

# Pushdown Automata and Context Free Language Parsing

*Prof. A. Morzenti*

**NB: up to slide 13 notions assumed to be known from other previous courses**

## PUSHDOWN AUTOMATA

- 1) stack auxiliary memory +  
input string with terminator  $\neg$

- 3) operations:

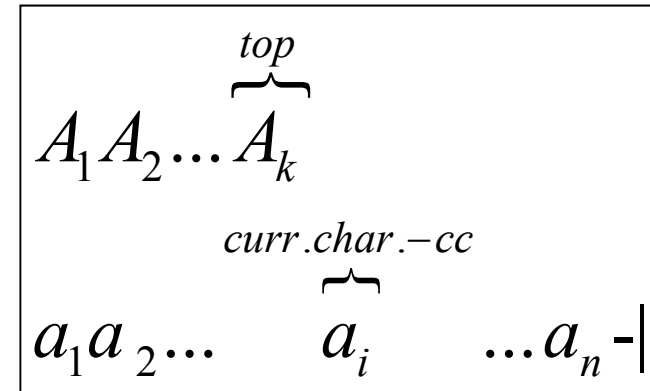
*push*( $B$ ), *push*( $B_1, B_2, \dots B_n$ ): push symbol(s) on top of the stack

*empty* test: a predicate that holds iff  $k = 0$

*pop*, if the stack is not empty, deletes  $A_k$

- 4)  $Z_0$  is the *initial (bottom)* stack symbol (can only be read)

- 5) configuration: current state, string portion from current character  $cc$ , stack content



## MOVE OF THE AUTOMATON:

- read  $cc$  and advance the head (*shift*), or (*spontaneous* move) do not advance the head
- read the top stack symbol (possibly  $Z_0$  if the stack is empty)
- based on the current char, state, and stack top symbol, go to a new state and replace the stack top symbol with a string (zero or more symbols)

## DEFINITION OF PUSHDOWN AUTOMATON

A pushdown automaton  $M$  (in general nondeterministic) is defined by:

1.  $Q$  *finite set of states of the control unit*
2.  $\Sigma$  *input alphabet*
3.  $\Gamma$  *stack alphabet*
4.  $\delta$  *transition function*
5.  $q_0 \in Q$  *initial state*
6.  $Z_0 \in \Gamma$  *initial stack symbol*
7.  $F \subseteq Q$  *set of final states*

TRANSITION FUNCTION:

domain:

range:

$Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma$

the powerset  $\wp(Q \times \Gamma^*)$  of  $Q \times \Gamma^*$

spontaneous move

nondeterminism

READING ( / scanning / shift) MOVE:  
(possibly nondeterministic)

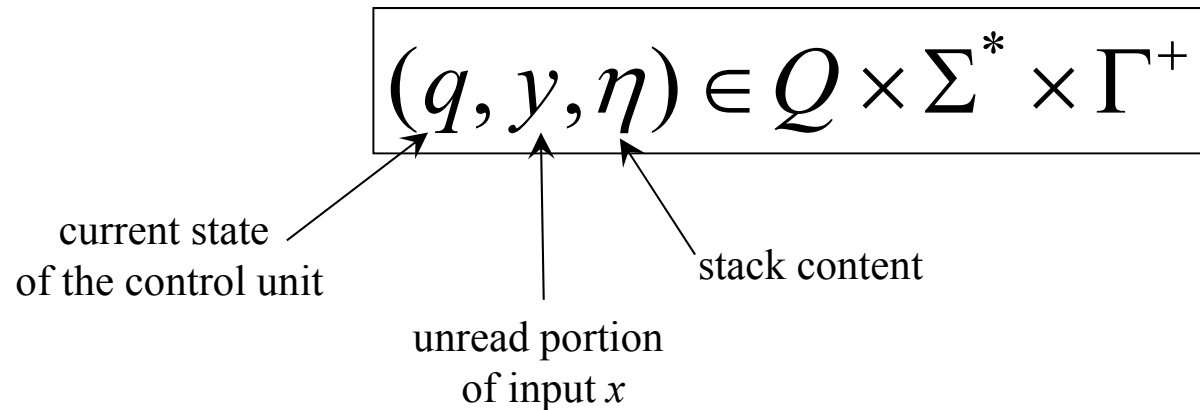
$$\delta(q, a, Z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots (p_n, \gamma_n)\}$$
$$\text{with } n \geq 1, a \in \Sigma, Z \in \Gamma, p_i \in Q, \gamma_i \in \Gamma^*$$

SPONTANEOUS MOVE:  
(possibly nondeterministic)

$$\delta(q, \varepsilon, Z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots (p_n, \gamma_n)\}$$
$$\text{with } n \geq 1, Z \in \Gamma, p_i \in Q, \gamma_i \in \Gamma^*$$

NONDETERMINISM: for a given triple (state, stack top, input) there are  $\geq 2$  possibilities among reading and spontaneous moves

INSTANTANEOUS CONFIGURATION OF MACHINE  $M$ : a triple



INITIAL CONFIGURATION:  $(q_0, x, Z_0)$

FINAL CONFIGURATION  $(q, \varepsilon, \eta)$  if  $q \in F$  (NB:  $\varepsilon$  means input completely scanned)

TRANSITION FROM ONE CONFIGURATION TO THE NEXT:

$$(q, y, \eta) \rightarrow (p, z, \lambda)$$

TRANSITION SEQUENCE:  $\xrightarrow{*}, \xrightarrow{+}$

current config.	next config	applied move
$(q, az, \eta Z)$	$(p, z, \eta \gamma)$	reading move $\delta(q, a, Z) = \{(p, \gamma), \dots\}$
$(q, az, \eta Z)$	$(p, az, \eta \gamma)$	spontaneous move $\delta(q, \varepsilon, Z) = \{(p, \gamma), \dots\}$

