule in } G\}$ ,  $T = T_G$  are the non-terminals and terminals with  $S', ret \in N \setminus N_G$  the start and end symbols;
- $\mathcal{U} = \{1, \dots, m\}$  where  $m$  the maximal rhs length in  $G$  is the set of labels used to identify threads.
- $\Theta$  is a finite set of transitions.

### TA for Simple RCG (3)

A configuration is a tree-shaped set of threads, one of them being the active thread, together with a position in the input that separates the part that has been recognized from the remaining part of the input.

**Definition 2 (Thread, Configuration)** Let  $M = \langle N, T, S', ret, \mathcal{U}, \Theta \rangle$  be a SRCG-TA.

- A thread is a pair  $p : A$  with  $p \in \mathcal{U}^*$ ,  $A \in N$ .  $p$  is the thread path, and  $A$  is the content of the thread.
- A thread store is a set of threads whose addresses are closed by prefix.
- A configuration of  $M$  is a tuple  $\langle i, p, \mathcal{S} \rangle$  where  $i$  is an input position,  $\mathcal{S}$  is a thread set and  $p$  is a thread path in  $\mathcal{S}$ .

---

### TA for Simple RCG (4)

The transitions defined within  $\Theta$  are the following:

- **Call** starts a new thread, either for the start predicate or for a daughter predicate:  
If active thread  $[p : S']$ , then add new thread  $[p1 : S]$  (where  $S$  start symbol of  $G$ ), set active thread to  $p1$ .  
If active thread  $[p : \gamma]$ ,  $\gamma$  a dotted rule with the dot preceding the first variable of the rhs non-terminal  $A$  and  $A$  is the  $i$ th rhs element, then add new thread  $[pi : A]$  and set active thread to  $pi$ .
- **Predict** predicts a new clause for a non-terminal:  
If active thread  $[p : A]$  with  $A \in N_G$  and if  $\gamma$  is a dotted  $A$ -rule with the dot at the leftmost position, then replace active thread with  $[p : \gamma]$ .

### TA for Simple RCG (5)

- **Scan** moves the dot over a terminal in the left-hand side while scanning the next input symbol:  
If input position  $i$  and active thread  $[p : \gamma]$  where  $\gamma$  a dotted rule such that the dot precedes a terminal that is the  $(i + 1)$ st input symbol, then move dot over this terminal in active thread and increment input position.
- **Publish** marks the end of a predicate:  
If active thread  $[p : \gamma]$  where  $\gamma$  a dotted rule such that the dot is at the end of the lhs, then replace active thread with  $[p : ret]$ .

---

### TA for Simple RCG (6)

- **Suspend** suspends a daughter thread and resumes the parent:
  1. If  $[pi : ret]$  is the active thread and in its parent thread  $[p : \gamma]$  the dot precedes the last variable of the  $i$ th rhs element, then remove active thread, move dot over this variable in  $p$  thread and set active thread to  $p$ .
  2. If  $[pi : \gamma]$  is the active thread with the dot at the end of the  $j$ th lhs component, the  $j$ th component is not the last, and if the parent thread is  $[p : \beta]$  with the dot preceding the variable that is the  $j$ th argument of the  $i$ th rhs element, then move the dot over this variable in the  $p$  thread and set active thread to  $p$ .

### TA for Simple RCG (7)

- **Resume** resumes an already present daughter thread:

If active thread is  $[p : \gamma]$  where the dot precedes a variable that is the  $j$ th ( $j > 1$ ) argument of the  $i$ th rhs element and if  $pi : \beta]$  is the daughter thread where the dot is at the end of the  $(j - 1)$ th lhs argument, then move the dot to the beginning of the  $j$ th lhs argument in the  $pi$  thread and set active thread to  $pi$ .

---

### TA for Simple RCG (8)

The set of possible configurations  $C(M, w)$  for a given input SRCG-TA  $M$  and a given input  $w$  contains all configurations that are reachable from  $\langle 0, \varepsilon, \{\varepsilon : S'\} \rangle$  via the reflexive transitive closure of the transitions.

The language of a SRCG-TA is the set of words that allow us, starting from the initial thread set  $\{\varepsilon : S'\}$ , to reach the set  $\{\varepsilon : S', 1 : ret\}$  after having scanned the entire input.

