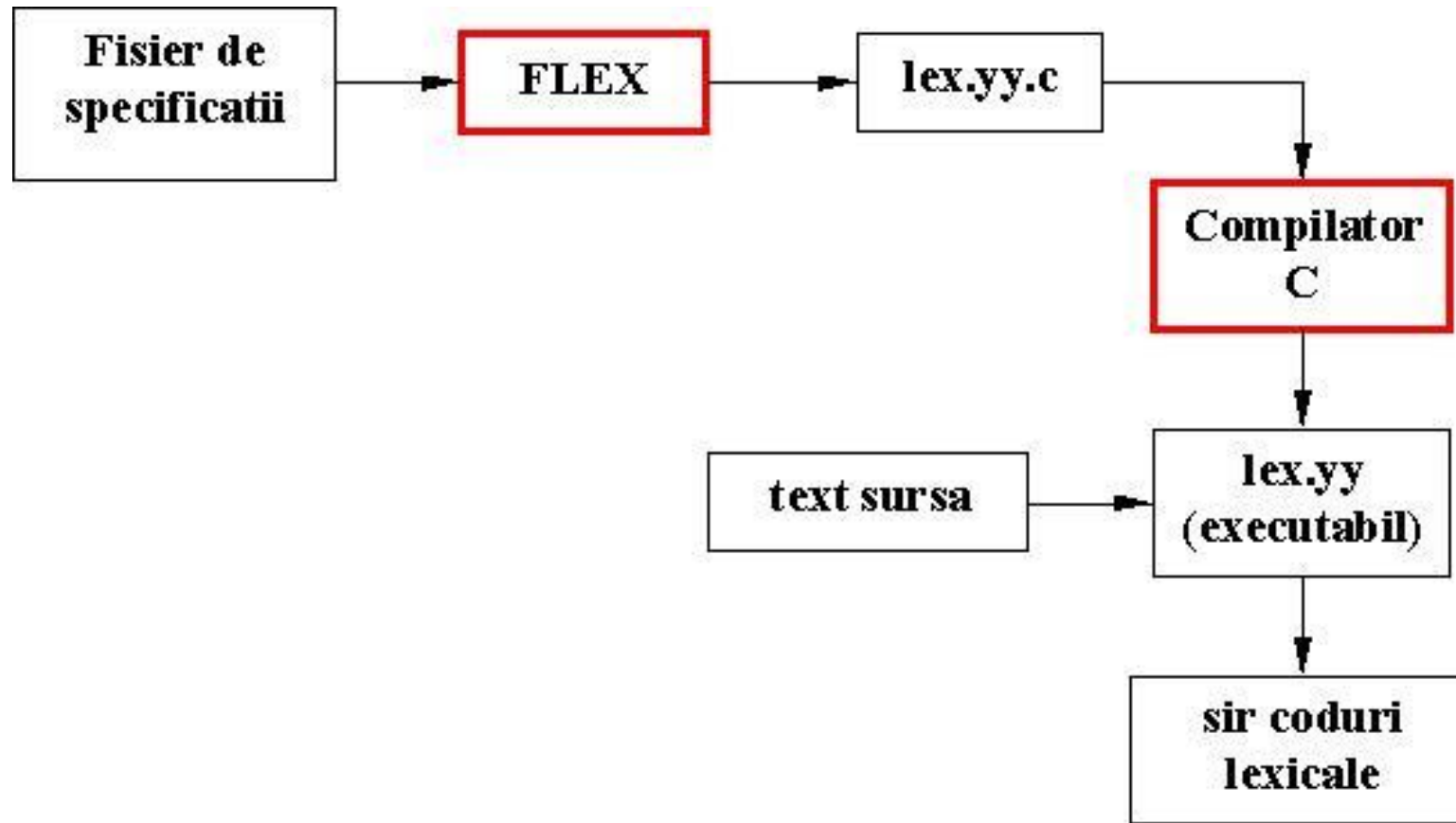# Course 3

# Back to compiler construction

# Scanning & Parsing Tools

- Scanning => lex

- Parsing => yacc //later

# Lex – Unix utilitary (flex – Windows version)

# INPUT FILE FORMAT

- The file containing the specification is a text file, that can have any name. Due to historic reasons we recommend the extension **.lxi**.