NB: A **string** is pushed: it can be  $\varepsilon$  (just a *pop* operation)  
or the same symbol previously on top (stack is unchanged)

A string  $x$  is recognized/accepted with final state if:

$$\boxed{(q_0, x, Z_0) \xrightarrow{+} (q, \varepsilon, \lambda)} \quad \begin{array}{l} q \in F \text{ and } \lambda \in \Gamma^* \\ \text{(no condition on } \lambda, \text{ it can be } \varepsilon, \text{ but not necessarily)} \end{array}$$

# STATE-TRANSITION DIAGRAM FOR PUSHDOWN AUTOMATA

Example: even-length palindromes accepted with final state by a (nondeterministic) PDA

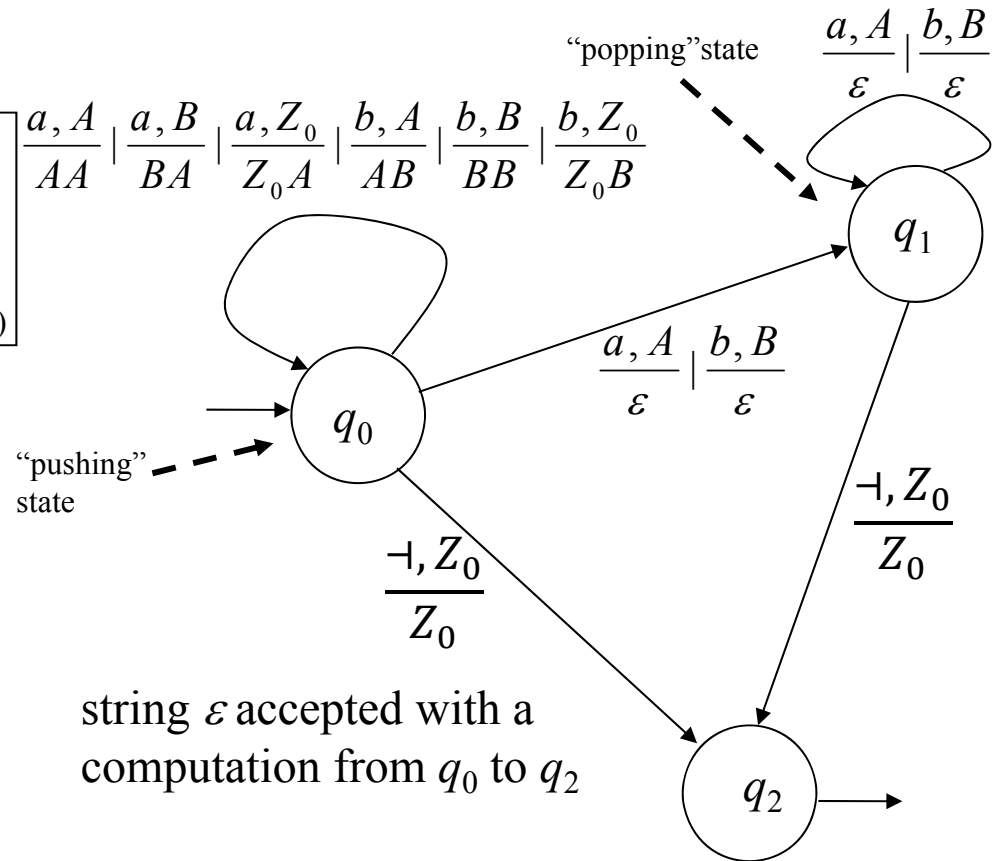
$$L = \{uu^R \mid u \in \{a,b\}^*\}$$

Stack	x	State	Comment
$Z_0$	$aa\vdash$	$q_0$	
$Z_0A$	$a\vdash$	$q_0$	
$Z_0AA$	$\vdash$	$q_0$	reject: no move defined for $(q_0, \vdash, A)$

“guesses” that  $|x|>2$

Stack	x	State	Comment
$Z_0$	$aa\vdash$	$q_0$	
$Z_0A$	$a\vdash$	$q_0$	
$Z_0$	$\vdash$	$q_1$	
$Z_0$	$\epsilon$	$q_2$	acceptance with final state

“guesses” that  $|x|=2$



## FROM THE GRAMMAR TO THE PUSHDOWN AUTOMATON

- 1) Grammar rules seen as instructions of a *non deterministic PDA* with a single state;  
*predictive* (goal oriented) analysis: stack used as a list of future actions.  
Only one state  $\Rightarrow$  called “*daisy automaton*”
- 2) Stack includes terminal and nonterminal symbols.  
Stack =  $A_1, \dots A_k$  goal is: read from input a string derived from  $A_k$ ,
- 3) The goal can be restated recursively in subgoals if  $A_k$  is a nonterminal from which other symbols are derived

**Initial goal is the axiom:** goal is to read a sentence, i.e., a string derived from axiom  $S$

Initially stack is  $Z_0 S$  and the input head on the first char of the input string

At each step the automaton chooses (nondeterministically) one of the applicable moves/rules

The string is recognized if the terminal  $\vdash$  is read with an empty stack



# FROM GRAMMAR RULES TO TRANSITIONS $A, B \in V, b \in \Sigma, A_i \in V \cup \Sigma$

Rule	Move	Comment
$A \rightarrow BA_1 \dots A_n \quad n \geq 0$	<b>if</b> $top = A$ <b>then</b> pop; push( $A_n \dots A_1 B$ ) <b>end if</b>	to recognize $A$ one must recognize $B A_1 \dots A_n$
$A \rightarrow bA_1 \dots A_n \quad n \geq 0$	<b>if</b> $cc = b \wedge top = A$ <b>then</b> pop; push( $A_n \dots A_1$ ); shift <b>end if</b>	$b$ is the first expected char and is read; $A_1 \dots A_n$ still to be accepted
$A \rightarrow \varepsilon$	<b>if</b> $top = A$ <b>then</b> pop <b>end if</b>	$\varepsilon$ deriving from $A$ is accepted
for each char. $b \in \Sigma$	<b>if</b> $cc = b \wedge top = b$ <b>then</b> pop; shift <b>end if</b>	$b$ is the first expected char and is read;
— — —	<b>if</b> $cc = \vdash \wedge \text{empty stack}$ <b>then</b> accept <b>end if</b> halt	string completely scanned, agenda completed

