

A formal definition of Bottom-up Embedded Push-Down Automata and their tabulation technique

Miguel A. Alonso, Eric de la Clergerie and Manuel Vilares

Abstract

We provide a formal definition of Bottom-up Embedded Push-Down Automata. By removing the finite-state control, we obtain an alternative definition which is adequate to define a tabulation framework for this class of automata and to show the equivalence with respect to other models of automata accepting tree adjoining languages.

Keywords: Automata, parsing, tabulation, tree adjoining grammars.

1 Introduction

Tree adjoining grammars (TAG) [7] are an extension of context-free grammars that use trees instead of productions as the primary representing structure. Formally, a TAG is a 5-tuple $\mathcal{G} = (V_N, V_T, S, \mathbf{I}, \mathbf{A})$, where V_N is a finite set of non-terminal symbols, V_T a finite set of terminal symbols, $S \in V_N$ is the axiom of the grammar, \mathbf{I} a finite set of *initial trees* and \mathbf{A} a finite set of *auxiliary trees*. $\mathbf{I} \cup \mathbf{A}$ is the set of *elementary trees*. Internal nodes are labeled by non-terminals and leaf nodes by terminals or ε , except for just one leaf per auxiliary tree (the *foot*) which is labeled by the same non-terminal used as the label of its root node. The path in an elementary tree from the root node to the foot node is called the *spine* of the tree.

New trees are derived by *adjoining*: let α be a tree containing a node N^α labeled by A and let β be an auxiliary tree whose root and foot nodes are also labeled by A . Then, the adjoining of β at the *adjunction node* N^α is obtained by excising the

Miguel A. Alonso Pardo, Departamento de Computación, Facultad de Informática, Universidad de La Coruña, Campus de Elviña s/n, 15071 La Coruña (Spain), <http://www.dc.fi.udc.es/~alonso/>, alonso@dc.fi.udc.es

Eric Villemonete de la Clergerie, Institut National de Recherche en Informatique et en Automatique, Domaine de Voluceau, Rocquencourt, B.P. 105, 78153 Le Chesnay Cedex (France), <http://atoll.inria.fr/~clerger/>, Eric.De.La.Clergerie@inria.fr

Manuel Vilares Ferro, Departamento de Computación, Facultad de Informática, Universidad de La Coruña, Campus de Elviña s/n, 15071 La Coruña (Spain), <http://www.dc.fi.udc.es/~vilares/>, vilares@dc.fi.udc.es

This research has been partially supported by the FEDER of EU (Grant 1FD97-0047-C04-02) and Xunta de Galicia (Grant PGIDT99XI10502B).

subtree of α with root N^α , attaching β to N^α and attaching the excised subtree to the foot of β .

The operation of *substitution* does not increase the generative power of the formalism but it is usually considered when we are dealing with lexicalized tree adjoining grammars. In this case, non-terminals can also label leaf nodes (called substitution nodes) of elementary trees. An initial tree can be substituted at a substitution node if its root is labeled by the same non-terminal that labels the substitution node.

Bottom-up embedded push-down automata (BEPDA) have been described in [11] as an extension of push-down automata adequate to implement parsing strategies for TAG in which adjunctions are recognized in a bottom-up way. In fact, BEPDA are the dual version of embedded push-down automata (EPDA) [12, 2], a model of automata in which adjunctions must be recognized in a top-down way. A less informal presentation of BEPDA has been shown in [10], with some inconsistencies between the set of allowed transitions and the set of configurations attainable.

Right-oriented linear indexed automata [9] and bottom-up 2-stack automata [5] can be used to implement parsing strategies similar to those of BEPDA. Both models of automata have associated a tabulation framework allowing their execution in polynomial time with respect to the size of the input string. In this paper we provide a formal definition for BEPDA and then we propose an alternative definition in order to provide a tabulation framework for this class of automata.

2 A formal definition of BEPDA

A bottom-up embedded push-down automata consists of a finite-state control, an input tape and a stack (that we call the main stack) made up of non-empty stacks containing stack symbols. Formally, a BEPDA is defined by a tuple $(Q, V_T, V_S, \delta, q_0, Q_F, \$_f)$ where:

- Q is a finite set of states.
- V_T is a finite set of terminal symbols.
- V_S is a finite set of stack symbols.
- $q_0 \in Q$ is the initial state.
- $Q_F \subseteq Q$ is the set of final states.
- $$_f \in V_S$ is the final stack symbol.$
- δ is a mapping from $Q \times V_T \cup \{\epsilon\} \times ([V_S^+]^* \times V_S^* \times ([V_S^+]^*$ into finite subsets of $Q \times V_S \cup \{[V_S\}$, where $[\notin V_S$ is a new symbol used as stack separator.