- Consists of 3 sections separated by a line containing %%:

```
definitions
%%
rules
%%
user code
```

*Example 1:*

```
%%

    username printf( "%s", getlogin() );
```

**specifies a scanner that, when finding the string
"`username`", will replace it with the user login name**

# Definition Section:

- - declarations of simple *name definitions* (used to simplify the scanner specification), of the form

```
name definition
```

- where:

- **name** is a word formed by one or more letters, digits, '_' or '-', with the remark that the first character MUST be letter or '_' and must be written on the FIRST POSITION OF THE LINE.

- **definition** is a regular expression and is starting with the first nonblank character after name until the end of line.

- declarations of *start conditions*.

# Rules Section

- to associate semantic actions with regular expressions. It may also contain user defined C code, in the following way:

**pattern action**

where:

- **pattern** is a regular expression, whose first character MUST BE ON THE FIRST POSITION OF THE LINE; see RegExp file

- **action** is a sequence of one or more C statements that MUST START ON THE SAME LINE WITH THE PATTERN. If there are more than one statements they will be nested between {}. In particular, the action can be a void statement.

# User Defined Code Section:

- Is optional (if is missing, then the separator %% following the rules section can also miss). If it exists, then its containing user defined C code is copied without any change at the end of the file lex.yy.c.

- Normally, in the user defined code section, one may have:

- function *main()* containing call(s) to *yylex()*, if we want the scanner to work autonomously (for ex., to test it);

- other called functions from *yylex()* (for ex. *yywrap()* or functions called during actions); in this case, the user code from definitions section must contain: either prototypes, either ***#include*** directives of the headers containing the prototypes

Launching the execution:

lex [*option*] [*name_specification _file*]


where *name_specification _file is an input file (implicitly, stdin)*
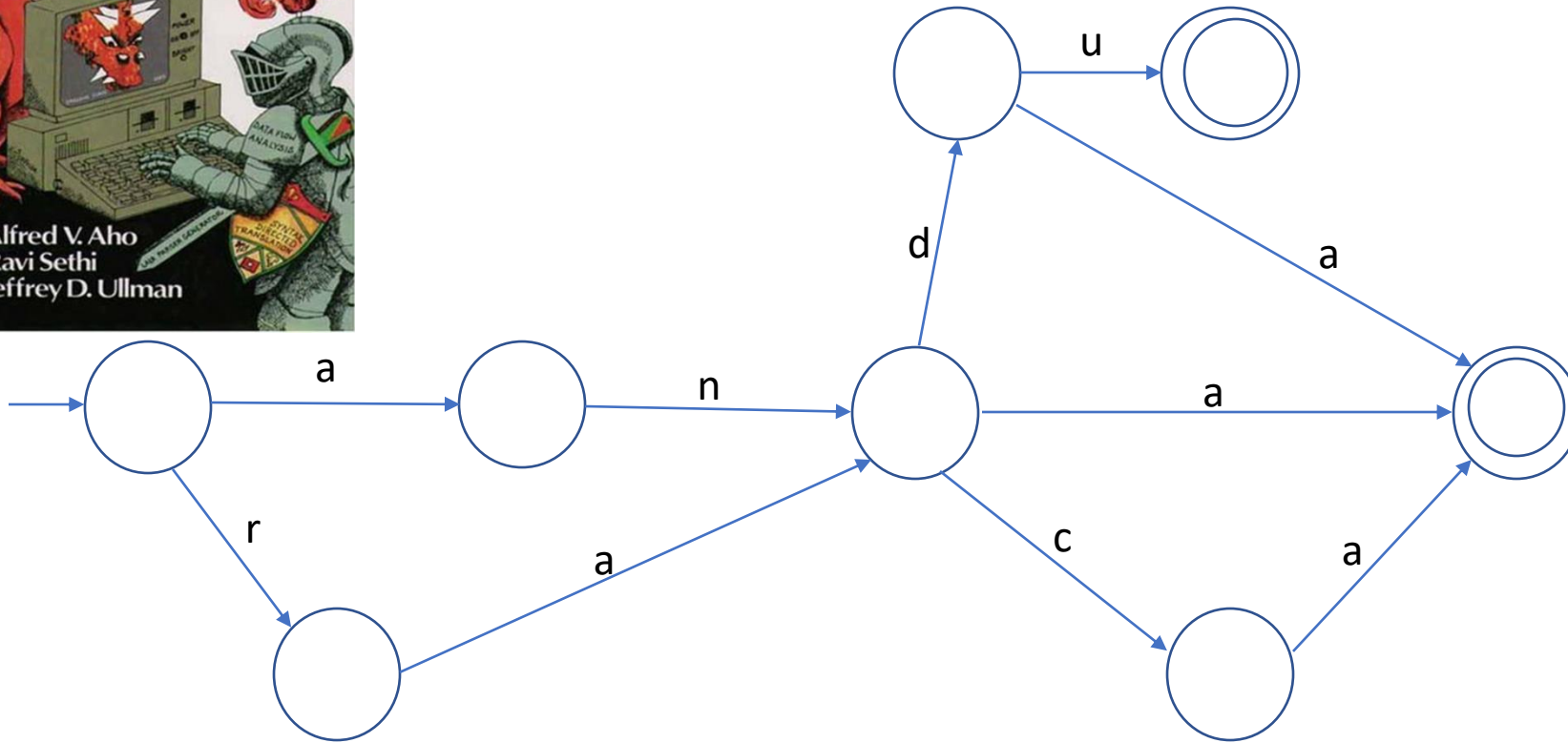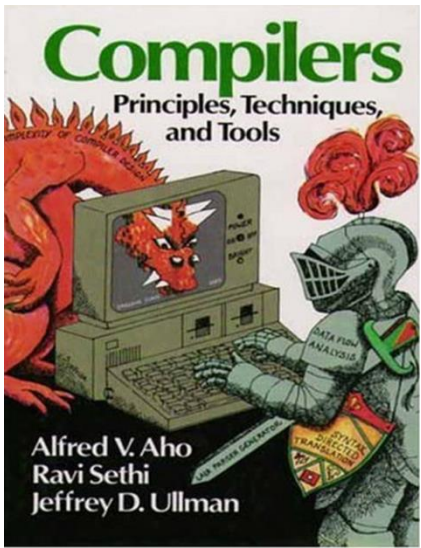