## Example – Rules and moves of the predictive recognizer

$$L = \{a^n b^m \mid n \geq m \geq 1\}$$

<u>Rule</u>	<u>Moves</u>
1. $S \rightarrow aS$	[ $\delta(q_0, a, S) = (q_0, S)$ ] <b>if</b> $cc = a \wedge top = S$ <b>then</b> pop; push(S); shift <b>end if</b>
2. $S \rightarrow A$	[ $\delta(q_0, \varepsilon, S) = (q_0, A)$ ] <b>if</b> $top = S$ <b>then</b> pop; push(A) <b>end if</b>
3. $A \rightarrow aAb$	[ $\delta(q_0, a, A) = (q_0, bA)$ ] <b>if</b> $cc = a \wedge top = A$ <b>then</b> pop; push(bA); shift <b>end if</b>
4. $A \rightarrow ab$	[ $\delta(q_0, a, A) = (q_0, b)$ ] <b>if</b> $cc = a \wedge top = A$ <b>then</b> pop; push(b); shift <b>end if</b>
5. --	[ $\delta(q_0, b, b) = (q_0, \varepsilon)$ ] <b>if</b> $cc = b \wedge top = b$ <b>then</b> pop; shift <b>end if</b>
6. --	[halt] <b>if</b> $cc = \perp \wedge empty\ stack$ <b>then</b> accept <b>end if</b>

Nondeterminism: between 1 and 2 (2 can also be chosen with input  $a$ ); between 3 and 4

String  $a^n b^m$ ,  $n \geq m \geq 1$  analyzed as  $a^{n-m} a^m b^m$ ,

“guessing nondeterministically” position where  $a^m b^m$  starts, by choice between moves 1 and 2

“guess” of position where  $a^m$  ends and  $b^m$  starts by choice between moves 3 and 4

this automaton accepts a string iff the grammar generates it  
 for every accepting computation there is a derivation and viceversa  
 the automaton *simulates the leftmost derivations*  
 the automaton must explore all derivations including nonaccepting ones  
 a string is accepted by several computations iff it is ambiguous

	$S \Rightarrow A \Rightarrow aAb \Rightarrow aabb$	
	accepting computation:	
	Stack	$x$
$\delta(q_0, \varepsilon, S) = (q_0, A)$	$Z_0 S$	$aabb-$
$\delta(q_0, a, A) = (q_0, bA)$	$Z_0 A$	$aabb-$
$\delta(q_0, a, A) = (q_0, b)$	$Z_0 bA$	$abb-$
$\delta(q_0, b, b) = (q_0, \varepsilon)$	$Z_0 bb$	$bb-$
$\delta(q_0, b, b) = (q_0, \varepsilon)$	$Z_0 b$	$b-$
$\delta(q_0, b, b) = (q_0, \varepsilon)$	$Z_0$	$-$

conversion grammar rules – PDA transitions is bidirectional:

It can be applied to obtain a grammar from a PDA

Therefore:

LANGUAGES DEFINED BY ***NONDETERMINISTIC*** PDA's  
AND CF LANGUAGES  
ARE A UNIQUE FAMILY

## VARIETIES OF PUSHDOWN AUTOMATA

PDA can be enriched in various ways, concerning internal states and acceptance conditions

3 possible accepting modes:

- with final state (stack content immaterial)
- empty stack (current state immaterial)
- combined: (final state and empty stack)

PROPERTY – These three accepting modes are equivalent
---

Another important PDA feature: absence of spontaneous loops and on-line functioning

PROPERTY: any PDA can be converted into an equivalent one

- with no cycles of spontaneous moves
- which can decide acceptance right after reading the last input symbol

## INTERSECTION OF REGULAR AND CONTEXT FREE LANGUAGES

We can now justify the statement:

Intersection of a *CF* and a *REG* language is *CF*

Given grammar *G* and automaton *A*, obtain PDA *M* accepting  $L(G) \cap L(A)$  as follows:

- 1) build PDA *N* accepting  $L(G)$  **with empty stack**
- 2) build machine *M* product of machines *N* and *A*,  
applying the known construction for finite automata adapted so that  
the product machine *M* manipulates the stack in the same way as component *N*

