

Pushdown Automata and Context Free Language Parsing

Prof. A. Morzenti

NB: up to slide 9 notions assumed to be known from other previous courses

PUSHDOWN AUTOMATA

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

- 3) operations:

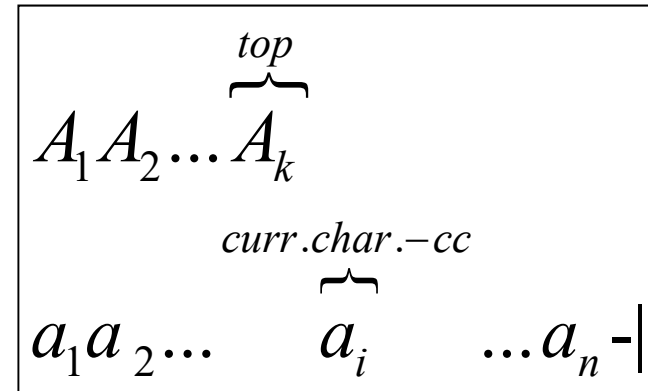
push(B), *push*($B_1, B_2, \dots B_n$): push symbol(s) on top of the stack

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. Σ *input alphabet*
3. Γ *stack alphabet*
4. δ *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)

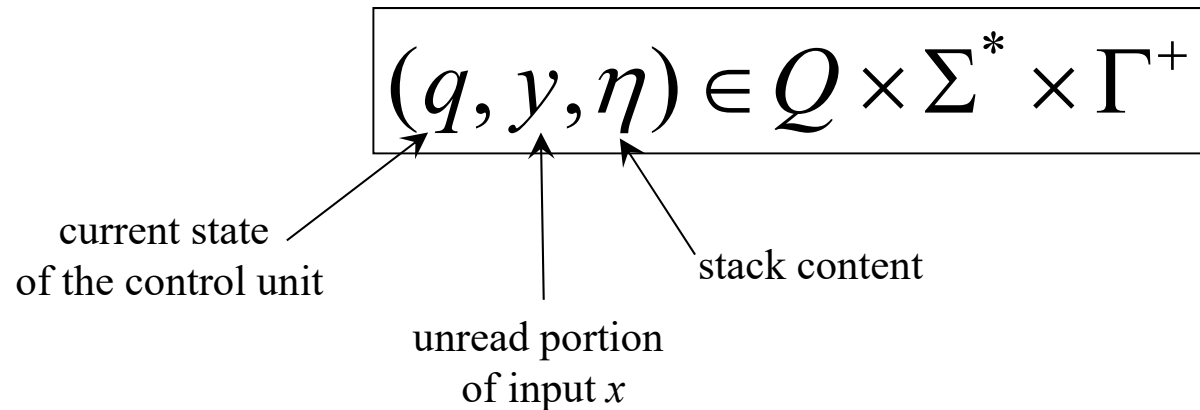
$$\delta(q, a, Z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots, (p_n, \gamma_n)\}$$
$$\text{with } n \geq 1, a \in \Sigma, Z \in \Gamma, p_i \in Q, \gamma_i \in \Gamma^*$$

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 ≥ 2 possibilities among reading and spontaneous moves

INSTANTANEOUS CONFIGURATION OF MACHINE M : a triple



INITIAL CONFIGURATION: (q_0, x, Z_0)

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

TRANSITION FROM ONE CONFIGURATION TO THE NEXT:

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

TRANSITION SEQUENCE: $\overset{*}{\rightarrow}, \overset{+}{\rightarrow}$

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 ε (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

