

A Tabular Interpretation of Bottom-up Automata for TAG

Eric de la Clergerie
INRIA
Domaine de Voluceau
Rocquencourt, B.P. 105
78153 Le Chesnay Cedex
France
Eric.Clergerie@inria.fr

Miguel A. Alonso Pardo
Depto. de Computación
Universidad de La Coruña
Campus de Elviña s/n
15071 La Coruña
Spain
alonso@dc.fi.udc.es

David Cabrero Souto
Centro Ramón Piñeiro para a
Investigación en Humanidades
Estrada Santiago-Noia km 3
15896 Santiago de Compostela
Spain
dcabrero@cirp.es

Abstract

We present a tabular interpretation for a class of 2-Stack Automata that may be used to describe bottom-up parsing strategies for TAGs. The results are also useful for tabulating other existing bottom-up automata models for this kind of languages.

1 Introduction

Several extensions of push-down automata has been proposed as operational devices for describing parsing strategies for TAGs. Embedded Push-Down Automata [EPDA] (Vijay-Shanker, 1988) and 2-Stack Automata [2-SA] (Becker, 1994) are suitable operational devices for top-down strategies. For bottom-up strategies, Bottom-up EPDA [BEPDA] (Schabes and Vijay-Shanker, 1990; Rambow, 1994) and Linear Indexed Automata [LIA] (Nederhof, 1998) have been proposed.

We classify parsing strategies for TAGs w.r.t. the way adjoining is recognized and regardless of how elementary trees are traversed. In *Top-Down* strategies, the auxiliary tree to be adjoined is predicted once the adjoining node has been reached. Examples are the Earley-like parsing algorithms which preserve the correct prefix property (Nederhof, 1997). Conversely, in *Bottom-Up* strategies, adjoining is considered only when a candidate auxiliary tree has been completely traversed. Examples are the popular CYK-like (Vijay-Shanker and Joshi, 1985) and Earley-like parsing algorithms without the valid prefix property (Schabes, 1991).

A TAG parser must handle elementary tree traversing as well as adjoining processing and keep some information about these two kinds of task. Then, a 2-stack automata is adequate to implement parsing algorithms for TAG.

Polynomial time complexity can be lost for a non deterministic grammar if redundant computations are not discarded using some kind of dynamic programming (tabular) techniques. For the above mentioned automata models, systematic tabulation is only available for LIA.

The automata model proposed in this paper for bottom-up parsing strategies presents the following

characteristics: separation of the tree traversal and adjunction information by using two stacks; systematic tabulation, achieving $\mathcal{O}(n^6)$ time complexity and $\mathcal{O}(n^4)$ space complexity; and results comparable with existing tabular algorithms for TAGs.

2 (Strongly-driven) bottom-up 2-Stack Automata

Strongly Driven 2-Stack Automata [SD 2-SA] has been introduced in (de la Clergerie and Alonso Pardo, 1998) to describe arbitrary parsing strategies for TAGs. They work on 2 stacks with some restrictions added to make them equivalent, w.r.t. the recognized languages, to the class of tree adjoining languages.

A SD 2-SA uses the **Master Stack MS** to drive the evaluation and the **Auxiliary Stack AS** for restricted bookkeeping. Actually, **AS** should be considered as a stack of stacks, each of them representing a **session**. Typically, in TAG parsing, a session contains a sequence of adjunctions done along the spines of auxiliary trees. A session starts in mode **w** (write) where pop action are forbidden on **MS** and switches at some point to mode **e** (erase) where push actions are forbidden on **MS**. The actions on **AS** in mode **e** should faithfully retrace the actions done in mode **w**. Exiting a session is only possible when reaching back (in **e** mode) the **MS** element that initiated the session and when the session stack on **AS** is empty.