**Definition 3 (Language)** Let  $M = \langle N, T, S', ret, \mathcal{U}, \Theta \rangle$  be a SRCG-TA. The language of  $M$  is defined as follows:

$$L(M) = \{w \mid \langle |w|, 1, \{\varepsilon : S', 1 : ret\} \rangle \in C(M, w)\}.$$

### Example (1)

Ordered simple RCG with the following rules:

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

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

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

We encode dotted rules as  $r_{i,j}$  where  $r$  the name of the rule,  $\langle i, j \rangle$  the position of the dot: the dot precedes the  $j$  element of the  $i$ th argument of the lhs.

Ex.:  $\beta_{2,0}$  encodes  $A(aX, \bullet aY, aZ) \rightarrow A(X, Y, Z)$

Input  $w = ababab$ .

---

**Example (2)**

thread set	rem. input	
$\varepsilon : S'$	<i>ababab</i>	
$\varepsilon : S', \mathbf{1} : S$	<i>ababab</i>	call
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}$	<i>ababab</i>	predict
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \mathbf{A}$	<i>ababab</i>	call
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,0}$	<i>ababab</i>	predict
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,1}$	<i>babab</i>	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,1}, \mathbf{111} : \mathbf{A}$	<i>babab</i>	call
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,1}, \mathbf{111} : \gamma_{1,0}$	<i>babab</i>	predict
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,1}, \mathbf{111} : \gamma_{1,1}$	<i>abab</i>	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,0}, \mathbf{11} : \beta_{1,2}, \mathbf{111} : \gamma_{1,1}$	<i>abab</i>	suspend
$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{1,2}, \mathbf{111} : \gamma_{1,1}$	<i>abab</i>	suspend

**Example (3)**

$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{2,0}, \mathbf{111} : \gamma_{1,1}$	<i>abab</i>	resume
$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{2,1}, \mathbf{111} : \gamma_{1,1}$	<i>bab</i>	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{2,1}, \mathbf{111} : \gamma_{2,0}$	<i>bab</i>	resume
$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{2,1}, \mathbf{111} : \gamma_{2,1}$	<i>ab</i>	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,1}, \mathbf{11} : \beta_{2,2}, \mathbf{111} : \gamma_{2,1}$	<i>ab</i>	suspend
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{2,2}, \mathbf{111} : \gamma_{2,1}$	<i>ab</i>	suspend
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,0}, \mathbf{111} : \gamma_{2,1}$	<i>ab</i>	resume
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,1}, \mathbf{111} : \gamma_{2,1}$	<i>b</i>	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,1}, \mathbf{111} : \gamma_{3,0}$	<i>b</i>	resume
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,1}, \mathbf{111} : \gamma_{3,1}$	$\varepsilon$	scan
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,1}, \mathbf{111} : \mathbf{ret}$	$\varepsilon$	publish
$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \beta_{3,2}$	$\varepsilon$	suspend

---

**Example (4)**

$\varepsilon : S', \mathbf{1} : \alpha_{1,2}, \mathbf{11} : \mathbf{ret}$	$\varepsilon$	publish
$\varepsilon : S', \mathbf{1} : \alpha_{1,3}$	$\varepsilon$	suspend
$\varepsilon : S', \mathbf{1} : \mathbf{ret}$	$\varepsilon$	publish

Input accepted with the last configuration.

**General Definition of TA (1)**

We will now give the general definition of TA.

**Definition 4 (Thread Automaton)** A Thread Automaton is a tuple  $\langle N, T, S, F, \kappa, \mathcal{K}, \delta, \mathcal{U}, \Theta \rangle$  where

- $N$  and  $T$  are non-terminal and terminal alphabets with  $S, F \in N$  the start and end symbols;
- $\kappa$ , the triggering function, is a partial function from  $N$  to some finite set  $\mathcal{K}$
- $\mathcal{U}$  is a finite set of labels used to identify threads.
- $\delta$  is a partial function from  $N$  to  $\mathcal{U} \cup \{\perp\}$  used to specify daughter threads that can be created or resumed at some point.
- $\Theta$  is a finite set of transitions.