An *instantaneous configuration* is a triple (q, Υ, w) , where q is the current state, $\Upsilon \in ([V_S^+]^*$ represents the contents of the automaton stack and w is the part of the input string that is yet to be read. It is important to remark that every individual stack contained in the main stack must store at least one stack symbol. The main stack will be empty only in the initial configuration (q_0, ϵ, w) .

Transitions in δ allow an automaton to derive a new configuration (q', Υ', w) from a configuration (q, Υ, aw) , which is denoted $(q, \Upsilon, aw) \vdash (q', \Upsilon', w)$. The reflexive and transitive closure of \vdash is denoted by \vdash^* . There exist two different types of transitions:

1. Transitions of the first type are of the form

$$(q', [Z] \in \delta(q, a, \epsilon, \epsilon, \epsilon)$$

and they can be applied to a configuration (q, Υ, aw) to yield a configuration $(q', \Upsilon[Z, w])$.

2. Transitions of the second type are of the form

$$(q', Z) \in \delta(q, a, [\alpha_k \dots [\alpha_{i+1}, Z_m \dots Z_1, [\alpha_i \dots [\alpha_1]$$

where $m \geq 0$ and $k \geq i \geq 0$. They can be applied to a configuration

$$(q, \Upsilon[\alpha_k \dots [\alpha_{i+1}[\alpha Z_m \dots Z_1[\alpha_i \dots [\alpha_1, aw]$$

to yield a configuration

$$(q', \Upsilon[\alpha Z, w)$$

The *language accepted by final state* by a BEPDA is the set $w \in V_T^*$ of input strings such that $(q_0, \epsilon, w) \vdash^*(p, \Upsilon, \epsilon)$, where $p \in Q_F$ and $\Upsilon \in ([V_S^+])^*$.

The *language accepted by empty stack* by a BEPDA is the set $w \in V_T^*$ of input strings such that $(q_0, \epsilon, w) \vdash^*(q, [\$_f, \epsilon)$ for any $q \in Q$. At this point it is interesting to observe the duality with respect to EPDA: EPDA computations start with a stack $[\$_0$ to finish with an empty stack while BEPDA computations start with an empty stack to finish with a stack $[\$_f$.

It can be proved that for any BEPDA accepting a language L by final state there exists a BEPDA accepting the same language by empty stack and vice versa.

Example 1 The bottom-up embedded push-down automaton defined by the tuple $(\{q_0, q_1, q_2, q_3\}, \{a, b, c, d\}, \{B, C, D\}, \delta, q_0, \emptyset, \$_f)$, with δ containing the transitions shown in Fig. 1 (left box), accepts the language $\{a^n b^n c^n d^n \mid n \geq 0\}$ by empty stack. The sequence of configurations for the recognition of the input string $aabbccdd$ is shown in Fig. 1 (right box), where the first column shows the transition applied, the second one the current state, the third one the contents of the stack and the fourth column shows the part of the input string to be read. \square

3 BEPDA without states

Finite-state control is not a fundamental component of push-down automata, as the current state in a configuration can be stored in the top element of the stack of the automaton [8]. Finite-state control can also be eliminated from bottom-up embedded push-down automata, obtaining a new definition that considers a BEPDA as a tuple $(V_T, V_S, \Theta, \$_0, \$_f)$ where V_T is a finite set of terminal symbols, V_S is a finite set of stack symbols, $\$_0 \in V_S$ is the initial stack symbol, $\$_f \in V_S$ is the final stack symbol and Θ is a finite set of six types of transition:

(a) $(q_0, [D] \in \delta(q_0, a, \epsilon, \epsilon, \epsilon)$	q_0 $aabbccdd$
(b) $(q_1, [C] \in \delta(q_0, b, \epsilon, \epsilon, \epsilon)$	(a) q_0 $[D]$ $abbccdd$
(c) $(q_1, [C] \in \delta(q_1, b, \epsilon, \epsilon, \epsilon)$	(a) q_0 $[D[D$ $bbccdd$
(d) $(q_2, B) \in \delta(q_1, c, \epsilon, C, \epsilon)$	(b) q_1 $[D[D[C$ $bccdd$
(e) $(q_2, B) \in \delta(q_2, c, [C, \epsilon, \epsilon)$	(c) q_1 $[D[D[C[C$ $ccdd$
(f) $(q_3, B) \in \delta(q_2, d, [D, BB, \epsilon)$	(d) q_2 $[D[D[C[B$ cdd
(g) $(q_3, B) \in \delta(q_3, d, [D, BB, \epsilon)$	(e) q_2 $[D[D[BB$ dd
(h) $(q_3, \$_f) \in \delta(q_3, d, [D, B, \epsilon)$	(f) q_3 $[D[B$ d
(i) $(q_0, [\$_f) \in \delta(q_0, a, \epsilon, \epsilon, \epsilon)$	(h) q_3 $[\$_f$

Figure 1: BEPDA accepting $\{a^n b^n c^n d^n \mid n > 0\}$ and configurations for $aabbccdd$

SWAP: Transitions of the form $C \xrightarrow{a} F$ that replace the top element of the top stack while scanning a . The application of such a transition on a stack $\Upsilon[\alpha C$ returns the stack $\Upsilon[\alpha F$.