The bottom-up “projection” of SD 2-SA, henceforth BU 2-SA, imposes an additional restriction: **AS** must remain empty in mode **w**. That means that adjunction can be only recognized when a complete auxiliary tree has been constructed. The different behaviors of SD 2-SA and BU 2-SA are obvious when comparing the shape of derivations as illustrated in Fig. 1, where the axis display the stack sizes.

More formally, a BU 2-SA \mathcal{A} is specified by a 6-tuple $(\Sigma, \mathcal{M}, \mathcal{X}, \$_0, \$_f, \Theta)$ where Σ denotes the finite set of terminals, \mathcal{M} the finite set of master stack elements and \mathcal{X} the finite set of auxiliary stack elements. The **init** symbol $\$_0$ and **final** symbol $\$_f$

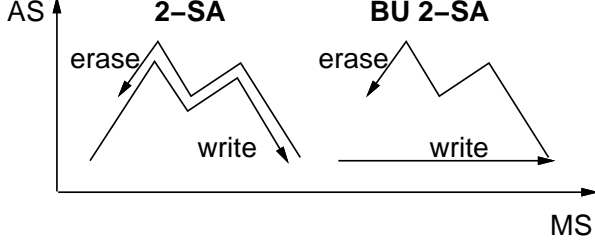


Figure 1: Derivation shapes for SD and BU 2-SA

are distinguished elements of \mathcal{M} . Θ is a finite set of transitions.

MS is a word in $(\mathcal{DM})^*$ where \mathcal{D} denotes the set $\{\triangleright, \models\}$ of **action marks**, projection of the larger action mark set $\{\nearrow, \rightarrow, \searrow, \models\}$ used for SD-2SA. Pushing an element on **MS** is either marked with \models if a “new session” starts at the same time, or by \triangleright otherwise.

AS is a word of $(\mathcal{KX}^*)^*$ where symbols in $\mathcal{K} = \{\models^{\mathbf{w}}, \models^{\mathbf{e}}\}$ are used to delimit session stacks and remember the mode of the previous session.

Given some input string $x_1 \dots x_n \in \Sigma^*$, a configuration of \mathcal{A} is a tuple (m, i, Ξ, ξ) where $m \in \{\mathbf{w}, \mathbf{e}\}$ denotes the current mode, i the current string position in $[0, n]$, Ξ the master stack and ξ the auxiliary stack. The initial configuration of \mathcal{A} is $(\mathbf{w}, 0, \models^{\$}_0, \models^{\mathbf{w}})$ and the final one $(\mathbf{e}, n, \models^{\$}_f, \models^{\mathbf{w}})$.

A transition τ is represented by a pair $(m, \Xi, \xi) \xrightarrow{z} (m', \Theta, \theta)$ where $m, m' \in \{\mathbf{w}, \mathbf{e}\}$, z in Σ^* , Ξ and Θ are suffixes of master stacks in $\mathcal{M}(\mathcal{DM})^*$, and ξ, θ suffixes of auxiliary stacks in $(\mathcal{X} \cup \mathcal{K})^*$. We denote $(m, i, \Psi \Xi, \psi \xi) \vdash (m', j, \Psi \Theta, \psi \theta)$ a valid derivation step using τ with $z = x_{i+1} \dots x_i$, and by \vdash^{\pm} the reflexive and transitive closure of \vdash . A string $a_1 \dots a_n$ is accepted by \mathcal{A} if $(\mathbf{w}, 0, \models^{\$}_0, \models^{\mathbf{w}}) \vdash^{\pm} (\mathbf{e}, n, \models^{\$}_f, \models^{\mathbf{w}})$.