---

### General Definition of TA (2)

In the simple RCG case,  $\kappa$  and  $\delta$  are used to indicate, for the dot preceding a given variable in the lhs, which of the rhs elements contains this variable as an argument. This determines the daughter thread that processes this variable. Furthermore, when starting a new daughter thread,  $\kappa$  indicates the corresponding non-terminal. Therefore,  $\mathcal{K} = N \cup \{\text{void}\}$  and

- $\kappa(\gamma_{k,i}) = A$  and  $\delta(\gamma_{k,i}) = j$  if  $A$  is the  $j$ th predicate in the rhs of  $\gamma$  and the dot precedes the first argument of  $A$ ,
- $\kappa(\gamma_{k,i}) = \text{void}$  and  $\delta(\gamma_{k,i}) = j$  if the dot precedes an argument of  $A$  that is not its first argument,
- $\kappa(\gamma_{k,i}) = \text{void}$  and  $\delta(\gamma_{k,i}) = \perp$  if the dot is at the end of the  $k$ th argument and, instead of moving into a daughter thread, we have to suspend this thread and resume the parent.

### General Definition of TA (3)

#### Definition 5 (TA Transitions)

Let  $M = \langle N, T, S, F, \kappa, \mathcal{K}, \delta, \mathcal{U}, \Theta \rangle$  be a TA. All transitions in  $\Theta$  have one of the following forms:

- $B \xrightarrow{\alpha} C$  with  $B, C \in N, \alpha \in T^*$  (**SWAP** operation)
- $b \rightarrow [b]C$  with  $b \in \mathcal{K}, C \in N$  (**PUSH** operation)
- $[B]C \rightarrow D$  with  $B, C, D \in N$  (**POP** operation)
- $b[C] \rightarrow [b]D$  with  $b \in \mathcal{K}, C, D \in N$  (**SPUSH** operation)
- $[B]c \rightarrow D[c]$  with  $c \in \mathcal{K}, B, D \in N$  (**SPOP** operation)

---

### General Definition of TA (4)

The set of configurations for  $w$ ,  $C(M, w)$ , is then defined by the following deduction rules:

- Initial configuration:  $\frac{}{\langle 0, \varepsilon, \{\varepsilon : S\} \rangle}$
- Swap:  $\frac{\langle i, p, \mathcal{S} \cup p : B \rangle}{\langle i + |\alpha|, p, \mathcal{S} \cup p : C \rangle} B \xrightarrow{\alpha} C, w_{i+1} \dots w_{i+|\alpha|} = \alpha$
- Push:  $\frac{\langle i, p, \mathcal{S} \cup p : B \rangle}{\langle i, pu, \mathcal{S} \cup p : B \cup pu : C \rangle} b \rightarrow [b]C, \kappa(B) = b, \delta(B) = u, pu \notin \text{dom}(\mathcal{S})$

### General Definition of TA (5)

- Pop:  $\frac{\langle i, pu, \mathcal{S} \cup p : B \cup pu : C \rangle}{\langle i, p, \mathcal{S} \cup p : D \rangle} [B]C \rightarrow D, \delta(C) = \perp, pu \notin \text{dom}(\mathcal{S})$
- Spush:  $\frac{\langle i, p, \mathcal{S} \cup p : B \cup pu : C \rangle}{\langle i, pu, \mathcal{S} \cup p : B \cup pu : D \rangle} b[C] \rightarrow [b]D, \kappa(B) = b, \delta(B) = u$
- Spop:  $\frac{\langle i, pu, \mathcal{S} \cup p : B \cup pu : C \rangle}{\langle i, p, \mathcal{S} \cup p : D \cup pu : C \rangle} [B]c \rightarrow D[c], \kappa(C) = c, \delta(C) = \perp$

---

### General Definition of TA (6)

The language of a TA is the set of words that allow us, starting from the initial thread set  $\{\varepsilon : S\}$ , to reach the set  $\{\varepsilon : S, \delta(S) : F\}$  after having scanned the entire input.