PUSH: Transitions of the form $C \xrightarrow{a} C F$ that push F onto C . The application of such a transition on a stack $\Upsilon[\alpha C$ returns the stack $\Upsilon[\alpha C F$.

POP: Transitions of the form $C F \xrightarrow{a} G$ that replace C and F by G . The application of such a transition on $\Upsilon[\alpha C F$ returns the stack $\Upsilon[\alpha G$.

WRAP: Transitions of the form $C \xrightarrow{a} C, [F$ that push a new stack $[F$ on the top of the main stack. The application of such a transition on a stack $\Upsilon[\alpha C$ returns the stack $\Upsilon[\alpha C [F$.

UNWRAP-A: Transitions *unwrap-above* of the form $C, [F \xrightarrow{a} G$ that delete the top stack $[F$ and replace the new top element by G . The application of such a transition on a stack $\Upsilon[\alpha C [F$ returns the stack $\Upsilon[\alpha G$.

UNWRAP-B: Transitions *unwrap-below* of the form $[C, F \xrightarrow{a} G$ that delete the stack $[C$ placed just below the top stack and replace the top element by G . The application of such a transition on a stack $\Upsilon[C [F$ returns the stack $\Upsilon[\alpha G$.

(a)	$\$0 \xrightarrow{a} \$0, [D$	$[\$0$	$aabbccdd$
(b)	$D \xrightarrow{a} D, [D$	(a) $[\$0[D$	$aabbccdd$
(c)	$D \mapsto D, [C$	(b) $[\$0[D[D$	$abbccdd$
(d)	$C \xrightarrow{b} C, [C$	(c) $[\$0[D[D[C$	$bbccdd$
(e)	$C \mapsto B$	(d) $[\$0[D[D[C[C$	$bccdd$
(f)	$B \mapsto BE$	(d) $[\$0[D[D[C[C[C$	$ccdd$
(g)	$[C, E \xrightarrow{c} C$	(e) $[\$0[D[D[C[C[B$	$ccdd$
(h)	$BC \mapsto B$	(f) $[\$0[D[D[C[C[BE$	$ccdd$
(i)	$[D, B \xrightarrow{d} D$	(g) $[\$0[D[D[C[BC$	cdd
(j)	$BD \mapsto B$	(e) $[\$0[D[D[C[BB$	cdd
(k)	$\$0, [D \mapsto \f	(f) $[\$0[D[D[C[BBE$	cdd
		(g) $[\$0[D[D[BBC$	dd
		(h) $[\$0[D[D[BB$	dd
		(i) $[\$0[D[BD$	d
		(j) $[\$0[D[B$	d
		(i) $[\$0[D$	
		(k) $[\$f$	

Figure 2: BEPDA without finite-state control for $\{a^n b^n c^n d^n \mid n > 0\}$

where $C, F, G \in V_S$, $\Upsilon \in ([V_S^*])^*$, $\alpha \in V_S^*$ and $a \in V_T \cup \{\epsilon\}$. It can be proved that transitions of a BEPDA with states can be emulated by transitions in Θ and vice versa.

An *instantaneous configuration* is a pair (Υ, w) , where Υ represents the contents of the automaton stack and w is the part of the input string that is yet to be read. A configuration (Υ, aw) derives a configuration (Υ', w) , denoted $(\Upsilon, aw) \vdash (\Upsilon', w)$, if and only if there exist a transition that applied to Υ gives Υ' and scans a from the input string. We use \vdash^* to denote the reflexive and transitive closure of \vdash . An input string is accepted by an BEPDA if $([\$0, w) \vdash^* ([\$f, \epsilon)$. The language accepted by an BEPDA is the set of $w \in V_T^*$ such that $([\$0, w) \vdash^* ([\$f, \epsilon)$.