For BU 2-SA, we consider the following kinds of transitions (which enforce that the **AS** topmost session remains empty in **w** mode), namely **SWAP** to change the top element of the **MS**; \models -**WRITE** and \models -**ERASE** to start and end sessions; and \triangleright -**WRITE** and δ -**ERASE** ($\delta \in \{\nearrow, \rightarrow, \searrow\}$) to push to and pop from **MS** while acting on **AS**:

$$\mathbf{SWAP1} \quad (p, A, \epsilon) \xrightarrow{z} (p, B, \epsilon)$$

$$\mathbf{SWAP2} \quad (\mathbf{w}, A, \models^o) \xrightarrow{z} (\mathbf{e}, B, \models^o)$$

$$\models\text{-WRITE} \quad (m, A, \epsilon) \xrightarrow{z} (\mathbf{w}, A \models B, \models^m)$$

$$\models\text{-ERASE} \quad (\mathbf{e}, A \models B, \models^m) \xrightarrow{z} (m, C, \epsilon)$$

$$\triangleright\text{-WRITE} \quad (\mathbf{w}, A, \epsilon) \xrightarrow{z} (\mathbf{w}, A \triangleright B, \epsilon)$$

$$\delta\text{-ERASE} \quad (\mathbf{e}, A \triangleright B, c) \xrightarrow{z} (\mathbf{e}, C, c') \text{ with}$$

$(\delta = \rightarrow \text{ and } c = c' = \epsilon)$ or $(\delta = \nearrow \text{ and } c = \epsilon)$
or $(\delta = \searrow \text{ and } c' = \epsilon)$.

3 TAG parsing with BU 2-SA

We present a BU 2-SA that simulates a Earley-like parsing algorithm without the valid-prefix property (Schabes, 1991). The automata performs full prediction on the context-free backbone but no prediction on the adjunctions during the descent phase.

Each elementary tree is represented by a set of context free productions of the form $\nu_{k,0} \rightarrow \nu_{k,1} \dots \nu_{k,n_k}$, where $\nu_{k,0}$ denotes some non-leaf node k and $\nu_{k,i}$ the i^{th} son of k , and a set of terminal productions $\nu_{k,0} \rightarrow a_k$, where $\nu_{k,0}$ denotes some leaf node k with terminal label a_k .

The 6-tuple $(V_T, \mathcal{M}, \mathcal{X}, \nu_{0,0}, \nu_{0,0}', \Theta)$ defines the automata \mathcal{A} , with $\mathcal{M} = \{\nabla_{k,i}^\alpha\} \cup \{\nu_{k,i}^\alpha\} \cup \{\nu_{k,i}'\}$ and $\mathcal{X} = \{\nabla_{k,i}^\alpha\}$, where symbols $\nabla_{k,i}$ denote dotted productions and $\nu_{k,i}^\alpha$ (resp. $\nu_{k,i}'$) denote the prediction (resp. successful recognition) of a node. The transitions are given by the following rules:

- Call / Return for a node not on a spine. The call starts a new session, exited at return.

$$\mathbf{CALL} : (m, \nabla_{k,i}, \epsilon) \mapsto (\mathbf{w}, \nabla_{k,i} \models \nu_{k,i+1}, \models^m)$$

$$\mathbf{RET} : (\mathbf{e}, \nabla_{k,i} \models \nu_{k,i+1}', \models^m) \mapsto (m, \nabla_{k,i+1}, \epsilon)$$

- Call / Return for an adjunction on node $\nu_{k,0}$. The computation is diverted to parse some acceptable auxiliary tree β with root node r_β . At return we check if the subtree attached to the foot node of β corresponds to the subtree rooted by $\nu_{k,0}$.

$$\mathbf{ACALL} : (\mathbf{w}, \nu_{k,0}, \epsilon) \mapsto (\mathbf{w}, \nu_{k,0} \triangleright r_\beta, \epsilon)$$

$$\mathbf{ARET} : (\mathbf{e}, \nu_{k,0} \triangleright r_\beta', \nabla_{k,n_k}) \mapsto (\mathbf{e}, \nu_{k,0}', \epsilon)$$

- Call / Return for a node $\nu_{k,i+1}$ on a spine. The adjunction stack is propagated bottom-up along the spine.