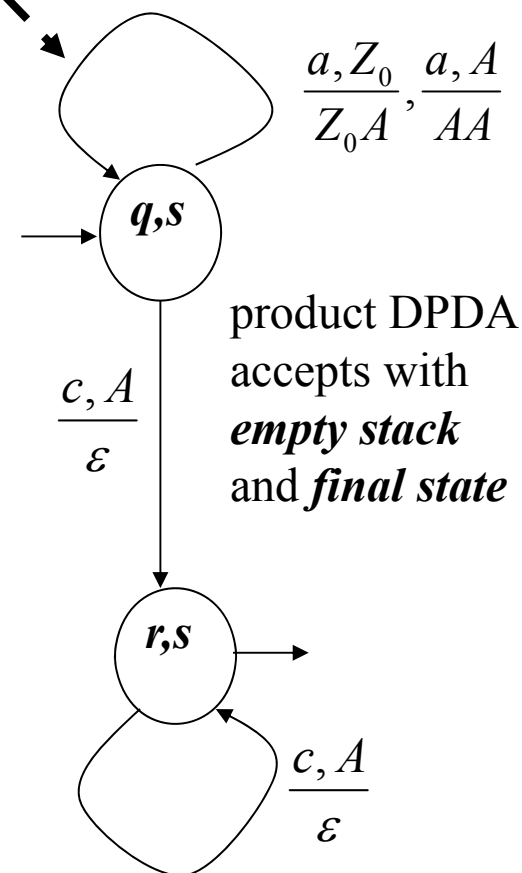
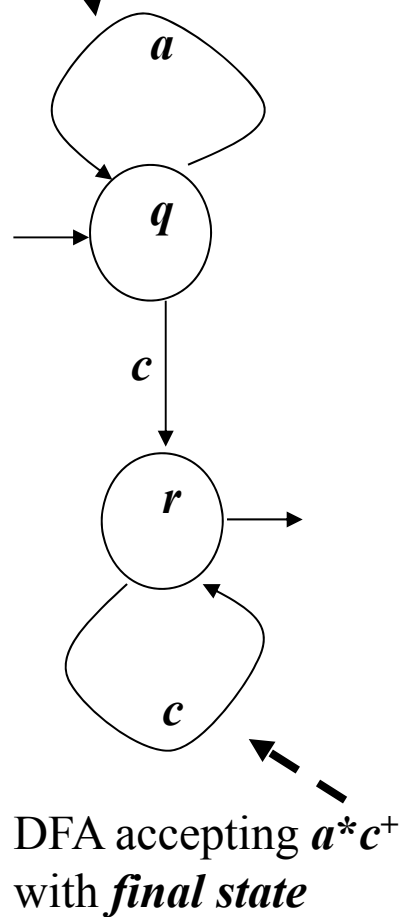
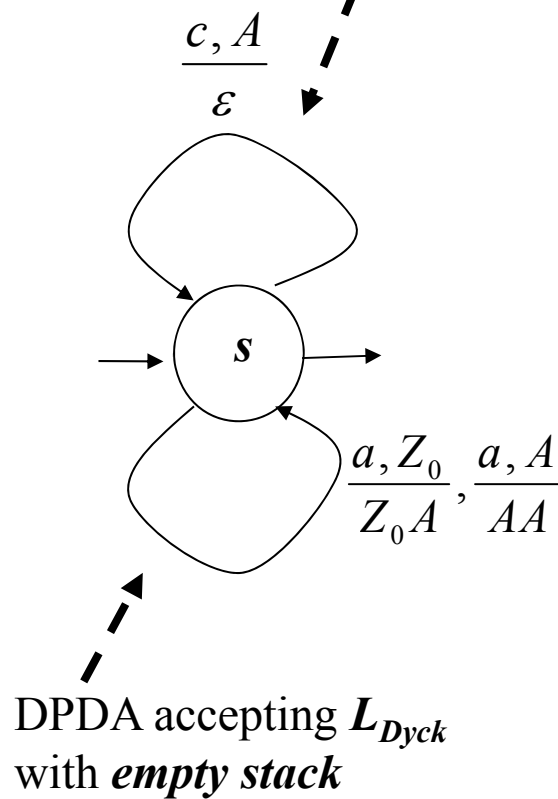
The machine thus constructed:

- internal states are the product of the state sets of the component machines
- accepts with **final state and empty stack**
- final states are those including a final state of the finite automaton *A*
- is deterministic if so are both machines *N* and *A*
- accepts exactly the strings of the intersection of the two languages

Example

$$L_{Dyck} \cap (a^*c^+) = \{a^n c^n \mid n \geq 1\}$$

Dyck language with one nest



## PUSHDOWN AUTOMATA AND DETERMINISTIC LANGUAGES (DET)

Only deterministic CF languages (those accepted by a deterministic PDA) are considered in language and compiler design due to efficiency reasons

Nondeterminism is absent if  $\delta$  is one-valued and also

if  $\delta(q, a, A)$  is defined for some  $a \in \Sigma$ , then  $\delta(q, \varepsilon, A)$  is not defined  
if  $\delta(q, \varepsilon, A)$  is defined then  $\delta(q, a, A)$  is not defined for any  $a \in \Sigma$

NB: therefore a **deterministic PDA** CAN have spontaneous moves



Relation between classe CF (Context Free Languages) and DET (deterministic ones)

Example: nondeterministic union of deterministic languages

$$L = \{a^n b^n \mid n \geq 1\} \cup \{a^n b^{2n} \mid n \geq 1\} = L' \cup L''$$

a PDA accepting  $x$  must push all  $a$ 's and, if  $x \in L'$  (e.g.,  $aabb$ ), pop one  $a$  for each  $b$ ;  
but if  $x \in L''$  (e.g.,  $aabbbb$ ), two  $b$ 's must be popped for each  $a$  in the stack  
The PDA does not know which alternative holds, it must try both

$L', L'' \in \text{DET}$ ,  $L', L'' \in \text{CF}$ ,  $L = L' \cup L''$ ,  $L \in \text{CF}$  but  $L \notin \text{DET}$ ,

hence **DET  $\subseteq$  CF and DET  $\neq$  CF**

## CLOSURE PROPERTIES OF DETERMINISTIC CF LANGUAGES

We denote as  $L$ ,  $D$ , and  $R$  a language in the family  $CF$ ,  $DET$  and  $REG$

Operation	Property	(Property already known)
Reflection	$D^R \notin DET$	$D^R \in CF$
Star	$D^* \notin DET$	$D^* \in CF$
Complement	$\neg D \in DET$	$\neg L \notin CF$
Union	$D_1 \cup D_2 \notin DET, D \cup R \in DET$	$D_1 \cup D_2 \in CF$
Concatenation	$D_1.D_2 \notin DET, D.R \in DET$	$D_1.D_2 \in CF$
Intersection	$D \cap R \in DET$	$D_1 \cap D_2 \notin CF$

NB: the typical operations on languages ( $R$ ,  $*$ ,  $\cup$ ,  $\cdot$ ) **DO NOT** preserve determinism

## SYNTAX ANALYSIS

Given a grammar  $G$ , the syntax analyzer (*parser*)

- reads the source string  $x$  and
- if  $x \in L(G)$ ,
  - accepts and possibly outputs a syntax tree or (equivalently) a derivation;
- otherwise it stops signalling an error (diagnosis)

## TOP-DOWN AND BOTTOM-UP ANALYSIS

One given tree in general corresponds to various derivations (left, right, ....)