**$ lex spec.lxi**

**$ gcc lex.yy.c -o your_lex**

**$ your_lex<input.txt**

**options:** http://dinosaur.compilertools.net/flex/manpage.html

# Example

# Formal Languages
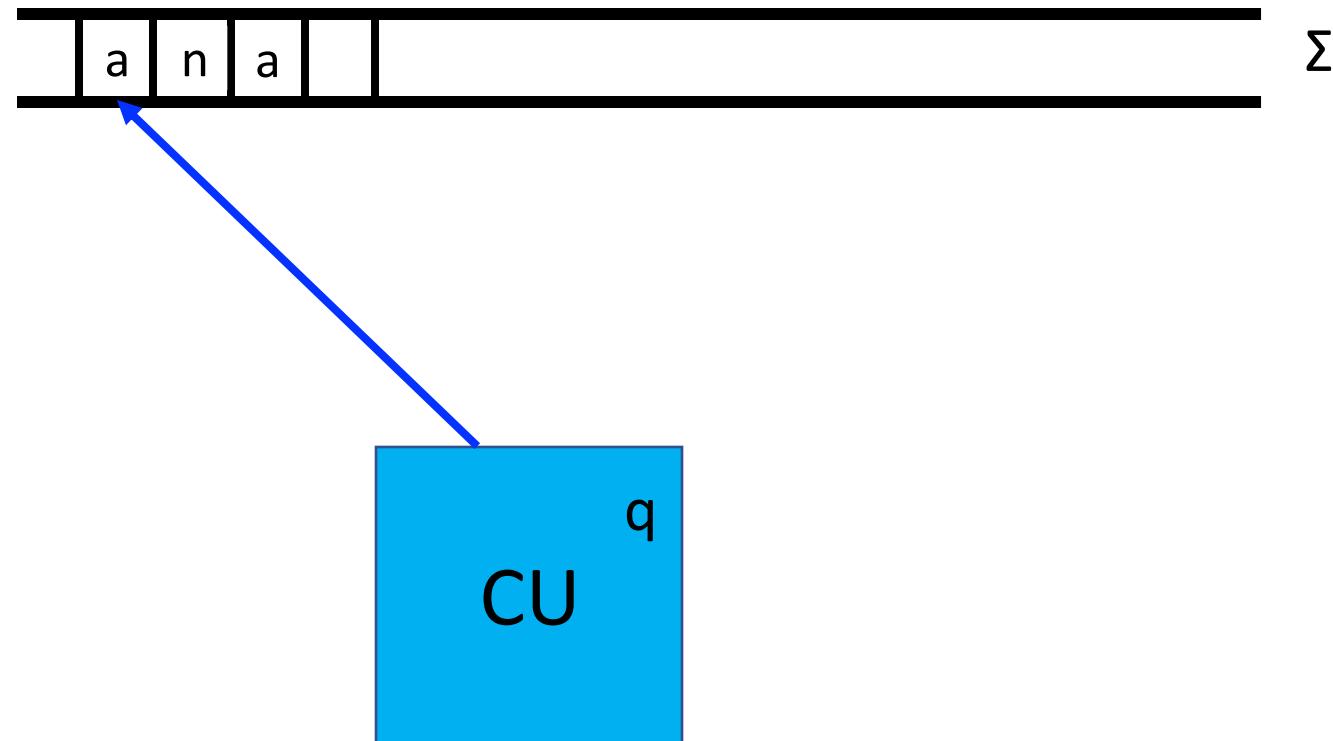## - *basic notions-*

**Problem**: The door to the tower is closed by the Red Dragon, using a complicated machinery. Prince Charming has managed to steal the plans and is asking for your help. Can you help him determining all the person names that can unlock the door

# Finite Automata (sing. Finite automaton; abrev. FA; transl = automat finit)

- Intuitive model

***Definition***: A ***finite automaton (FA)*** is a 5-tuple

$$M = (Q, \Sigma, \delta, q0, F)$$

where:

- Q - finite set of states $(|Q| < \infty)$
- $\Sigma$ - finite alphabet $(|\Sigma| < \infty)$
- $\delta$ – transition function : $\delta : Q \times \Sigma \rightarrow P(Q)$
- $q_0$ – initial state $q_0 \in Q$

- $F \subseteq Q$ – set of final states

# *Remarks*

1. $Q \cap \Sigma = \emptyset$

2. $\delta: Q \times \Sigma \rightarrow P(Q)$, $\varepsilon \in \Sigma^0$ - relation $\delta(q, \varepsilon) = p$ **NOT** allowed

3. If $|\delta(q,a)| \leq 1$ => deterministic finite automaton (DFA)

4. If $|\delta(q,a)| > 1$ (more than a state obtained as result) => nondeterministic finite automaton (NFA)

*Property*: For any NFA *M* there exists a DFA *M'* equivalent to *M*