$$\mathbf{SCALL} : (\mathbf{w}, \nabla_{k,i}, \epsilon) \mapsto (\mathbf{w}, \nabla_{k,i} \triangleright \nu_{k,i+1}, \epsilon)$$

$$\mathbf{SRET} : (\mathbf{e}, \nabla_{k,i} \triangleright \nu_{k,i+1}', \epsilon) \mapsto (\mathbf{e}, \nabla_{k,i+1}, \epsilon)$$

- Call / Return for a foot node f_β . A candidate adjunction node for β is predicted. At return we remember what node was considered.

$$\mathbf{FCALL} : (\mathbf{w}, f_\beta, \epsilon) \mapsto (\mathbf{w}, f_\beta \triangleright \nabla_{k,0}, \epsilon)$$

$$\mathbf{FRET} : (\mathbf{e}, f_\beta \triangleright \nabla_{k,n_k}, \epsilon) \mapsto (\mathbf{e}, f_\beta', \nabla_{k,n_k})$$

- Production Selection

$$\mathbf{SEL} : (\mathbf{w}, \nu_{k,0}, \epsilon) \mapsto (\mathbf{w}, \nabla_{k,0}, \epsilon)$$

- Production Publishing

$$\mathbf{PUB} : (m, \nabla_{k,n_k}, \epsilon) \mapsto (\mathbf{e}, \nu_{k,0}', \epsilon)$$

- Scanning

$$\mathbf{SCAN} : (\mathbf{w}, \nu_{k,0}, \models^m) \xrightarrow{a_k} (\mathbf{e}, \nu_{k,0}', \models^m)$$

4 Tabulation

In a tabular framework, items store essential information about characteristics “points” of elementary derivations. Tabulation of SD 2-SA (de la Clergerie and Alonso Pardo, 1998), that achieves

$\mathcal{O}(n^6)$ time and $\mathcal{O}(n^5)$ space complexity, needs two kinds of items, namely 3-point Context-Free [CF] items and 5-point escaped Context-free [XCF] items. Each point is either a *mini configuration* $\langle i, A, a \rangle$ or a *micro configuration* $\langle i, A \rangle$ that stores some relevant information about a configuration, namely the position i in the input string, the top **MS** element A , and optionally the top **AS** element a . The uppermost curve of Fig. 2 illustrates a 3-point CF item $[\langle h, A, - \rangle, \langle i, B, - \rangle, \delta, \langle j, C, c \rangle]$, also denoted $\mathbf{B}\delta\bar{\mathbf{C}}\mathbf{w}$ where \mathbf{A} and \mathbf{B} are micro configurations and $\bar{\mathbf{C}}$ is a mini configuration. The uppermost curve of Fig. 3 illustrates a 5-point XCF item $[\langle h, A, - \rangle, \langle i, B, - \rangle, \delta, \langle p, D, d \rangle, \langle q, E, - \rangle, \langle j, C, c \rangle]$, also denoted $\mathbf{A}\mathbf{B}\delta[\bar{\mathbf{D}}\bar{\mathbf{E}}]\bar{\mathbf{C}}\mathbf{e}$ where \mathbf{A} , \mathbf{B} , \mathbf{E} (resp. $\bar{\mathbf{D}}$, $\bar{\mathbf{C}}$) are micro (resp. mini) configurations.

BU 2-SA restrictions imply that **AS** remains empty in **w** mode, so the points \mathbf{A} , \mathbf{B} and $\bar{\mathbf{C}}$ of a CF item and the points \mathbf{A} , \mathbf{B} and $\bar{\mathbf{D}}$ of a XCF item are “projected” w.r.t. the top element of the **AS**. Furthermore, it may be shown that point \mathbf{A} is actually redundant and can be discarded. The bottom curve of Fig. 2 illustrates a BU 2-SA CF item $[\langle i, B, - \rangle, \triangleright, \langle j, C, c \rangle]$, also denoted as $\mathbf{B}\triangleright\bar{\mathbf{C}}\mathbf{w}$. The bottom curve of Fig. 3 illustrates a BU 2-SA XCF item $[\langle i, B, - \rangle, \triangleright, \langle p, D, \models^o \rangle, \langle q, E, - \rangle, \langle j, C, c \rangle]$, also denoted as $\mathbf{B}\triangleright[\bar{\mathbf{D}}\bar{\mathbf{E}}]\bar{\mathbf{C}}\mathbf{e}$. In both figures, the projection is materialized by the dashed arrows.