The two most important types of parsers characterized by the type of identified derivation: **left** or **right**, and **order** of the tree and derivation construction

TOP-DOWN ANALYSIS: builds  
a **left** derivation in **direct order**  
syntax tree: **from root to leaves**, through *expansions*

BOTTOM-UP ANALYSIS: builds  
**right** derivation in **reverse order**  
syntax tree: **from leaves to root**, through *reductions*

Example – top-down analysis of sentence:

***a b b b a a***

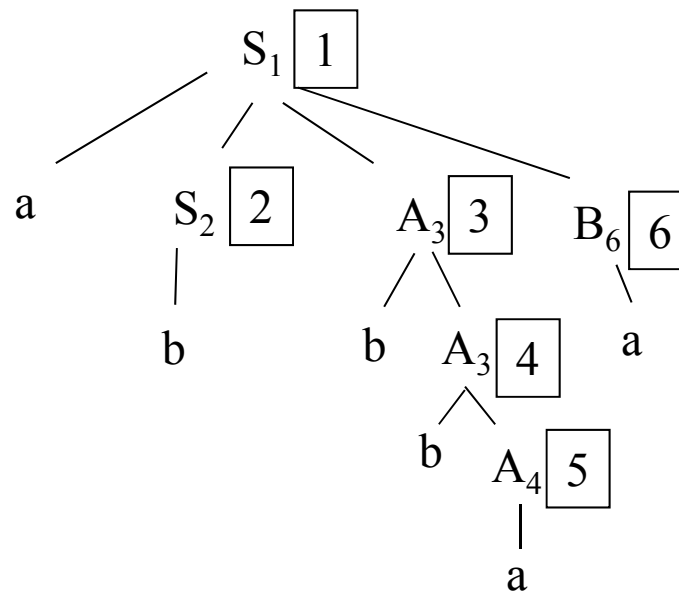
framed numbers = order in rule application

subscripts of nonterminals in the tree = applied rule

Leftmost: always expanded first nonterm. from the left

1. $S \rightarrow aSAB$	2. $S \rightarrow b$
3. $A \rightarrow bA$	4. $A \rightarrow a$
5. $B \rightarrow cB$	6. $B \rightarrow a$

$S \Rightarrow aSAB \Rightarrow abAB \Rightarrow abbAB \Rightarrow abbbAB \Rightarrow abbbaB \Rightarrow abbbbaa$



Example – bottom-up analysis of sentence:

***a b b b a a***

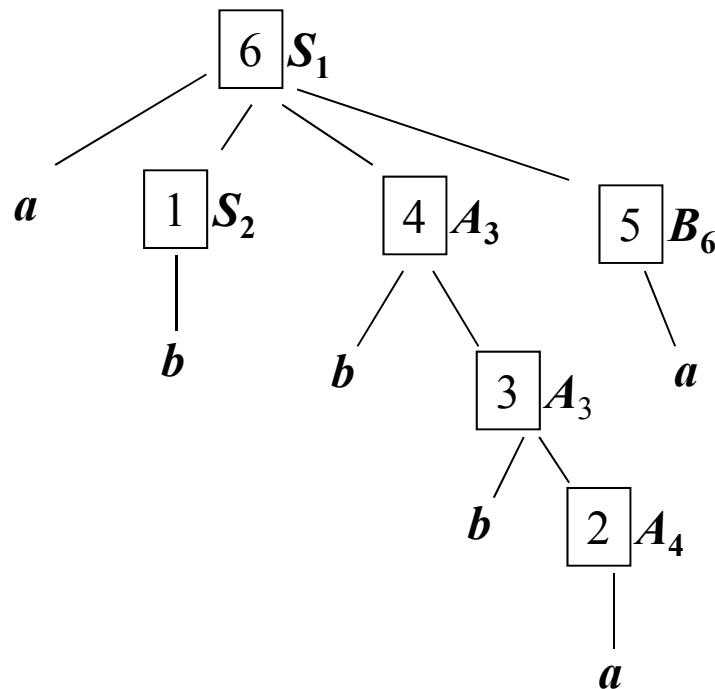
1. $S \rightarrow aSAB$	2. $S \rightarrow b$
3. $A \rightarrow bA$	4. $A \rightarrow a$
5. $B \rightarrow cB$	6. $B \rightarrow a$

Derivation (rightmost n.t. always expanded):

