

PARSING MILDLY CONTEXT-SENSITIVE RMS

Tilman Becker Dominik Heckmann

German Research Center for Artificial Intelligence (DFKI GmbH)

{becker,heckmann}@dfki.de

Abstract

We introduce Recursive Matrix Systems (RMS) which encompass mildly context-sensitive formalisms and present efficient parsing algorithms for linear and context-free variants of RMS. The time-complexities are $O(n^{2h+1})$, and $O(n^{3h})$ respectively, where h is the height of the matrix. It is possible to represent Tree Adjoining Grammars (TAG [1], MC-TAG [2], and R-TAG [3]) as RMS uniformly.

1 Recursive Matrix Systems (RMS)

$RMS = (G, I)$ is a two-step formalism. In a first step, a grammar $G = (N, T, S, P)$ generates a set of recursive matrices. In a second step, a yield function maps, according to the interpretation I , the recursive matrices to a set of strings $L(G, I)$. The recursive matrix generating grammars are ordinary phrase structure grammars, where the terminal symbols are replaced by vertical vectors. These vertical vectors are filled with terminal symbols, nonterminal symbols, or they are left empty.

The two-step formalism: $S \in N \xrightarrow{*}_G \text{derive} m \in \mathcal{M} \xrightarrow{I} \text{yield} w \in T^*$.

Example: $G_1 = (N, T, S, P) = (\{S\}, \{a, b, c\}, S, \{S \rightarrow \begin{smallmatrix} a \\ b \\ c \end{smallmatrix} S, S \rightarrow \begin{smallmatrix} a \\ b \\ c \end{smallmatrix}\})$. Every time the first rule is applied, a new column with terminals is added to the matrix. The last rule terminates the process.

A possible derivation: $S \xrightarrow{1}_{G_1} \begin{smallmatrix} a \\ b \\ c \end{smallmatrix} S \xrightarrow{1}_{G_1} \begin{smallmatrix} a & a \\ b & b \\ c & c \end{smallmatrix} S \xrightarrow{1}_{G_1} \begin{smallmatrix} a & a & a \\ b & b & b \\ c & c & c \end{smallmatrix}$.

The product of the derivation process is a matrix with terminals as elements. Finally, these terminal symbols are combined into one string. There are many possible interpretation functions. However, it seems reasonable to read the terminal symbols within the matrix row by row from top to bottom, and each row alternating from left to right and from right to left. This interpretation condition for height 3 could be visualized by $\overleftrightarrow{\updownarrow}$. The grammar G_1 together with this interpretation generates strings of the form $aaa...bbb...ccc...$. Thus, the generated language is $L(G_1, \overleftrightarrow{\updownarrow}) = \{a^n b^n c^n \mid n \in \mathbb{N}\}$.

Definition: h is a positive integer. A *recursive matrix grammar* is a four-tuple $G = (N, T, S, P)$, where N and T are disjoint alphabets, $S \in N$, P is a finite set of ordered pairs $(u, v) \in N \times (N \cup C_h)^*$, where G is *context-free*, if $P \subset N \times (N \cup C_h)^*$, *linear*, if $P \subset N \times C_h N^e C_h^e$, *regular*, if $P \subset N \times C_h N^e$.

$C_h = \left\{ \begin{smallmatrix} v_1 \\ \vdots \\ v_h \end{smallmatrix} \mid \forall_{1 \leq x \leq h} : v_x \in V^e \right\}$. C_h is the set of *columns* of height h over $V^e = N \cup T \cup \{\epsilon\}$.

An *interpretation* $I \in \mathcal{I}_h$ is a vertical vector with left and right arrows as elements, where

$\mathcal{I}_h = \left\{ \begin{smallmatrix} d_1 \\ \vdots \\ d_h \end{smallmatrix} \mid \forall_{1 \leq x \leq h} : d_x \in \{\rightarrow, \leftarrow\} \right\}$. \mathcal{I}_h is the *set of Interpretations* of height h .

2 Parsing Schemata for Recursive Matrix Systems (RMS)

The notation for the parsing schemata is close to [4]. An Earley parser for RMS can be found in [2].

RMS offer brief vector notations: Let $\vec{i} = \begin{matrix} i_1 \\ \vdots \\ i_h \end{matrix}$, $\vec{i}' = \begin{matrix} i_2 \\ \vdots \\ i_{h+1} \end{matrix}$, $\vec{j} = \begin{matrix} j_1 \\ \vdots \\ j_h \end{matrix}$, $\vec{k} = \begin{matrix} k_1 \\ \vdots \\ k_h \end{matrix}$, $\vec{v} = \begin{matrix} v_1 \\ \vdots \\ v_h \end{matrix}$.

The set of hypothesis: $\mathcal{H} = \{ \boxed{a \mid i-1 \mid i} \mid a = a_i, 1 \leq i \leq n \} \cup \{ \boxed{\epsilon \mid i \mid i} \mid 0 \leq i \leq n \}$.