Formally, we identify two kinds of items for BU 2-SA, associated to two different kinds of derivations:

Bottom-up CF [buCF]

items correspond to context-free derivations that depend only on the topmost element of **MS**

$$\begin{aligned} & (\mathbf{w}, i, \Xi B, \xi \models^o) \mid^* (m, j, \Xi B \triangleright C, \xi \models^o) \\ \text{or} & (o, i, \Xi B, \xi) \mid^* (\mathbf{w}, j, \Xi B \models C, \xi \models^o) \end{aligned}$$

and are denoted by $\mathbf{B}\delta\bar{\mathbf{C}}\mathbf{m}$, where $\mathbf{B} = \langle i, B \rangle$, $\bar{\mathbf{C}} = \langle j, C, \models^o \rangle$, and $\delta \in \mathbf{D}$.

Bottom-up Escaped CF [buXCF] items correspond to escaped context-free derivations of the form:

$$\begin{aligned} (\mathbf{w}, i, \Xi B, \xi \models^o) & \mid^* (\mathbf{w}, p, \Xi \Phi D, \xi \models^o) \\ & \mid^* (\mathbf{e}, q, \Xi \Phi D \triangleright E, \xi \models^o \phi) \\ & \mid^* (\mathbf{e}, j, \Xi B \triangleright C, \xi \models^o \phi c) \end{aligned}$$

and are denoted by $\mathbf{B}\triangleright[\bar{\mathbf{D}}\bar{\mathbf{E}}]\bar{\mathbf{C}}\mathbf{e}$, where $\mathbf{B} = \langle i, B \rangle$, $\mathbf{D} = \langle p, D \rangle$, $\mathbf{E} = \langle q, E \rangle$, $\bar{\mathbf{C}} = \langle j, C, c \rangle$.

A set of rules combines items and transitions in order to retrieve all possible derivations. Due to space limitations, we only describe the most complex rule (see Fig. 4), used to apply a transition $\tau = (\mathbf{e}, B \triangleright C, c) \xrightarrow{z} (\mathbf{e}, F, \epsilon)$, omitting the scanning constraint z on the input string:

