

# Teoria Compilării

Drăgulici Dumitru Daniel

Facultatea de matematică și informatică,  
Universitatea București

2011

# Bibliografie



Aho A.V., Ullman J.D.:  
The Theory of Parsing, Translation and Compiling,  
Prentice Hall Inc., vol. I (1973), vol. II (1974)



Aho A.V., Lam M.S., Sethi R., Ullman J.D.:  
Compilers. Principles, Techniques, & Tools  
Second Edition,  
Pearson Education Inc., 2006



Adrian Atanasiu:  
Bazele matematice ale scrierii compilatoarelor,  
Ed. Olimp, 1996



Adrian Atanasiu:  
Bazele informaticii, suport de curs pentru anul II seral,  
Tipogr. Univ. din București, 1987



Internet:  
[www.wikipedia.org](http://www.wikipedia.org)  
[www.scribd.com](http://www.scribd.com)

# Cuprins

## 1 Etapele compilării

## 2 Analiza lexicală

# Etapele compilării

Cod sursă (șir de caractere)



Analiza lexicală



Șir de tokeni



Analiza sintactică



Arbore de derivare



Analiza semantică



Arbore cu decorații



Generare cod



Cod obiect

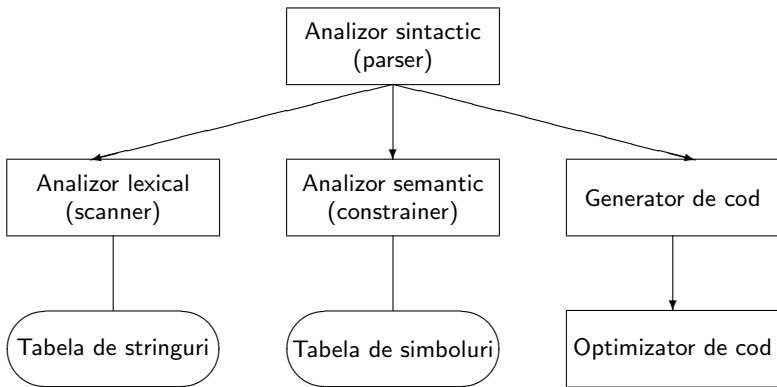


Optimizare cod



Cod obiect optimizat

# Structura unui compilator



# Cuprins

1 Etapele compilării

2 Analiza lexicală

# Analiza lexicală

**Analiza lexicală** (scanarea) are ca scop identificarea **tokenilor** (**unități lexicale**, **atomi lexicali**) succesivi din care este alcătuit codul sursă analizat.

Fiecare token are:

- un **tip** (**clasă**); exemple de tipuri de token în limbajele de programare uzuale:

*identificator*: `ionescu`, `Nelu`, `_abc23x`;

*literal întreg*: `123`, `0215`, `0x32f1`;

*literal flotant*: `12.34`, `12.34e-5`, `12e5`;

*literal caracter*: `'a'`, `'A'`, `'\n'`, `'\0'`;

*literal string*: `"abc De12"`;

*operator*: `+=`, `+`;

*separator*: `;`, `{`, `}`;

*comentariu*: `/*abc def*/`, `//abc def`;

*spațiu*;

(exemplele sunt din limbajele C/C++).

- o **valoare**; valoarea este secvența de caractere care formează tokenul.

De exemplu tokenul `_abc23x` are tipul *identificator* și valoarea șirul de caractere `_abc23x`.

# Analiza lexicală

Analiza lexicală se realizează folosind un algoritm liniar, bazat pe instrumente de limbaje regulate.

O modalitate ușoară de lucru este să definim tipurile de token folosind expresii regulate și să recunoaștem tokenii folosind automate finite deterministe.



# Expresii regulate

## Definiție: Expresii regulate (ER):

- Fie  $V, V' = \{+, *, \cdot, (, )\}$  alfabete disjuncte.