Example 2 The bottom-up embedded push-down automaton without states defined by the tuple $(\{a, b, c, d\}, \{B, C, D, E, F, \$0, \$f\}, \Theta, \$0, \$f)$, with Θ containing the transitions shown in Fig. 2 (left box), accepts the language $\{a^n b^n c^n d^n \mid n \geq 0\}$. The sequence of configurations to recognize the input string $aabbccdd$ is also shown in Fig. 2 (right-box), where the first column shows the transition applied in each step, the second one shows the contents of the stack and the third column shows the part of the input string to be read. \square

[INIT]	$\$0 \mapsto \$0 \ [\overrightarrow{\top}^\alpha]$	$\alpha \in \mathbf{I}, S = \text{label}(\mathbf{R}^\alpha)$
[FINAL]	$\$0 \ [\overleftarrow{\top}^\alpha \mapsto \f	$\alpha \in \mathbf{I}, S = \text{label}(\mathbf{R}^\alpha)$
[CALL]	$\nabla_{r,s}^\gamma \mapsto \nabla_{r,s}^\gamma \ [\overrightarrow{N_{r,s+1}^\gamma}]$	$N_{r,s+1}^\gamma \notin \text{spine}(\gamma), \mathbf{nil} \in \text{adj}(N_{r,s+1}^\gamma)$
[SCALL]	$\nabla_{r,s}^\beta \mapsto \nabla_{r,s}^\beta \ [\overrightarrow{N_{r,s+1}^\beta}]$	$N_{r,s+1}^\beta \in \text{spine}(\beta), \mathbf{nil} \in \text{adj}(N_{r,s+1}^\beta)$
[SEL]	$\overrightarrow{N_{r,0}^\gamma} \mapsto \nabla_{r,0}^\gamma$	
[PUB]	$\nabla_{r,n_r}^\gamma \mapsto \overleftarrow{N_{r,0}^\gamma}$	
[RET]	$\nabla_{r,s}^\gamma \ [\overleftarrow{N_{r,s+1}^\gamma} \mapsto \nabla_{r,s+1}^\gamma$	$N_{r,s+1}^\gamma \notin \text{spine}(\gamma), \mathbf{nil} \in \text{adj}(N_{r,s+1}^\gamma)$
[SRET]	$[\nabla_{r,s}^\beta \ [\overleftarrow{N_{r,s+1}^\beta} \mapsto \nabla_{r,s+1}^\beta$	$N_{r,s+1}^\beta \in \text{spine}(\beta), \mathbf{nil} \in \text{adj}(N_{r,s+1}^\beta)$
[SCAN]	$\overrightarrow{N_{r,0}^\gamma} \xrightarrow{a} \overleftarrow{N_{r,0}^\gamma}$	$N_{r,0}^\gamma \rightarrow a$
[AdjCALL]	$\nabla_{r,s}^\gamma \mapsto \nabla_{r,s}^\gamma \ [\overrightarrow{\top}^\beta]$	$\beta \in \text{adj}(N_{r,s+1}^\gamma)$
[AdjRET-a]	$[\nabla_{r,s}^\gamma \ [\overleftarrow{\top}^\beta \mapsto \top$	$\beta \in \text{adj}(N_{r,s+1}^\gamma)$
[AdjRET-b]	$\Delta_{r,s}^\gamma \ \top \mapsto \nabla_{r,s+1}^\gamma$	
[FootCALL]	$\nabla_{f,0}^\beta \mapsto \nabla_{f,0}^\beta \ [\overrightarrow{N_{r,s+1}^\gamma}]$	$N_{f,0}^\beta = \mathbf{F}^\beta, \beta \in \text{adj}(N_{r,s+1}^\gamma)$
[FootRET-a]	$[\nabla_{f,0}^\beta \ [\overleftarrow{N_{r,s+1}^\gamma} \mapsto \Delta_{r,s}^\gamma$	$N_{f,0}^\beta = \mathbf{F}^\beta, \beta \in \text{adj}(N_{r,s+1}^\gamma)$
[FootRET-b]	$\Delta_{r,s}^\gamma \mapsto \Delta_{r,s}^\gamma \ \nabla_{f,1}^\beta$	$N_{f,0}^\beta = \mathbf{F}^\beta, \beta \in \text{adj}(N_{r,s+1}^\gamma)$
[SubsCALL]	$\nabla_{r,s}^\gamma \mapsto \nabla_{r,s}^\gamma \ [\overrightarrow{\top}^\alpha]$	$\alpha \in \text{subs}(N_{r,s+1}^\gamma)$
[SubsRET]	$\nabla_{r,s}^\gamma \ [\overleftarrow{\top}^\alpha \mapsto \nabla_{r,s+1}^\gamma$	$\alpha \in \text{subs}(N_{r,s+1}^\gamma)$

Figure 3: Generic compilation schema for TAG

4 Compiling TAG into BEPDA

Automata are interesting for parsing because they allow us to separate two different problems that arise during the definition of parsing algorithms: the description of the parsing strategy and the execution of the parser. By means of automata, a parsing strategy for a given grammar can be translated into a set of transitions defining a (possibly non deterministic) automaton and then the automaton can be executed using some standard technique.