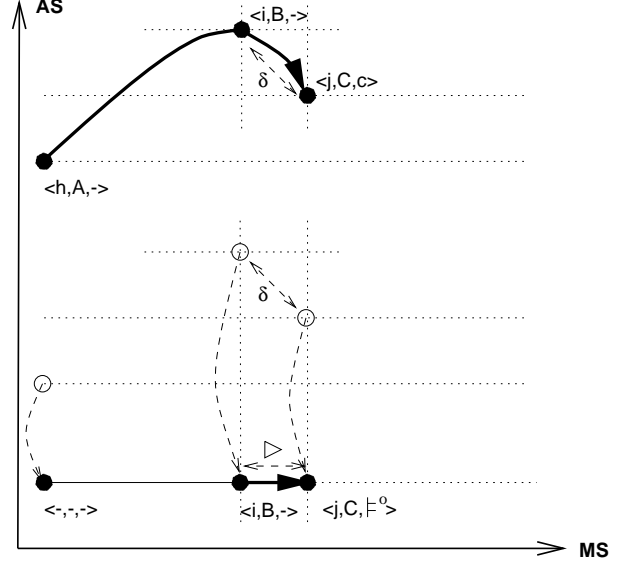


Figure 2: CF items for SD 2-SA and BU 2-SA

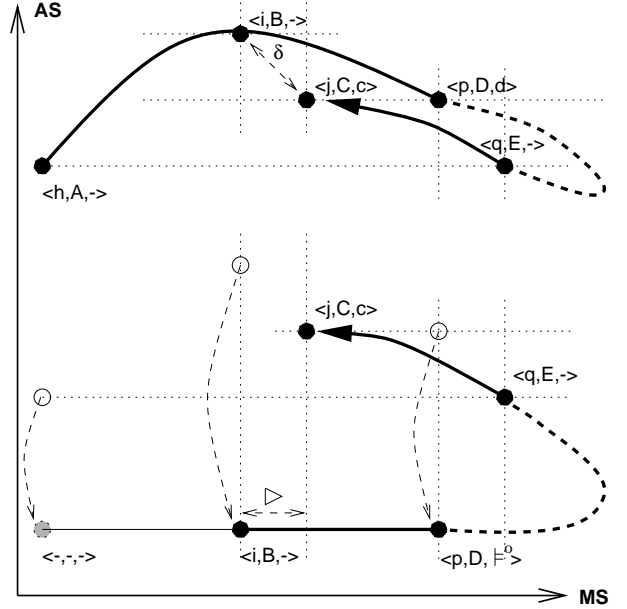


Figure 3: XCF items for SD 2-SA and BU 2-SA

$$\left. \begin{array}{l} \bar{\mathbf{B}}^o \triangleright [\bar{\mathbf{D}}\bar{\mathbf{E}}]^o \bar{\mathbf{C}}\mathbf{e} \\ N\delta\bar{\mathbf{B}}\mathbf{w} \\ \mathbf{D} \triangleright [\mathbf{O}\mathbf{P}]\bar{\mathbf{E}}\mathbf{e} \end{array} \right\} \xrightarrow{\tau} N\delta[\mathbf{O}\mathbf{P}]\bar{\mathbf{F}}\mathbf{e} \quad (1)$$

where $\bar{\mathbf{C}} = \langle j, C, c \rangle$, $\mathbf{B} = \langle i, B, \models^o \rangle$, $\bar{\mathbf{F}} = \langle k, F, b \rangle$, and $\bar{\mathbf{B}}^o = \langle i, B \rangle$ the projection of $\bar{\mathbf{B}}$ to a micro configuration.

The time complexity of this rule is $\mathcal{O}(n^7)$ but may be reduced to $\mathcal{O}(n^6)$ by partially applying the rule on the first two items to build an intermediary structure where \mathbf{B} is discarded.

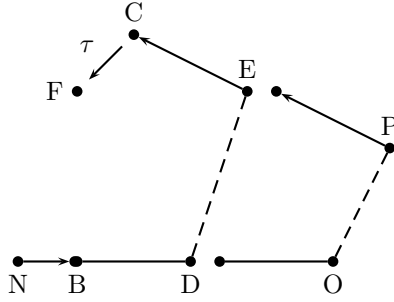


Figure 4: Application of Rule 1

Space complexity of the tabular technique for BU 2-SA is obviously $\mathcal{O}(n^4)$ as at most 4 indices are stored in buXCF items.

5 Related work

Our tabular interpretation may be used to re-interpret other existing tabular algorithms for TAGs, based on some automata model or not.

Linear Indexed Automata [LIA] (Nederhof, 1998) is the only other automata model we are aware of that has an associated tabular algorithm. This algorithm considers items $((B, C, i, j), (\diamond, \square, \square, 0, 0))$ corresponding to buCF items $B\delta C m$, as well as items $((B, C, i, j), (c, D, E, p, q))$ corresponding to buXCF items $B \triangleright [DE] \bar{C} e$. Because LIAs work on a stack of stacks, the empty stack markers we use are useless, the \models mark being implicit when the second part of an item is equal to $(\diamond, \square, \square, 0, 0)$.