Mulțimea  $E(V)$  a ER peste  $V$  se definește recursiv astfel:

- 1)  $V \cup \{\lambda\} \cup \{\emptyset\} \subseteq E(V)$ ;
- 2) dacă  $\alpha, \beta \in E(V)$  atunci  $\alpha + \beta, \alpha \cdot \beta, \alpha^*, (\alpha) \in E(V)$ ;
- 3)  $E(V)$  nu conține alte elemente.

Vom considera următoarea ordine a priorităților:  $+ < \cdot < ^*$  și vom omite uneori parantezele, dacă nu sunt necesare; de asemenea, uneori în loc de  $\cdot$  vom folosi juxtapunerea.

- Fiecărei ER  $e \in E(V)$  îi corespunde limbajul  $\bar{e} \subseteq V^*$  definit astfel:

- 1)  $\overline{\emptyset} = \emptyset$ ;
- 2) dacă  $a \in V \cup \{\lambda\}$  atunci  $\bar{a} = \{a\}$ ;
- 3) dacă  $\alpha, \beta \in E(V)$  atunci  $\overline{\alpha + \beta} = \bar{\alpha} \cup \bar{\beta}$ ,  $\overline{\alpha \beta} = \bar{\alpha} \bar{\beta}$ ,  $\overline{\alpha^*} = \bar{\alpha}^*$ ,  
 $\overline{(\alpha)} = \bar{\alpha}$ .

# Expresii regulate

## Exemplu:

Fie  $V = \{\underline{A}, \dots, \underline{Z}, \underline{a}, \dots, \underline{z}, \underline{0}, \dots, \underline{9}, \underline{+}, \underline{*}, \underline{/}, \underline{(}, \underline{)}, \underline{=}, \underline{;}, \underline{ }\}$ .

Considerăm următoarele ER:

$idf = lit ( lit + cif ) ^*$

$int = cif cif ^*$

$com = \underline{/} \underline{*} ( notstar + \underline{*} notslash ) ^* \underline{*} \underline{/}$

$spa = \underline{ } \underline{ } ^*$

$op = \underline{+} + \underline{*} + \underline{/} + \underline{(} + \underline{)} + \underline{=}$

$delim = \underline{;}$

$lit = (\underline{A} + \dots + \underline{Z} + \underline{a} + \dots + \underline{z})$

$cif = (\underline{0} + \dots + \underline{9})$

$notstar = (lit + cif + \underline{+} + \underline{/} + \underline{(} + \underline{)} + \underline{=} + \underline{;} + \underline{ })$

$notslash = (lit + cif + \underline{+} + \underline{*} + \underline{(} + \underline{)} + \underline{=} + \underline{;} + \underline{ })$

Presupunem că doar primele 6 ER definesc clase (tipuri) de tokeni (puteam să nu mai dăm un nume ultimelor ER și să inserăm direct formula lor).

# Expresii regulate

## Exemplu:

Notăm clasele de token cu respectiv: 1 (*idf*), 2 (*int*), 3 (*com*), 4 (*spa*), 5 (*op*), 6 (*delim*).

# Expresii regulate

## Exemplu:

Notăm clasele de token cu respectiv: 1 (*idf*), 2 (*int*), 3 (*com*), 4 (*spa*), 5 (*op*), 6 (*delim*).

Fie codul sursă (cuvânt  $w \in V^*$ ):

abc=(12+x30)+1;\_y\_=abc;

# Expresii regulate

## Exemplu:

Notăm clasele de token cu respectiv: 1 (*idf*), 2 (*int*), 3 (*com*), 4 (*spa*), 5 (*op*), 6 (*delim*).