## Configuration $C=(q,x)$

where:

- q state

- x unread sequence from input: $x \in \sum^*$

Initial configuration : $(q_0,w)$ , w - whole sequence

Final configuration: $(q_f ,\varepsilon)$ , $q_f \in F$, $\varepsilon$ –empty sequence

       (corresponds to accept)

# Relations between configurations

- ⊢ **move** / **transition** (simple, one step)
  $(q,ax) \vdash (p,x)$ , $p \in \delta(q,a)$


- $\overset{k}{\vdash}$ **k move** = a sequence of k simple transitions) $C_0 \vdash C_1 \vdash \ldots \vdash C_k$

- $\overset{+}{\vdash}$ **+ move**
  $C \overset{+}{\vdash} C'$ : $\exists$ k>0  such that        $C \overset{k}{\vdash} C'$

- $\overset{*}{\vdash}$ **\* move (star move)**
  $C \overset{*}{\vdash} C'$ : $\exists$ k≥0 such that        $C \overset{k}{\vdash} C'$

**Definition** : **Language** accepted by FA M = (Q,Σ,δ,q0,F) is:

$$L(M)=\{ w \in \Sigma^* \mid (q_0,w) \vdash^* (q_f ,\varepsilon) , q_f \in F \}$$

**Remarks**

1. 2 finite automata $M_1$ and $M_2$ are equivalent if and only if they accept the same language

$$L(M_1)=L(M_2)$$

1. $\varepsilon \in L(M) \Leftrightarrow q_0 \in F$ (initial state is final state)

# Representing FA



1. List of all elements
2. Table
3. Graphical representation

M=(Q,Σ,δ,p,F)
Q = {p,q,r}
Σ = {a,b}
δ(p,a) = q
δ(q,a)=q
δ(q,b)=r
δ(p,b)=r
F = {r}