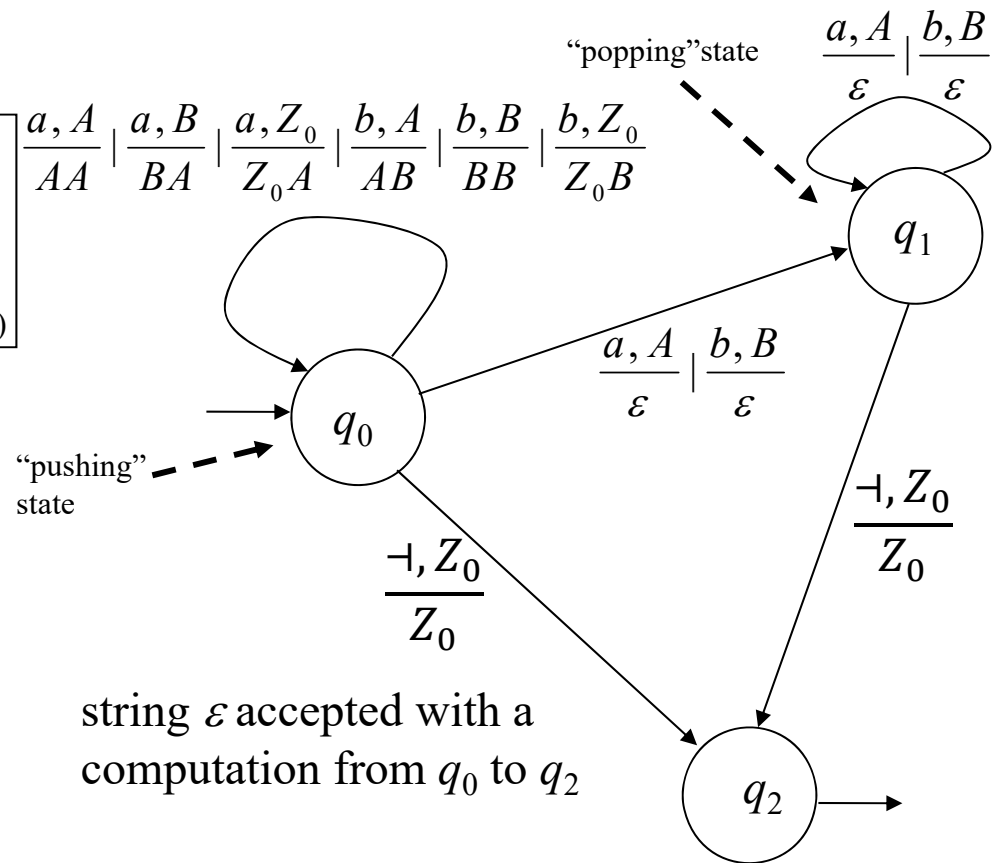
$$L = \{uu^R \mid u \in \{a,b\}^*\}$$

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	ε	q_2	acceptance with final state

“guesses” that $|x|=2$



There exist conversion procedures

Context-free grammar \Leftrightarrow Pushdown Automaton (possibly nondeterministic)

Therefore:

LANGUAGES
ACCEPTED BY *NONDETERMINISTIC* PDA's
AND
GENERATED BY CF GRAMMARS
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

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 δ is one-valued and also, $\forall q \in Q$ and $\forall A \in \Gamma$

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

But do not worry: there is a procedure to determine if a given free *grammar* $\in \text{DET}$

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

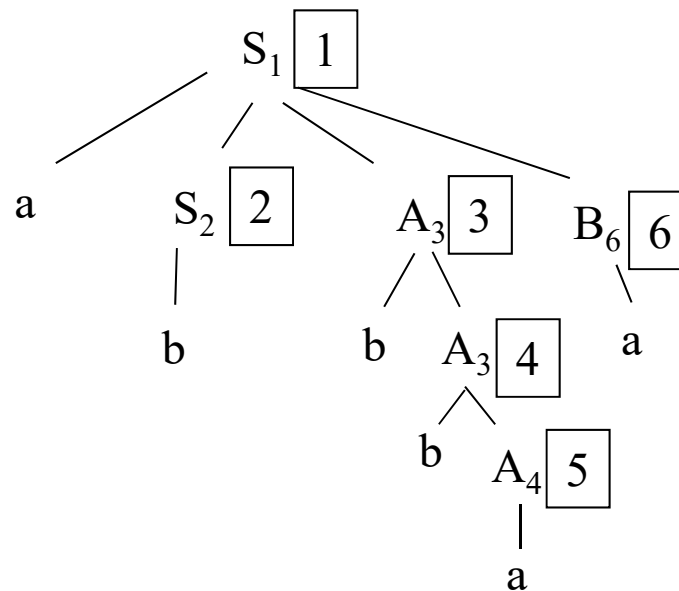
framed numbers = order in rule application

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

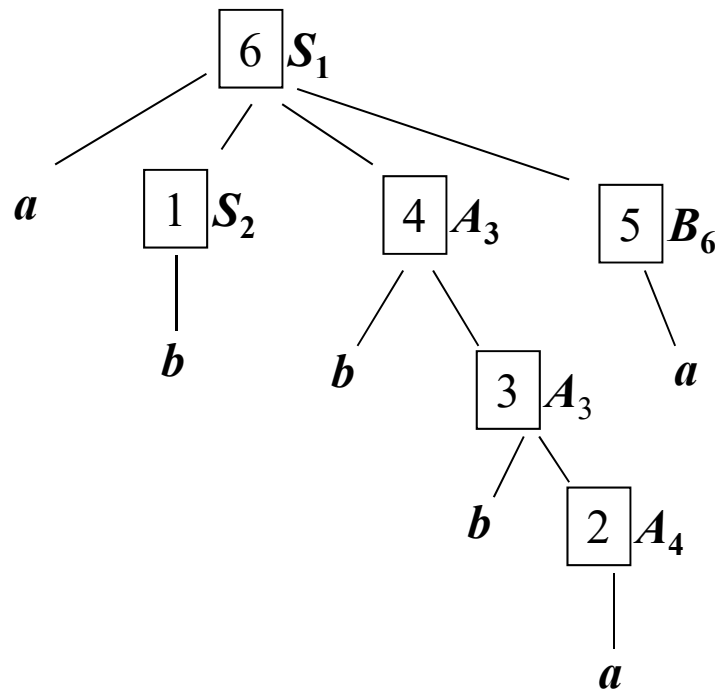
Derivation (rightmost n.t. always expanded):