$S \Rightarrow aSAB \Rightarrow aSAa \Rightarrow aSbAa \Rightarrow aSbbAa \Rightarrow aSbbaa \Rightarrow a b b b a a$

framed numbers = order in reduction sequence

subscripts of nonterminals in the tree = applied rule



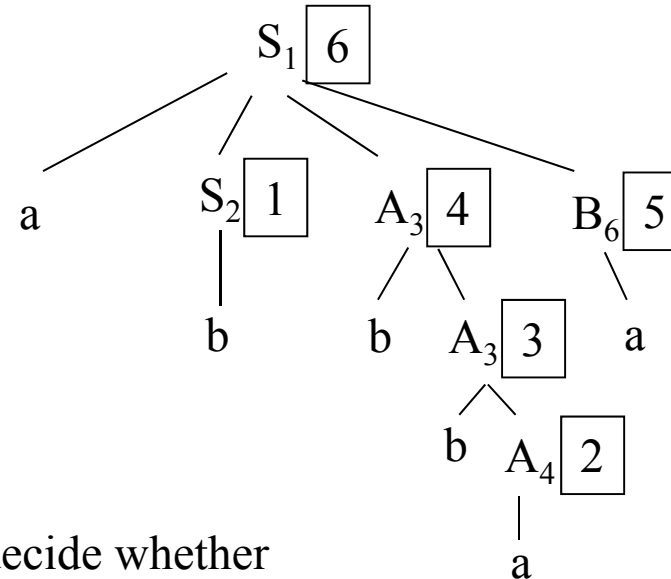
right parts of rules are ***reduced*** as they are scanned, first in the sentence, then in the phrase forms obtained after the ***reductions***. The process terminates when the entire string is reduced to the axiom

NB: reconstruction of a ***right derivation*** in ***reverse order***

Bottom-up analysis: the **reduction** operations ( $\text{>rd>}$  below) transform part of a phrase form into a string  $\alpha \in (\Sigma \cup V)^*$ , called «*viable prefix*», that may include the result of previous reductions (and is stored in the parser's stack). Here below, the '^' shows the head position (right of ^ is the unread string), **underline** shows the part to be reduced, called «*handle*»

$\underline{a}bbbaa \text{>rd>} \quad aSbb\underline{aa} \text{>rd>} \quad aSb\underline{bA}a \text{>rd>} \quad aS\underline{bA}a \text{>rd>} \quad aSA\underline{a} \text{>rd>} \quad \underline{aSAB} \text{>rd>} \quad S$

1. $S \rightarrow aSAB$	2. $S \rightarrow b$
3. $A \rightarrow bA$	4. $A \rightarrow a$
5. $B \rightarrow cB$	6. $B \rightarrow a$



At each step of the analysis, the parser must decide whether

- to continue and read the next symbol (shift), or
- to build a subtree from a portion of the viable prefix

it chooses based on the symbol(s) coming after the current one (*lookahead*)

# GRAMMARS AS NETWORKS OF FINITE AUTOMATA

Suppose  $G$  in extended form: every nonterminal has a unique rule

$A \rightarrow \alpha$  with  $\alpha$  regular expression on terminals and nonterminals

$\alpha$  defines a regular language, hence there exists a finite automaton  $M_A$  that accepts  $L(\alpha)$

any transition of  $M_A$  labeled by a nonterminal  $B$

is interpreted as a «call» of an automaton  $M_B$  (if  $B=A$  then recursive call)

let us call

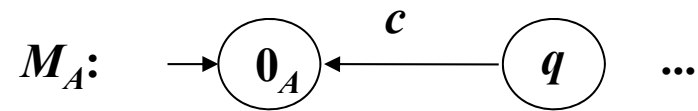
- “***machines***” the finite automata of the various nonterminals,
- “***automaton***” the PDA that accepts and parses the language  $L(G)$
- “***net***” the set of all machines
  
- $L(q)$  the set of terminal strings generated along a path of a machine, starting from state  $q$  (possibly including calls of other machines) and reaching a final state (examples follow)



We set a further requirement on machines corresponding to nonterminals

The initial state  $0_A$  of machine  $A$  is not visited after the start of the computation

No machine  $M_A$  includes a transition like



Requirement very easy to “implement” in case it is not satisfied ...

... it suffices to add one state (to be the new initial state) and a few arcs ...

$\Rightarrow$  the machine is not minimal, but only one state has been added

Automata satisfying this condition are called

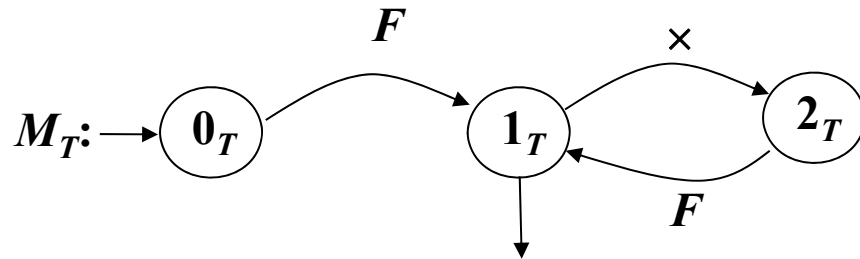
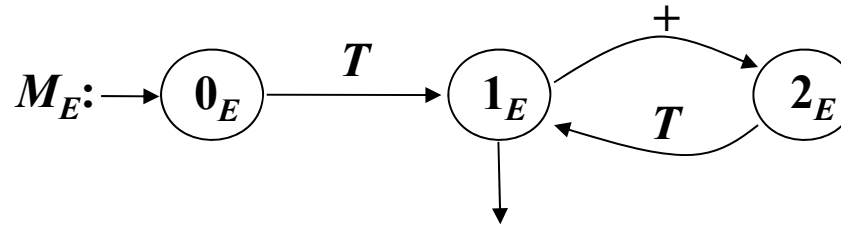
**normalized** or with initial state **non recirculating** or **non reentrant**

NB: it is **NOT forbidden** that the initial state be also final

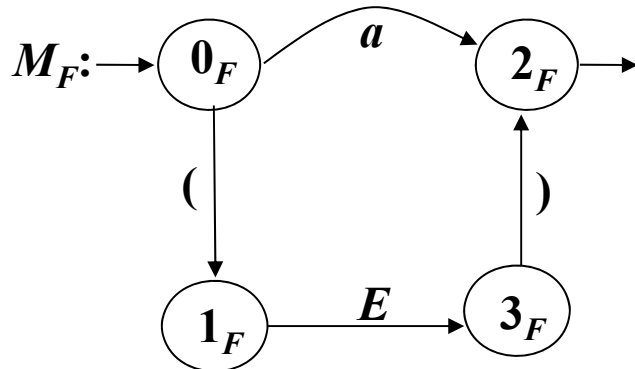
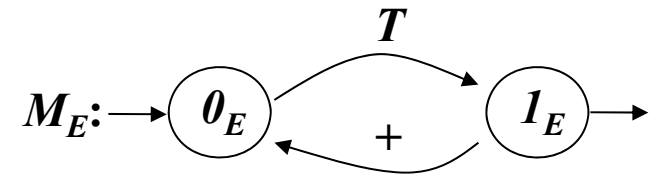
(which is necessary if  $\varepsilon \in L(A)$  )