If we now consider the tabular algorithm of (Vijay-Shanker and Weir, 1994), which is not based on an automata model, we find that, using their terminology, our buXCF items $B \triangleright [DE] \bar{C} e$ correspond to a head $B\bar{C}$ with a terminator pointer $[DE]$ and buCF items to a head, without terminator pointer.

In both cases, marks and modes (\mathbf{w} and \mathbf{e}) are absent from the proposed items, but one may show that they are actually implicitly present. They may be also be discarded from our items when considering specific parsing strategies, but are needed if one wishes to exploit the full potentiality of BU-2SA, for instance for more complex parsing strategies.

6 Conclusion

Bottom-up 2-SA may be seen as the projection of a subclass of strongly-driven 2-SA, specialized to describe parsing strategies for TAG where adjunction is recognized in a bottom-up way (i.e. when being in mode erase). A tabular interpretation of BU 2-SA is straightforwardly derived by “projecting” the tabular interpretation for SD 2-SA. So, a buXCF item $B \triangleright [DE] \bar{C} e$ is the projection of a XCF item $AB\delta[DE]\bar{C}e$ and a buCF item $B\delta C m$ is the projection of a CF item $AB\delta C m$. For SD 2-SA, \mathbf{A} is needed to handle popping on \mathbf{AS} in \mathbf{w} mode, but

it may be safely removed for BU 2-SA because of the extra condition on the emptiness of \mathbf{AS} in \mathbf{w} mode. While the worst case time complexity remains $\mathcal{O}(n^6)$, the worst case space complexity decreases from $\mathcal{O}(n^5)$ for 2-SA to $\mathcal{O}(n^4)$ for BU 2-SA. Of course, the drawback is the violation of the valid-prefix property and it remains to investigate whether or not this is a good thing for TAG grammars used in Natural Language Processing.

7 Acknowledgements

This work has been partially supported by the European Union (1FD97-0047-C04-02), Government of Spain (HF97-223) and Xunta de Galicia (XUGA10505B96 and XUGA20402B97).

References

- Tilman Becker. 1994. A new automaton model for TAGs: 2-SA. *Computational Intelligence*, 10(4):422–430.
- Eric de la Clergerie and Miguel A. Alonso Pardo. 1998. A tabular interpretation of a class of 2-Stack Automata. In *Proc. COLING/ACI'98*, Montreal, Canada, August.
- Mark-Jan Nederhof. 1997. Solving the correct-prefix property for TAGs. In T. Becker and H.-V. Krieger, editors, *Proc. of the Fifth Meeting on Mathematics of Language*, pages 124–130, Schloss Dagstuhl, Saarbruecken, Germany, August.
- Mark-Jan Nederhof. 1998. Linear indexed automata and tabulation of TAG parsing. In *Proc. of First Workshop on Tabulation in Parsing and Deduction (TAPD'98)*, pages 1–9, Paris, France, April.
- Owen Rambow. 1994. *Formal and Computational Aspects of Natural Language Syntax*. Ph.D. thesis, University of Pennsylvania.
- Yves Schabes and K. Vijay-Shanker. 1990. Deterministic left to right parsing of tree adjoining languages. In *Proc. of 28th Annual Meeting of the Association for Computational Linguistics*, pages 276–283, Oittsburgh, Pennsylvania, USA, June.
- Yves Schabes. 1991. The valid prefix property and left to right parsing of tree-adjoining grammar. In *Proc. of II International Workshop on Parsing Technologies, IWPT'91*, pages 21–30, Cancún, Mexico.
- K. Vijay-Shanker and Aravind K. Joshi. 1985. Some computational properties of tree adjoining grammars. In *23rd Annual Meeting of the Association for Computational Linguistics*, pages 82–93, Chicago, IL, USA, July.
- K. Vijay-Shanker and David J. Weir. 1994. Parsing some constrained grammar formalisms. *Computational Linguistics*, 19(4):591–636.
- K. Vijay-Shanker. 1988. *A Study of Tree Adjoining Grammars*. Ph.D. thesis, University of Pennsylvania, January.