In this section we define a generic compilation schema for tree adjoining grammars

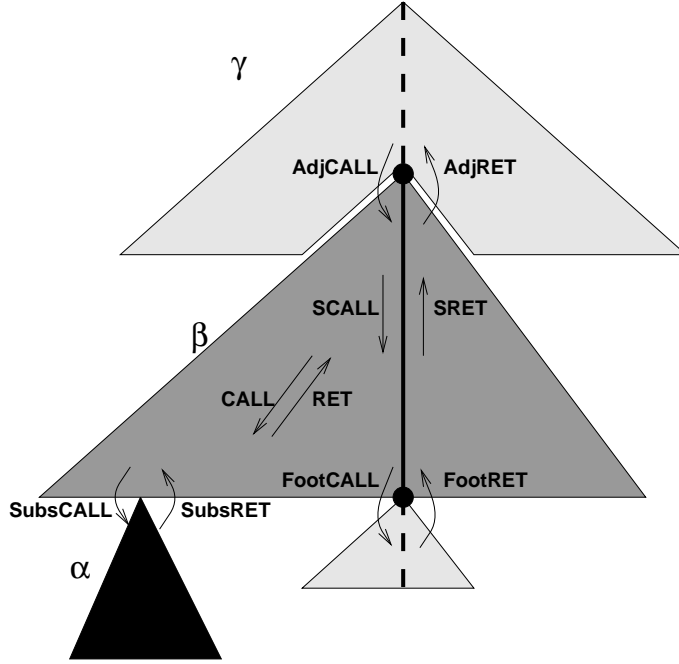


Figure 4: Meaning of compilation rules

based on a call/return model [6]. We consider each elementary tree γ of a TAG as formed by a set of context-free productions $\mathcal{P}(\gamma)$: a node N^γ and its g children $N_1^\gamma \dots N_g^\gamma$ are represented by a production $N^\gamma \rightarrow N_1^\gamma \dots N_g^\gamma$. The elements of the productions are the nodes of the tree, except for the case of elements belonging to $V_T \cup \{\varepsilon\}$ in the right-hand side of production. Those elements may not have children and are not candidates to be adjunction nodes, so we identify such nodes labeled by a terminal with that terminal. We use $\beta \in \text{adj}(N^\gamma)$ to denote that a tree $\beta \in \mathbf{A}$ may be adjoined at node N^γ . If adjunction is not mandatory at N^γ , then $\text{nil} \in \text{adj}(N^\gamma)$. If a tree $\alpha \in \mathbf{I}$ may be substituted at node N^γ , then $\alpha \in \text{subs}(N^\gamma)$. We consider the additional productions $\top^\alpha \rightarrow \mathbf{R}^\alpha$, $\top^\beta \rightarrow \mathbf{R}^\beta$ and $\mathbf{F}^\beta \rightarrow \perp$ for each initial tree α and each auxiliary tree β , where \mathbf{R}^α is the root node of α and \mathbf{R}^β and \mathbf{F}^β are the root node and foot node of β , respectively.

Fig. 3 shows the compilation rules from TAG to BEPDA, where symbols $\nabla_{r,s}^\gamma$ has been introduced to denote dotted productions. The meaning of each compilation rule is graphically shown in Fig. 4. This schema is parameterized by $\overrightarrow{N}^\gamma$, the information propagated top-down with respect to the node N^γ , and by \overleftarrow{N}^γ , the information propagated bottom-up. When the schema is used to implement a top-down traversal of elementary trees $\overrightarrow{N}^\gamma = N^\gamma$ and $\overleftarrow{N}^\gamma = \square$, where \square is a fresh stack symbol. A bottom-up traversal requires $\overrightarrow{N}^\gamma = \square$ and $\overleftarrow{N}^\gamma = N^\gamma$. For a mixed traversal of elementary trees, $\overrightarrow{N}^\gamma = \overline{N}^\gamma$ and $\overleftarrow{N}^\gamma = \overleftarrow{\overline{N}^\gamma}$, where \overline{N}^γ and $\overleftarrow{\overline{N}^\gamma}$ are used to distinguish the top-down prediction from the bottom-up propagation of a node.