M=(Q,Σ,δ,p,F)
F = {r}

|   | a | b |
|---|---|---|
| p | q | r |
| q | q | r |
| r | - | - |

# Example

M=(Q,Σ,δ,p,F)
Q = {p,q,r,s}
Σ = {0,1}
$\delta(p,1) = q$
$\delta(q,0)=q$
$\delta(q,1)=r$
$\delta(p,0)=s$
$\delta(s,1)=s$
$\delta(s,0)=r$
F = {r}



s

1      q

p

r

(p,101)|-(q,01)|-(q,1)|-(r, ε) accepted
(p,110) |-(q,10)|-(r,0) –not accepted

F = {p,r}

# Regular languages

# Why?

1. Search engine – succes of Google
2. Unix commands
3. Programming languages – new feature
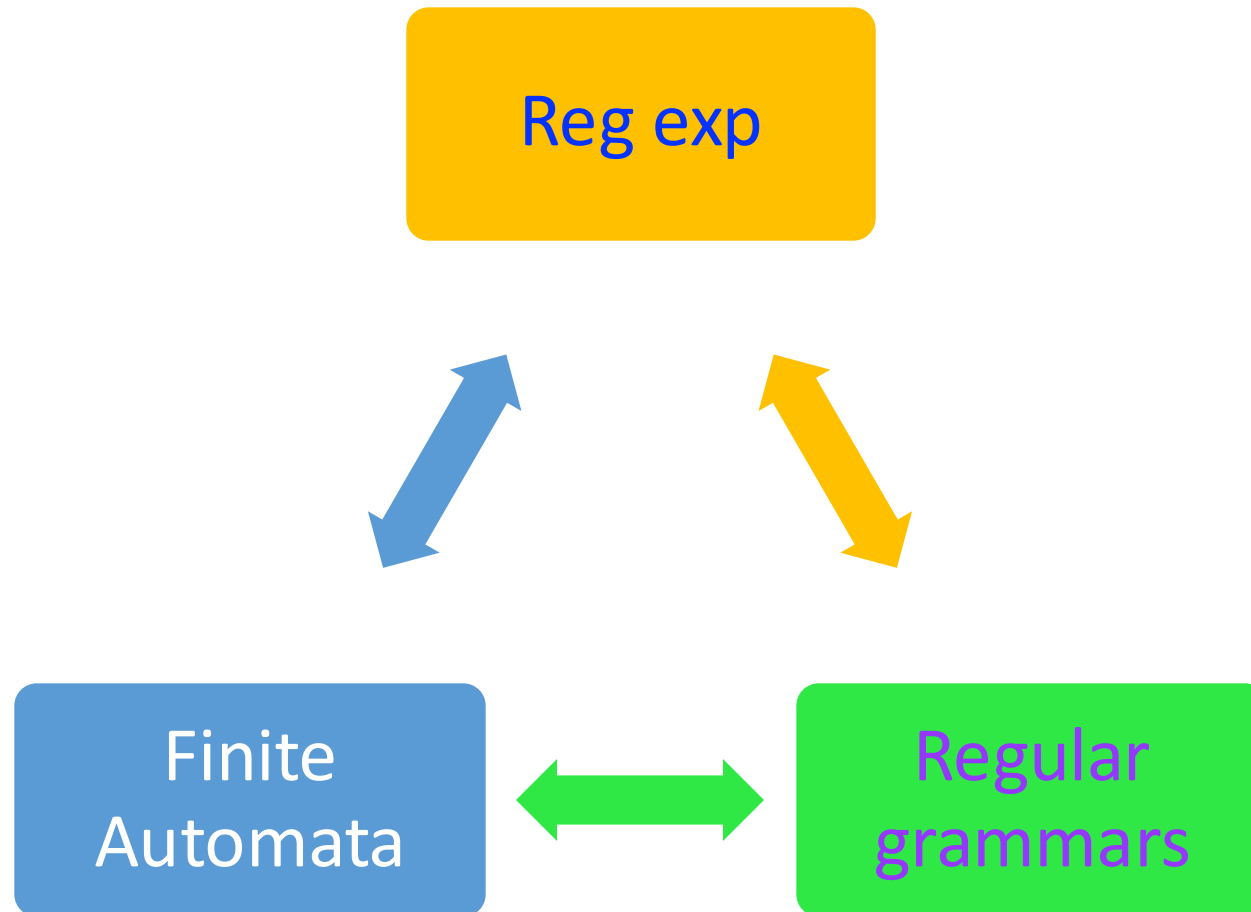
# Remember