**Definition 6 (Language)** Let  $M = \langle N, T, S, F, \kappa, \mathcal{K}, \delta, \mathcal{U}, \Theta \rangle$  be a TA,

The language of  $M$  is defined as follows:

$$L(M) = \{w \mid \langle n, \delta(S), \{\varepsilon : S, \delta(S) : F\} \rangle \in C(M, w)\}.$$

### General Definition of TA (7)

Take again the TA equivalent to the simple RCG with the following clauses:

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

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

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

Transitions of the corresponding TA (start symbol  $S'$ ):

$$\text{Call: } S' \rightarrow [S']S \quad \alpha_{1,0} \rightarrow [\alpha_{1,0}]A \quad \beta_{1,1} \rightarrow [\beta_{1,1}]A$$

$$\text{Predict: } S \rightarrow \alpha_{1,0} \quad A \rightarrow \beta_{1,0} \quad A \rightarrow \gamma_{1,0}$$

$$\text{Scan: } \beta_{1,0} \xrightarrow{a} \beta_{1,1} \quad \beta_{2,0} \xrightarrow{a} \beta_{2,1} \quad \beta_{3,0} \xrightarrow{a} \beta_{3,1}$$

$$\gamma_{1,0} \xrightarrow{b} \gamma_{1,1} \quad \gamma_{2,0} \xrightarrow{b} \gamma_{2,1} \quad \gamma_{3,0} \xrightarrow{b} \gamma_{3,1}$$

---

### General Definition of TA (8)

Suspend:

$$[\alpha_{1,0}]\beta_{1,2} \rightarrow \alpha_{1,1}[\beta_{1,2}] \quad [\alpha_{1,1}]\beta_{2,2} \rightarrow \alpha_{1,2}[\beta_{2,2}] \quad [\alpha_{1,2}]ret \rightarrow \alpha_{1,3}$$

$$[\alpha_{1,0}]\gamma_{1,1} \rightarrow \alpha_{1,1}[\gamma_{1,1}] \quad [\alpha_{1,1}]\gamma_{2,1} \rightarrow \alpha_{1,2}[\gamma_{2,1}]$$

$$[\beta_{1,1}]\beta_{1,2} \rightarrow \beta_{1,2}[\beta_{1,2}] \quad [\beta_{2,1}]\beta_{2,2} \rightarrow \beta_{2,2}[\beta_{2,2}] \quad [\beta_{3,1}]ret \rightarrow \beta_{3,2}$$

$$[\beta_{1,1}]\gamma_{1,1} \rightarrow \beta_{1,2}[\gamma_{1,1}] \quad [\beta_{2,1}]\gamma_{2,1} \rightarrow \beta_{2,2}[\gamma_{2,1}]$$

Resume

$$\alpha_{1,1}[\beta_{1,2}] \rightarrow [\alpha_{1,1}]\beta_{2,0} \quad \beta_{2,1}[\beta_{1,2}] \rightarrow [\beta_{2,1}]\beta_{2,0}$$

$$\alpha_{1,1}[\gamma_{1,0}] \rightarrow [\alpha_{1,1}]\gamma_{2,0} \quad \beta_{2,1}[\gamma_{1,0}] \rightarrow [\beta_{2,1}]\gamma_{2,0}$$

$$\alpha_{1,2}[\beta_{2,2}] \rightarrow [\alpha_{1,2}]\beta_{3,0} \quad \beta_{2,1}[\beta_{2,2}] \rightarrow [\beta_{2,1}]\beta_{3,0}$$

$$\alpha_{1,2}[\gamma_{2,0}] \rightarrow [\alpha_{1,2}]\gamma_{3,0} \quad \beta_{3,1}[\gamma_{2,0}] \rightarrow [\beta_{3,1}]\gamma_{3,0}$$

Publish:

$$\alpha_{1,3} \rightarrow ret$$

$$\beta_{3,2} \rightarrow ret$$

$$\gamma_{3,1} \rightarrow ret$$

## References

- [Kallmeyer and Maier, 2009] Kallmeyer, L. and Maier, W. (2009). An incremental early parser for simple range concatenation grammar. In *Proceedings of IWPT 2009*.
- [Villemonde de La Clergerie, 2002] Villemonde de La Clergerie, E. (2002). Parsing mildly context-sensitive languages with thread automata. In *Proc. of COLING'02*.