Fie codul sursă (cuvânt  $w \in V^*$ ):

abc=(12+x30)+1;\_y\_=abc;

Atunci, în urma analizei lexicale va rezulta următoarea **interpretare** a lui  $w$ :

$\langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{\_}, 4 \rangle$   
 $\langle \underline{\text{y}}, 1 \rangle \langle \underline{\_}, 4 \rangle \langle \underline{=}, 5 \rangle \langle \underline{\text{abc}}, 1 \rangle \langle \underline{;}, 6 \rangle$

# Expresii regulate

## Exemplu:

Observații:

- În loc de  $\langle \underline{abc}, 1 \rangle$  *puteam găsi*:  $\langle \underline{a}, 1 \rangle \langle \underline{b}, 1 \rangle \langle \underline{c}, 1 \rangle$ ,  
sau:  $\langle \underline{ab}, 1 \rangle \langle \underline{c}, 1 \rangle$ , etc.;
- în general dorim să identificăm cei mai lungi tokeni care se pot forma de la poziția curentă încolo (o interpretare cu această proprietate s.n. **orientată către dreapta**);
- Pentru identificarea tokenilor ca mai sus se poate folosi un AFD.

# Automate finite deterministe

## Definiție: Automate finite deterministe (AFD):

- Un AFD (parțial) este un sistem  $A = \langle Q, V, \delta, q_0, F \rangle$ , unde:
  - $Q$  este o mulțime finită, nevidă (mulțimea stărilor);
  - $V$  este o mulțime finită, nevidă (alfabetul de intrare);
  - $\delta : Q \times V \rightarrow Q$  este o funcție parțială (funcția de tranziție);
  - $q_0 \in Q$  (starea inițială);
  - $F \subseteq Q$  (mulțimea stărilor finale).
- Extindem  $\delta$  la o funcție parțială  $\hat{\delta} : Q \times V^* \rightarrow Q$  astfel:
  - $\hat{\delta}(q, \lambda) = q$  pentru orice  $q \in Q$ ,
  - $\hat{\delta}(q, \alpha a) = \delta(\hat{\delta}(q, \alpha), a)$ , dacă există  $\hat{\delta}(q, \alpha)$  și  $\delta(\hat{\delta}(q, \alpha), a)$ .
- Notăm  $\hat{\delta}$  tot cu  $\delta$ . Va rezulta:
  - $\delta(q, \alpha\beta) = \delta(\delta(q, \alpha), \beta)$  pentru orice  $q \in Q$  și  $\alpha, \beta \in V^*$ .
- Limbajul recunoscut de  $A$  este următoarea submulțime a lui  $V^*$ :
$$L(A) = \{\alpha \in V^* : \delta(q_0, \alpha) \in F\}.$$
- Un AFD  $A$  este total, dacă  $\delta$  este total definită (ca funcție  $\delta : Q \times V \rightarrow Q$ ); va rezulta că și  $\hat{\delta}$  este total definită (ca funcție  $\hat{\delta} : Q \times V^* \rightarrow Q$ ).

# Automate finite deterministe

## Observație:

*Pentru orice AFD  $A$  există un AFD total  $A'$  echivalent cu el, adică a.î.  
 $L(A) = L(A')$ .*

*El se poate construi din  $A$  adăugând o nouă stare nefinală și ducând în ea toate tranzițiile nedefinite, inclusiv tranziții de la ea la ea însăși, cu toate simbolurile alfabetului.*



# Automate finite deterministe

Un AFD se poate reprezenta printr-o diagramă de tip graf orientat, unde:

- vârfurile reprezintă stări; se desenează sub forma unor discuri ce conțin o etichetă de identificare a vârfului; starea inițială are o săgeată de intrare, stările finale au conturul dublat (sau o săgeată de ieșire);
- arcele reprezintă tranziții; se desenează prin săgeți ce unesc stările sursă și destinație corespunzătoare, în dreptul cărora se scrie simbolul cu care se face tranziția; pentru simplitate, dacă mai multe tranziții au aceeași sursă și destinație, ele se pot desena printr-o singură săgeată, în dreptul căreia se scriu toate simbolurile asociate lor (sau o mulțime, ER, etc.).

# Automate finite deterministe

## Propoziție:

Pentru orice limbaj  $L \subseteq V^*$  sunt echivalente:

- 1) Există un AFD  $A$  cu alfabetul  $V$  a.î.  $L = L(A)$ ;
- 2) Există o ER  $e \in E(V)$  a.î.  $L = \bar{e}$ .