With respect to adjunctions, we can observe in Fig. 3 that each stack stores pending adjunctions with respect to the node placed on the top of the stack in a bottom-up treatment of adjunctions: when a foot node is reached, the adjunction node is stored on the top of the stack (**[FootCALL-a]**); the traversal of the elemen-

Transition	BEPDA	R-LIA
SWAP	$C \xrightarrow{a} F$	$C[\circ\circ] \xrightarrow{a} F[\circ\circ]$
PUSH	$C \xrightarrow{a} CF$	$C[\circ\circ] \xrightarrow{a} F[\circ\circ C]$
POP	$CF \xrightarrow{a} G$	$F[\circ\circ C] \xrightarrow{a} G[\circ\circ]$
UNWRAP-A	$C, [F \xrightarrow{a} G$	$C[\circ\circ] F[] \xrightarrow{a} G[\circ\circ]$
UNWRAP-B	$[C, F \xrightarrow{a} G$	$C[] F[\circ\circ] \xrightarrow{a} G[\circ\circ]$
WRAP	$C \xrightarrow{a} C, [F$	$C[\circ\circ] \xrightarrow{a} C[\circ\circ] F[]$

Figure 5: Equivalence between BEPDA and L-LIA

tary tree is suspended to continue with the traversal of the adjoined auxiliary tree ([**FootCALL-b**]); the adjunction stack is propagated through the spine ([**SRET**]) up to the root node ([**AdjRET-a**]); and then the stack element corresponding to the auxiliary tree is eliminated to resume the traversal of the elementary tree ([**AdjRET-b**]). To avoid confusion, we store $\Delta_{r,s}^\gamma$ instead of $\nabla_{r,s}^\gamma$ to indicate that an adjunction was started at node $N_{r,s+1}^\gamma$. A symbol Δ can be seen as a symbol ∇ waiting an adjunction to be completed.

Substitution is managed through transitions generated by compilation rules [**SubsCALL**], which start the traversal of the substituted trees, and [**SubsRET**], which resume the traversal of the tree containing the substitution node once the substituted tree has been completely traversed.

5 BEPDA and other automata for TAG

5.1 Right-oriented linear indexed automata

Linear indexed automata (LIA) [3] are an extension of push-down automata in which each stack symbol has been associated to a list of indices. Right-oriented linear indexed automata (R-LIA) are a subclass of linear indexed automata that can be used to implement parsing strategies for TAG in which adjunctions are recognized in a bottom-up way. BEPDA and R-LIA are equivalent classes of automata. Given a BEPDA, the equivalent R-LIA is obtained by means of a simple change in the notation: the top element of a stack is considered a stack symbol, and the rest of the stack is considered the indices list associated to it, as is shown in Fig. 5. The same procedure also serves to obtain the BEPDA equivalent to a given R-LIA.

5.2 Bottom-up 2-stack automata

Strongly-driven 2-stack automata (SD-2SA) [4] are an extension of push-down automata working on a pair of asymmetric stacks, a master stack and an auxiliary

	BEPDA transition	BU-2SA transition	
SWAP	$C \xrightarrow{a} F$	$(m, C, \epsilon) \xrightarrow{a} (m, F, \epsilon)$	SWAP1
	$[C \xrightarrow{a} [F$	$(\mathbf{w}, C, \models^m) \xrightarrow{a} (\mathbf{e}, F, \models^m)$	SWAP2
WRAP	$C \xrightarrow{a} C, [F$	$(m, C, \epsilon) \xrightarrow{a} (\mathbf{w}, C \models^m F, \models^m)$	\models WRITE
	$[C \xrightarrow{a} [C, [F$	$(\mathbf{w}, C, \epsilon) \xrightarrow{a} (\mathbf{w}, C \triangleright F, \epsilon)$	\triangleright WRITE
UNWRAP-A	$C, [F \xrightarrow{a} G$	$(\mathbf{e}, C \models^m F, \models^m) \xrightarrow{a} (m, G, \epsilon)$	\models ERASE
UNWRAP-B	$[C, F \xrightarrow{a} G$	$(\mathbf{e}, C \triangleright F, \epsilon) \xrightarrow{a} (\mathbf{e}, G, \epsilon)$	\rightarrow ERASE
	$[C, XF \xrightarrow{a} G$	$(\mathbf{e}, C \triangleright F, X) \xrightarrow{a} (\mathbf{e}, G, \epsilon)$	\nearrow ERASE
	$[C, F \xrightarrow{a} XG$	$(\mathbf{e}, C \triangleright F, \epsilon) \xrightarrow{a} (\mathbf{e}, G, X)$	\searrow ERASE

Figure 6: Correspondence between BEPDA and BU-2SA

stack. These stacks are partitioned into sessions. Computations in each session are performed in one of two modes *write* and *erase*. A session starts in mode *write* and switches at some point to mode *erase*. In mode *write* (resp. *erase*), no element can be popped from (resp. pushed to) the master stack. Switching back from *erase* to *write* mode is not allowed. Bottom-up 2-stack automata (BU-2SA) [5] are a projection of SD-2SA requiring the emptiness of the auxiliary stack during computations in mode *write*. When a new session is created in a BU-2SA, a mark \models is left on the master stack, other movements performed in *write* mode leaving a mark \triangleright . These marks are popped in *erase* mode.