$S \Rightarrow aSAB \Rightarrow aSAa \Rightarrow aSbAa \Rightarrow aSbbAa \Rightarrow aSbbaa \Rightarrow abbbbaa$

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

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

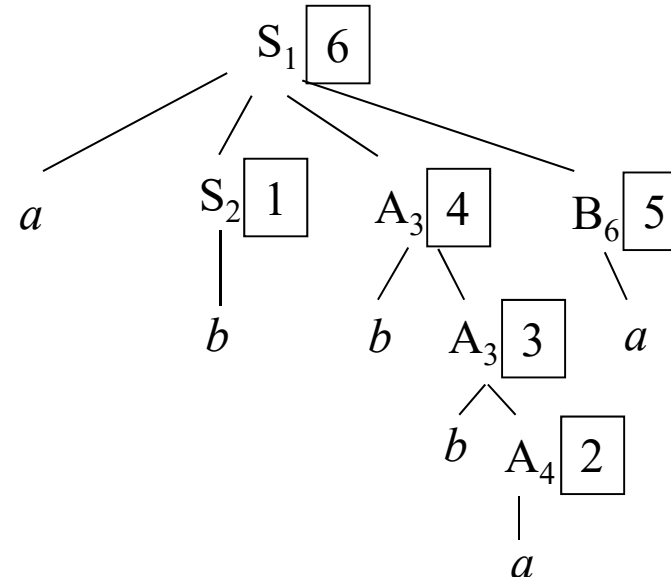
NB: reconstruction of a ***right derivation*** in ***reverse order***

Why in reverse order? Because input string is ***read from left***

Bottom-up analysis: the **reduction** operations (>rd> below) transform prefix 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 reduced part, called «*handle*»

$a\underline{b}bbaa \text{>rd>} aSbb\underline{a}a \text{>rd>} aSb\underline{bA}a \text{>rd>} aS\underline{bA}a \text{>rd>} aSA\underline{a} \text{>rd>} \underline{aSAB} \text{>rd>} S$

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*)

NB: the reason for the choice is not explained here: it will be in coming lessons

EXTENDED GRAMMARS REPRESENTED AS NETWORKS OF FINITE AUTOMATA

We use G in extended form: every nonterminal A has a unique rule

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

α 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

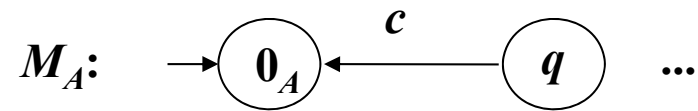
- “***machines***” the finite automata of the various nonterminals,
- “***automaton***” the PDA that accepts and parses the language $L(G)$
- “***net***” the set of all machines

- $L(q)$ 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 (\Rightarrow if q is final then $L(q)=\{\epsilon\}$)

We set a further requirement on machines corresponding to nonterminals

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

No machine M_A may include 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 transitions from it ...

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

Automata satisfying this condition are called

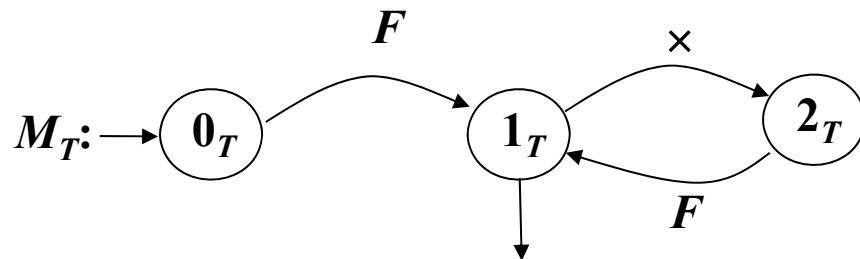
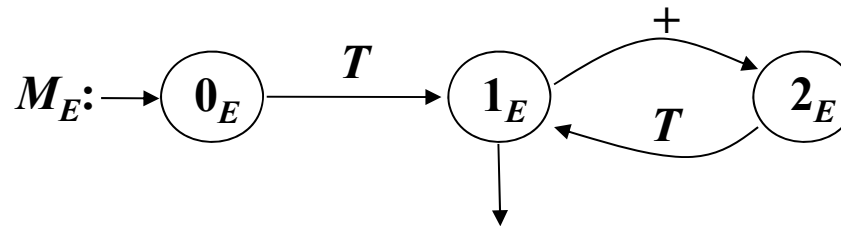
normalized or with initial state **non recirculating** or **non reentrant**