# Example – Arithmetic expressions

$$\begin{aligned} E &\rightarrow T (+ T)^* \\ T &\rightarrow F ( \times F )^* \\ F &\rightarrow a \mid ' ( E ) ' \end{aligned}$$



NB:  $M_E$  minimal not normalized



$$L(2_F) = \{ \varepsilon \}$$

$$L(3_F) = \{ ) \}$$

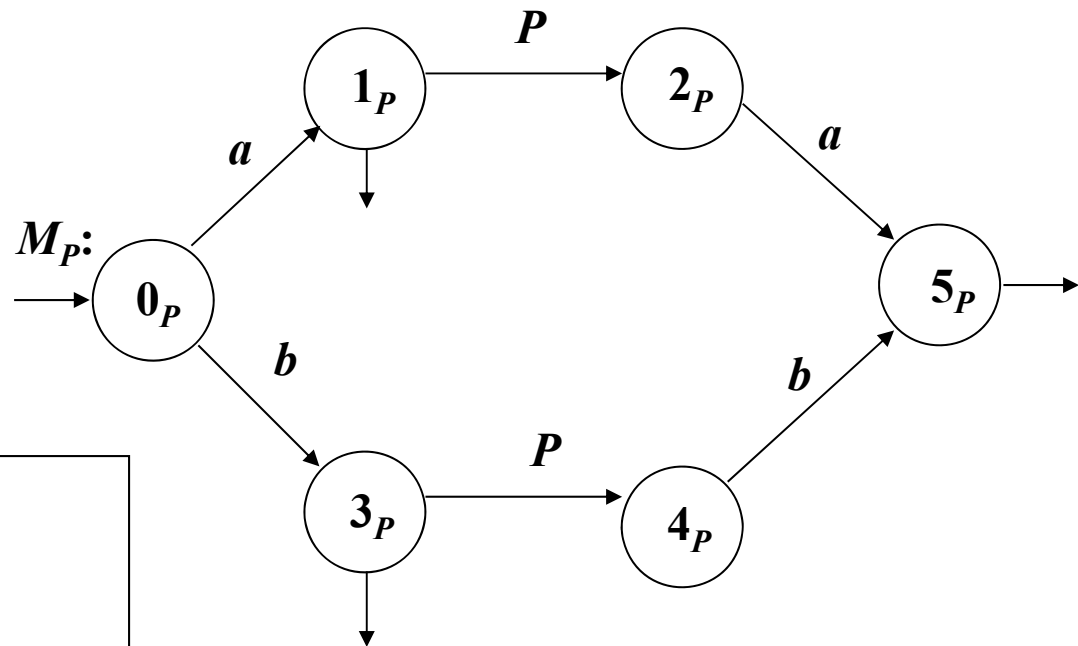
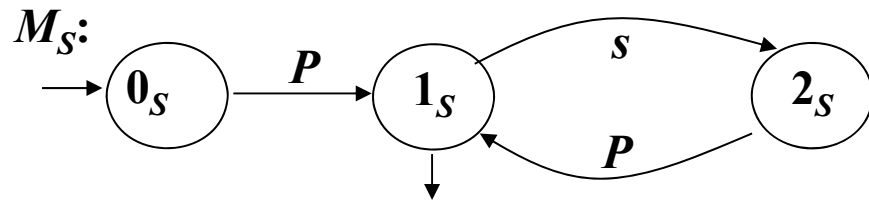
$$L(0_F) = L(F) \text{ the language generated from } F$$

$$L(1_F) = L(E) \cdot \{ ) \}$$

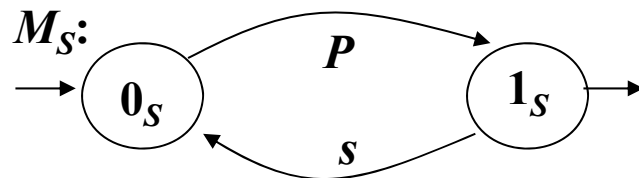
Example – List of odd-length palindromes, separated by ‘s’

$$S \rightarrow P ( s P )^*$$

$$P \rightarrow aPa \mid bPb \mid a \mid b$$



**NB:  $M_S$  minimal not normalized**



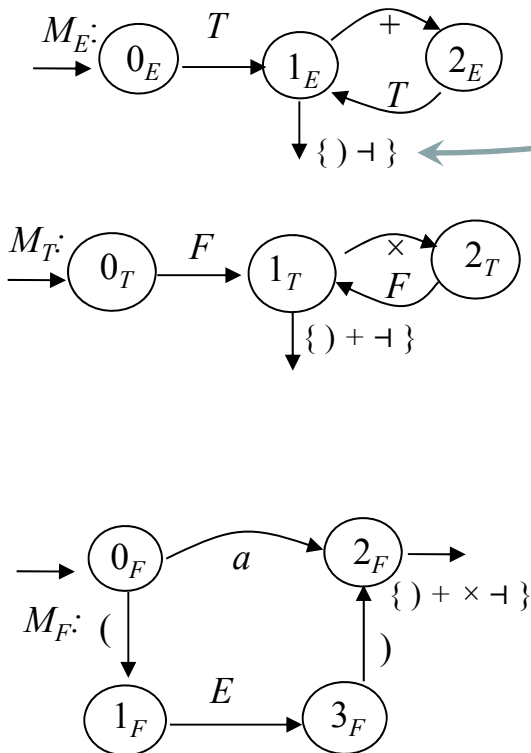
# TOP-DOWN ANALYSIS WITH RECURSIVE PROCEDURES