- **Grammar**

$$G=(N,\Sigma,P,S)$$

$$L(G)=\{w\in\Sigma^* \mid S \overset{*}{\Rightarrow} w\}$$

- Finite automaton

$$M = (Q,\Sigma,\delta,q_0,F)$$

$$L(M)=\{ w \in \Sigma^* \mid (q_0,w) \vdash (q_f ,\varepsilon) , q_f \in F \}$$

# Regular grammars

- G = (N, $\Sigma$, P, S) <span style="color:red">right linear grammar</span> if

    $\forall p \in P$: A$\rightarrow$aB or A $\rightarrow$b, where A,B $\in$N and a,b $\in \Sigma$

<span style="color:white; background-color:#4472C4">A-&gt;aA|a ok   ✔<br>S-&gt;aA| ε and A-&gt;b ok ✔<br>S-&gt;aA| ε and A-&gt; ε NOT ok ✘<br>S-&gt;aA| ε and A-&gt;bS|a NOT ok✘</span>

- G = (N, $\Sigma$, P, S) <span style="color:red">regular grammar</span> if
    - G is right linear grammar

    and

    - A$\rightarrow\varepsilon \notin$ P, with the exception that S $\rightarrow\varepsilon \in$ P, in which case S does not appear in the rhs (right hand side) of any other production

- L(G) = {w $\in \Sigma$* | S $\overset{*}{=>}$ w}  - right linear language

**_Theorem 1_**: For any regular grammar G=(N, **Σ**, P, S) there exists a FA M=(Q, **Σ**, $\delta$, $q_0$,F) such that L(G) = L(M)

Proof: construct M based on G

Q = N U {K}, K $\notin$ N

$q_0$ = S

F = {K} U {S| if S→$\boldsymbol{\varepsilon}$ $\in$ P}

$\delta$: if A →aB $\in$ P then $\delta$(A,a) = B

if A →a $\in$ P then $\delta$(A,a) = K

*Theorem 1*: For any regular grammar G=(N, **Σ**, P, S) there exists a FA M=(Q, **Σ**, $\delta$, $q_0$,F) such that L(G) = L(M)

Proof: construct M based on G

Q = N U {K}, K $\notin$ N

$q_0$ = S

F = {K} U {S| if S→**ε** ∈ P}

$\delta$: if A →aB ∈ P then $\delta$(A,a) = B

if A →a ∈ P then $\delta$(A,a) = K

*Theorem 1*: For any regular grammar G=(N, **Σ**, P, S) there exists a FA M=(Q, **Σ**, $\delta$, $q_0$,F) such that L(G) = L(M)

Proof: construct M based on G

Q = N ∪ {K}, K ∉ N

$q_0$ = S

F = {K} ∪ {S| if S→$\varepsilon$ ∈ P}

$\delta$: if A →aB ∈ P then $\delta$(A,a) = B

if A →a ∈ P then $\delta$(A,a) = K

*Prove that L(G) = L(M)    (w∈L(G) ⟺ w∈L(M)):*

$S \overset{*}{\Rightarrow} w \Leftrightarrow (S,w) \overset{*}{\vdash} (qf, \varepsilon)$

w= $\varepsilon$: $S \overset{*}{\Rightarrow} \varepsilon \Leftrightarrow (S, \varepsilon) \overset{*}{\vdash} (S, \varepsilon)$ – true

w=$a_1a_2...a_n$: $S \overset{*}{\Rightarrow} w \Leftrightarrow (S,w) \overset{*}{\vdash} (K, \varepsilon)$

$S \Rightarrow a_1A_1 \Rightarrow a_1a_2A_2 \Rightarrow ... \Rightarrow a_1a_2...a_{n-1}A_{n-1} \Rightarrow a_1a_2...a_{n-1}a_n$

$S \Rightarrow a_1A_1$ exists if $S \rightarrow a_1A_1$ and then $\delta(S,a_1)=A_1$

$A_1 \rightarrow a_2A_2 : \delta(A_1,a_2)=A_2 ...$

$A_{n-1} \rightarrow a_n : \delta(A_{n-1},a_n)=K$

$(S,a_1a_2...a_n) \vdash (A_1,a_2...a_n) \vdash (A_2,a_3...a_n) \vdash ... \vdash (A_{n-1},a_n) \vdash (K, \varepsilon)$ , K∈F

# EX 1

G = ({S,A}, {0,1}, P, S)

P:      S-> 0S|0A

        A -> 1A|1

M:

Q = {S,A,K}

$q_0$ =S

F ={K}

$\delta$

|  | 0 | 1 |
|---|---|---|
| S | S, A |  |
| A |  | A,K |
| K |  |  |