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

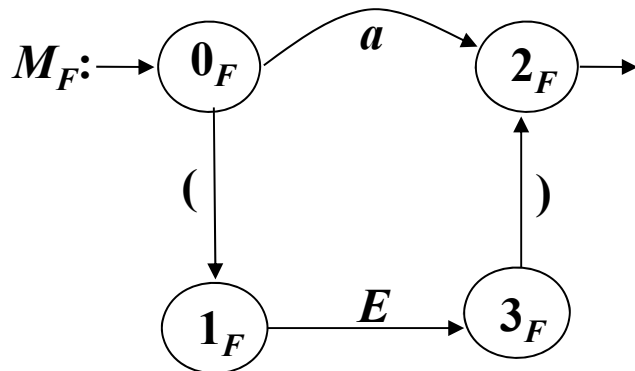
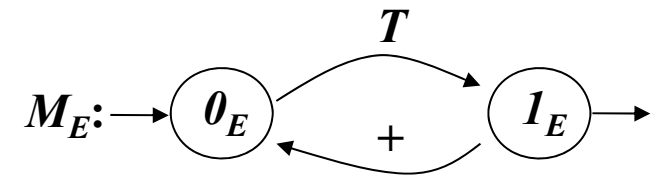
(which is necessary if $\varepsilon \in L(A)$)

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



$L(2_F) = \{ \varepsilon \}$ reason: state 2_F is final

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

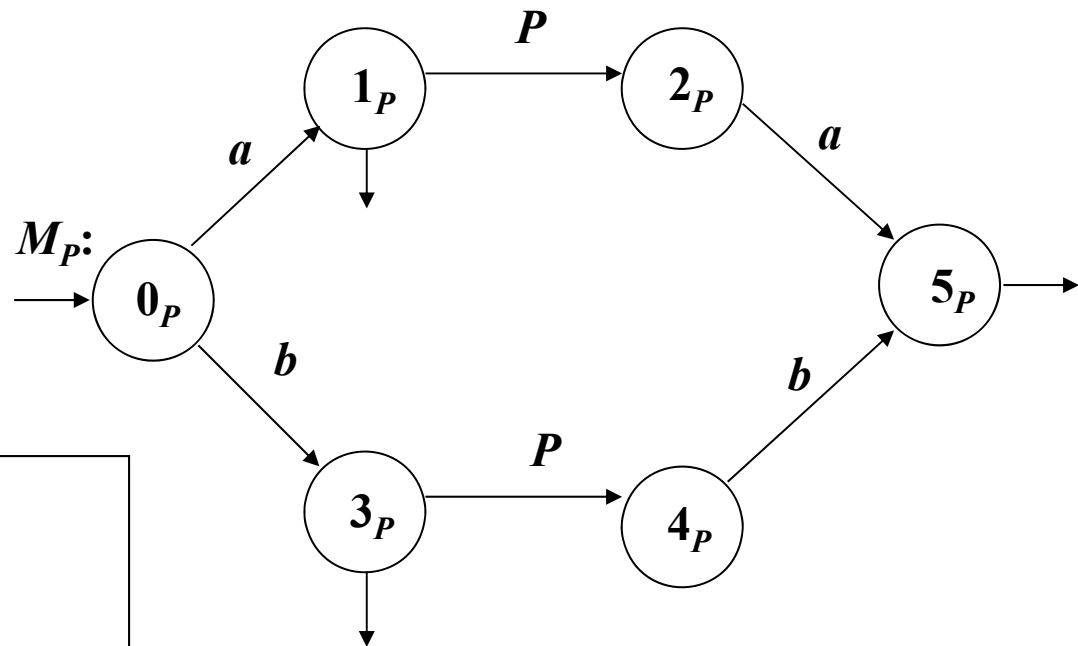
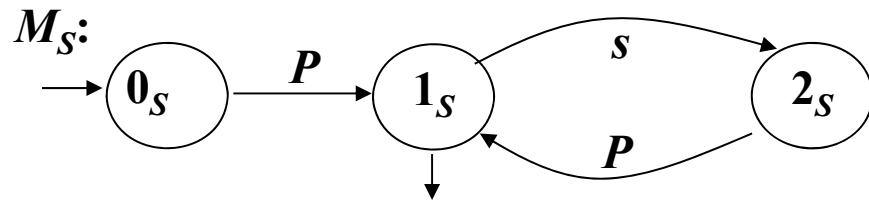
$L(0_F) = L(F)$ the language generated from F

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

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