Simple and elegant: the procedure code reflects the transitions of machines

$$\begin{aligned} E &\rightarrow T (+ T)^* \\ T &\rightarrow F (\times F)^* \\ F &\rightarrow a \mid ' ( E ) ' \end{aligned}$$

Example – Arithmetic expressions

When in the automaton there are *bifurcations*, to decide which one to follow one must consider the symbols encountered on every arc, including those encountered when «exiting the machine»



$cc$  is current char,  
symbol under input head  
updated by the *next* procedure

```

procedure E
  call T;
  while (cc = '+')
    cc:=next;
    call T;
  end;
  if (cc ∈ { } ) + { } )
    return;
  else error;
endif
end E;
    
```

```

procedure T
  call F;
  while (cc = '×')
    cc:=next;
    call F;
  end;
  if (cc ∈ { } ) + + { } )
    return;
  else error;
endif
end T;
    
```

```

procedure F
  if (cc = 'a')
    cc:=next;
  elseif (cc='(')
    cc:=next;
    call E;
    if (cc = ')')
      cc:=next;
    else error;
  endif;
  else error;
endif;
  if (cc ∈ { } ) + × + { } )
    return;
  else error;
endif;
end F;
    
```

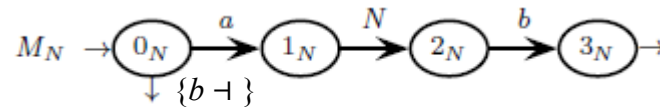
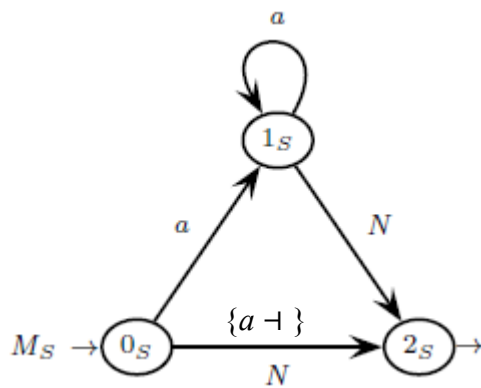
NB: this method **does not work** if

the sets of possible symbols on arcs outgoing from a single state (the «guide sets»)  
are **not disjoint**

Example:  $S \rightarrow a^* N$

$N \rightarrow aNb \mid \varepsilon$

$L = \{ a^n b^m \mid n \geq m \geq 0 \}$



```

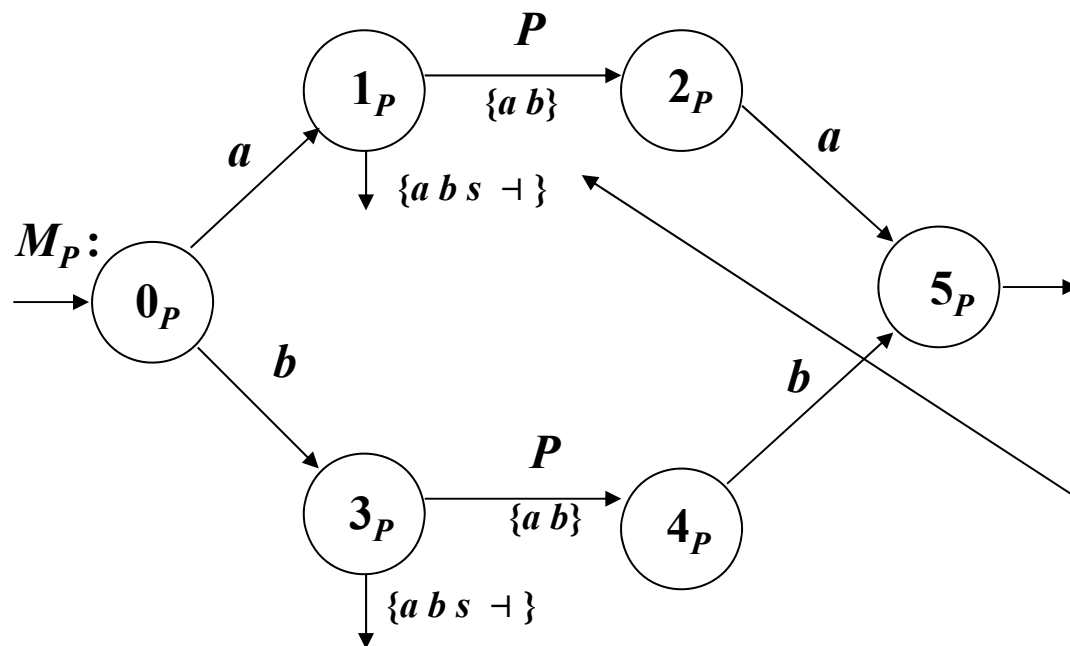
procedure S
  if (cc='a')
    ?? call N; ??
    ?? cc:=next; ??
  ...
end S;
  
```

Another Example – List of odd-length palindromes, separated by an ‘s’

(NB: it is a nondeterministic language: even the bottom-up method does not work)

$$S \rightarrow P ( s P )^*$$

$$P \rightarrow aPa \mid bPb \mid a \mid b$$



```

procedure P
  if ( cc = 'a' )
    cc:=next;
    if ( cc ∈ {a b} )
      ?? call P; ??
      ?? return; ??
    ...
  end P;
  
```

the recursive  
procedure cannot  
decide based on the  
next symbol