The set of items \mathcal{J} depends on the height h : $\mathcal{J}(h) = \{ \boxed{A \mid i \mid j}, \boxed{A \mid \vec{i} \mid \vec{j}} \}$.

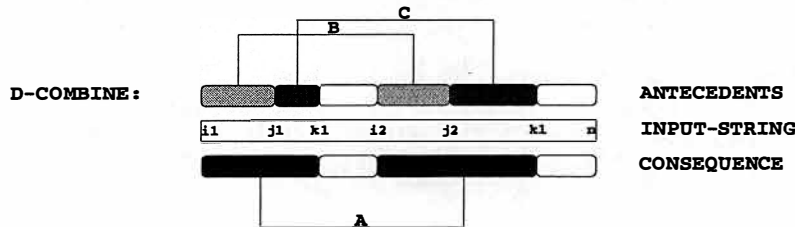
The parsing schema $\mathbb{P}(CYK)(hRMS(CFG, \rightrightarrows) \equiv MC-TAG)$ is defined by the deduction rules:

$$\begin{aligned} D^{init} &= \{ \boxed{v_1 \mid i_1 \mid k_1}, \dots, \boxed{v_h \mid i_h \mid k_h} \vdash \boxed{A \mid \vec{i} \mid \vec{k}} \mid (A \rightarrow \vec{v}) \in P \} \\ D^{combine} &= \{ \boxed{B \mid \vec{i} \mid \vec{j}}, \boxed{C \mid \vec{j} \mid \vec{k}} \vdash \boxed{A \mid \vec{i} \mid \vec{k}} \mid (A \rightarrow BC) \in P \} \\ D^{linearize} &= \{ \boxed{A \mid \vec{i} \mid \vec{i}'} \vdash \boxed{A \mid i_1 \mid i_{h+1}} \} \end{aligned}$$

The parsing schema $\mathbb{P}(CYK)(2RMS(LIN, \rightrightarrows) \equiv R-TAG)$ is defined by the deduction rules:

$$\begin{aligned} D^{init} &= \left\{ \begin{array}{|c|c|c|} \hline v_1 & i_1 & j_2 \\ \hline v_3 & i_3 & j_4 \\ \hline \end{array} \vdash \begin{array}{|c|c|c|} \hline A & i_1 & j_2 \\ \hline & i_3 & j_4 \\ \hline \end{array} \mid (A \rightarrow \begin{array}{|c|} \hline v_1 \\ \hline v_3 \\ \hline \end{array}) \in P \right\} \\ D^{adjoin} &= \left\{ \begin{array}{|c|c|c|} \hline v_1 & i_1 & j_1 \\ \hline v_3 & i_3 & j_3 \\ \hline \end{array} \begin{array}{|c|c|c|} \hline B & j_1 & i_2 \\ \hline & j_3 & i_4 \\ \hline \end{array} \begin{array}{|c|c|c|} \hline v_2 & i_2 & j_2 \\ \hline v_4 & i_4 & j_4 \\ \hline \end{array} \vdash \begin{array}{|c|c|c|} \hline A & i_1 & j_2 \\ \hline & i_3 & j_4 \\ \hline \end{array} \mid (A \rightarrow \begin{array}{|c|} \hline v_1 \\ \hline v_3 \\ \hline \end{array} B \begin{array}{|c|} \hline v_2 \\ \hline v_4 \\ \hline \end{array}) \in P \right\} \\ D^{linearize} &= \left\{ \begin{array}{|c|c|c|} \hline A & i_1 & i_2 \\ \hline & i_2 & i_3 \\ \hline \end{array} \vdash \begin{array}{|c|c|c|} \hline A & i_1 & i_3 \\ \hline \end{array} \right\} \end{aligned}$$

The time-complexity of the algorithms can be derived from the number of relevant indices, where each vector \vec{i} represents h indices. Step $D^{combine}$ in algorithm $\mathbb{P}(CYK)(hRMS(CFG, \rightrightarrows))$ has three relevant vectors $\vec{i}, \vec{j}, \vec{k}$, which each represent h indices, resulting in a time-complexity of $O(n^{3h})$. Step $D^{linearize}$ has a $O(n^{h+1})$ time-complexity. The correctness is inherited from the CFG cases. The following figure explains the variables of the deduction rule $D^{Combine}$ for $\mathbb{P}(CYK)(2RMS(CFG, \rightrightarrows))$:



References

- [1] Tilman Becker and Dominik Heckmann. 1998. Recursive Matrix Systems (RMS) and TAG. *Proceedings of the Fourth International Workshop on TAG (TAG+4)*, Philadelphia, PA.
- [2] Dominik Heckmann. 1999. Recursive Matrix Systems. *Diploma Thesis, University of Saarland*.
- [3] Giorgio Satta and William Schuler. 1998. Restrictions on Tree Adjoining Languages. *Proceedings of COLING-ACL'98*, pages 1176–1182.
- [4] Klaas Sikkels and Anton Nijhold. 1997. Parsing of Context-Free Languages. *Handbook of Formal Languages, Vol. 2, Springer, Berlin*.