The full set of BU-2SA transitions is shown in Fig. 6. Transitions of type **SWAP2** are equivalent to $[C \xrightarrow{a} [F$ in BEPDA, compound transitions obtained from the consecutive application of $C \mapsto C, [F'$ and $[C, F' \xrightarrow{a} F$, where F' is a fresh stack symbol. In a similar way, transitions of type \nearrow **ERASE** are translated into compound transitions formed by an UNWRAP-B and a POP transition, and transitions of type \searrow **ERASE** are translated into the composition of UNWRAP-B and PUSH transitions. Slightly different is the case for transitions of type \triangleright **WRITE**, equivalent to $[C \xrightarrow{a} [C, [F$ transitions in BEPDA, which are obtained as the consecutive application of $[C \xrightarrow{a} [C'$ and $C' \mapsto C', [F$, an additional transition $C', [G \xrightarrow{b} K$ for each transition $C, [G \xrightarrow{b} K$ in the automaton, and an additional transition $[C', G \xrightarrow{b} K$ for each transition $[C, G \xrightarrow{b} K$, where C' is a fresh stack symbol.

As a consequence, it is possible to build a BEPDA for any given BU-2SA. However, the reverse is not always true: PUSH and POP transitions can only be translated into BU-2SA if they are merged with an UNWRAP-B transition. So, a BEPDA implementing a shift-reduce strategy (requiring the use of PUSH and POP transitions in combination with UNWRAP-A transitions¹) can not be translated into a BU-2SA.

¹A linear indexed automata implementing a LR-like strategy for linear indexed grammars using this kind of transitions is described in [1].

6 Tabulation

The direct execution of BEPDA may be exponential with respect to the length of the input string and may even loop. To get polynomial complexity, we must avoid duplicating computations by tabulating traces of configurations called *items*. The amount of information to keep in an item is the crucial point to determine to get efficient executions.

In order to define items and attending to the form of the transitions, we can classify derivations of BEPDA into the following two types:

Call derivations. Correspond to the placement of an unitary stack onto the top of the main stack:

$$(\Upsilon [B, a_{i+1} \dots a_n] \vdash^* (\Upsilon [C, a_{j+1} \dots a_n])$$

where $B, C \in V_S$ and $\Upsilon \in ([V_S^*])^*$. These derivations are independent of Υ and so they can be represented by items

$$[B, i, C, j, - \mid -, -, -, -]$$

Return derivations. Correspond to the placement of a non-unitary stack onto the top of the main stack:

$$\begin{aligned} & (\Upsilon [B, a_{i+1} \dots a_n]) \\ & \vdash^* (\Upsilon [B \ \Upsilon_1 [D, a_{p+1} \dots a_n]) \\ & \vdash^* (\Upsilon [B \ \Upsilon_1 [\alpha E, a_{q+1} \dots a_n]) \\ & \vdash^* (\Upsilon [\alpha X C, a_{j+1} \dots a_n]) \end{aligned}$$

where $B, C, D, E, X \in V_S$, $\alpha \in V_S^*$, $\Upsilon, \Upsilon_1 \in ([V_S^*])^*$. The two occurrences of α denote the same stack in the sense that α is neither consulted nor modified through derivation. These derivations are independent of Υ but not with respect to the subderivation

$$([D, a_{p+1} \dots a_n] \vdash^* ([\alpha E, a_{q+1} \dots a_n])$$

so they can be represented in compact form by items

$$[B, i, C, j, X \mid D, p, E, q]$$

To combine items, we use the following set of inference rules. Each rule is of the form $\frac{\eta_1, \dots, \eta_k}{\eta'} \text{trans}$, meaning that if all antecedents η_i are present and there exists the transition *trans*, then the consequent item η' should be generated.

$$\frac{[B, i, C, j, X \mid D, p, E, q]}{[B, i, F, k, X \mid D, p, E, q]} C \xrightarrow{a} F$$

$$\frac{[B, i, C, j, X \mid D, p, E, q]}{[B, i, F, k, C \mid B, i, C, j]} C \xrightarrow{a} CF$$

$$\frac{[B, i, F, j, C \mid D, p, E, q]}{[B, i, G, k, X' \mid O, u, P, v]} \frac{[D, p, E, q, X' \mid O, u, P, v]}{[B, i, G, k, X' \mid O, u, P, v]} CF \xrightarrow{a} G$$

$$\frac{[B, i, C, j, X \mid D, p, E, q]}{[F, k, F, k, - \mid -, -, -, -]} C \xrightarrow{a} C, [F$$

$$\frac{[F, k, F', k', - \mid -, -, -, -]}{[B, i, C, j, X \mid D, p, E, q]} C \xrightarrow{a} C, [F$$

$$\frac{[B, i, G, l, X \mid D, p, E, q]}{[B, i, G, l, X \mid D, p, E, q]} C, [F' \xrightarrow{b} G$$

$$\frac{[F, k, F', k', X \mid D, p, E, q]}{[B, i, C, j, - \mid -, -, -, -]} C \xrightarrow{a} C, [F$$

$$\frac{[B, i, G, l, X \mid D, p, E, q]}{[B, i, G, l, X \mid D, p, E, q]} [C, F' \xrightarrow{b} G$$

where $k = j$ if $a = \epsilon$, $k = j + 1$ if $a = a_{j+1}$, $l = k'$ if $b = \epsilon$ and $l = k' + 1$ if $b = a_{k'+1}$.

Computation start with the initial item $[\$0, 0, \$0, 0, - \mid -, -, -, -]$. An input string $a_1 \dots a_n$ has been recognized if the final item $[\$0, 0, \$f, n, - \mid -, -, -, -]$ is present. It can be proved that handling items with the inference rules is equivalent to applying the transitions on the whole stacks.

The space complexity of the proposed tabulation technique with respect to the length of the input string is $\mathcal{O}(n^4)$, due to every item stores four positions of the input string. The worst case time complexity is $\mathcal{O}(n^6)$.

7 Conclusion

We have provided a formal definition of bottom-up embedded push-down automata. We have also shown that finite-state control can be eliminated, obtaining a new definition in which transitions are in a form useful to describe compilation schemata for TAG and suitable for tabulation. The resulting definition has been shown to be equivalent to right-oriented linear indexed automata and a superset of bottom-up 2-stack automata with respect to the parsing strategies that can be described in both models of automata.

References

- [1] Miguel A. Alonso, Eric de la Clergerie, and Manuel Vilares. Automata-based parsing in dynamic programming for Linear Indexed Grammars. In A. S. Narin'yani, editor, *Proc. of DIALOGUE'97 Computational Linguistics and its Applications International Workshop*, pages 22–27, Moscow, Russia, June 1997.

- [2] Miguel A. Alonso, Eric de la Clergerie, and Manuel Vilares. A redefinition of Embedded Push-Down Automata. In *Proc. of 5th International Workshop on Tree Adjoining Grammars and Related Formalisms (TAG+5)*, pages 19–26, Paris, France, May 2000.
- [3] Miguel A. Alonso, Mark-Jan Nederhof, and Eric de la Clergerie. Tabulation of automata for tree adjoining languages. *Grammars*, Forthcoming.
- [4] Eric de la Clergerie and Miguel A. Alonso. A tabular interpretation of a class of 2-Stack Automata. In *COLING-ACL'98, 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Proceedings of the Conference*, volume II, pages 1333–1339, Montreal, Quebec, Canada, August 1998. ACL.
- [5] Eric de la Clergerie, Miguel A. Alonso, and David Cabrero Souto. A tabular interpretation of bottom-up automata for TAG. In *Proc. of Fourth International Workshop on Tree-Adjoining Grammars and Related Frameworks (TAG+4)*, pages 42–45, Philadelphia, PA, USA, August 1998.
- [6] Eric de la Clergerie and François Barthélemy. Information flow in tabular interpretations for generalized Push-Down Automata. *Theoretical Computer Science*, 199(1–2):167–198, 1998.
- [7] Aravind K. Joshi and Yves Schabes. Tree-adjoining grammars. In Grzegorz Rozenberg and Arto Salomaa, editors, *Handbook of Formal Languages. Vol 3: Beyond Words*, chapter 2, pages 69–123. Springer-Verlag, Berlin/Heidelberg/New York, 1997.
- [8] Bernard Lang. Towards a uniform formal framework for parsing. In Masaru Tomita, editor, *Current Issues in Parsing Technology*, pages 153–171. Kluwer Academic Publishers, Norwell, MA, USA, 1991.
- [9] Mark-Jan Nederhof. 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 1998.
- [10] Owen Rambow. *Formal and Computational Aspects of Natural Language Syntax*. PhD thesis, University of Pennsylvania, 1994. Available as IRCS Report 94-08 of the Institute of Research in Cognitive Science, University of Pennsylvania.
- [11] Yves Schabes and K. Vijay-Shanker. 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 1990. ACL.
- [12] K. Vijay-Shanker. *A Study of Tree Adjoining Grammars*. PhD thesis, University of Pennsylvania, January 1988. Available as Technical Report MS-CIS-88-03 LINC LAB 95 of the Department of Computer and Information Science, University of Pennsylvania.