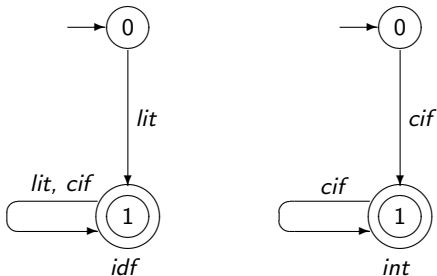
Un limbaj ce satisface oricare din proprietățile echivalente de mai sus s.n. **limbaj regulat**.

Există și alte proprietăți echivalente cu cele două, anume:

- $L$  este recunoscut de un automat finit nedeterminist (AFN);
  - $L$  este recunoscut de un automat finit nedeterminist cu  $\lambda$ -tranziții (AFNL);
  - $L$  este generat de o gramatică regulată (GR);
- (a se vedea cursul de limbaje formale, anul I).

# Automate finite deterministe

De exemplu, AFD echivalente cu primele două ER din exemplul anterior sunt:



(prin "lit", "cif" și "lit, cif" am notat acum totalitatea simbolurilor din  $V$  descrise de respectiv aceste ER).

Riguros, trecerea de la ER la AFD echivalent se face folosind un algoritm ...

# Algoritm ER $\longrightarrow$ AFD

## Algoritm (ER $\longrightarrow$ AFD):

Intrare:  $e \in E(V)$ .

Ieșire:  $A = \langle Q, V, \delta, q_0, F \rangle$  AFD a.î.  $L(A) = \bar{e}$ .

Notăție:

Considerăm un simbol nou  $\# \notin V$ .

Numim **extensia** lui  $e$  expresia regulată  $(e)\# \in E(V \cup \{\#\})$ .

# Algoritm ER $\longrightarrow$ AFD

## Pasul 1:

Se construiește arborele asociat ER extinse; în acest arbore nu vom mai considera și nodurile asociate parantezelor (nu mai sunt necesare).

Nodurile neterminale ale arborelui vor corespunde unor operatori  $+$ ,  $\cdot$  sau  $\star$ ; cele terminale vor corespunde unor simboluri din  $V \cup \{\lambda\} \cup \{\#\}$ .

Asociem fiecărui nod (neterminal sau terminal) un identificator distinct, de ex.  $n1$ ,  $n2$ , etc.

Asociem frunzelor, de la stânga la dreapta, numere întregi consecutive, numărând de la 1; aceste numere se vor numi **poziții**.

Obs: frunze diferite vor avea asociate numere (poziții) diferite, chiar dacă simbolul din alfabet corespunzător este aceeași.

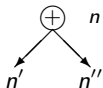
# Algorithm ER $\rightarrow$ AFD

Deci nodurile arborelui vor fi de forma:

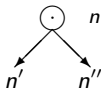
$$a \in V \cup \{\lambda\}$$



$$e_1 + e_2$$



$$e_1 \cdot e_2$$



$$e^*$$



( $n$ ,  $n'$ ,  $n''$  sunt identificatori,  $i$  sunt poziții (doar în cazul frunzelor)).

# Algoritm ER $\rightarrow$ AFD

## Pasul 2:

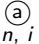
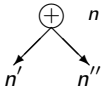
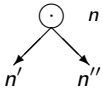
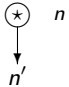
Se calculează mulțimile:

$firstpos(n)$ ,  $lastpos(n)$ ,  $nullable(n)$ , pentru fiecare nod  $n$ ;

$followpos(i)$ , pentru fiecare poziție  $i$ .

## Algoritm ER $\rightarrow$ AFD

Mai întâi construim mulțimile  $firstpos(n)$ ,  $lastpos(n)$ ,  $nullable(n)$ , astfel:

tip nod $n$	$firstpos(n)$	$lastpos(n)$	$nullable(n)$
 $a \in V \cup \{\lambda\}$ $n, i$	$\{i\}$	$\{i\}$	true d.d $a = \lambda$
	$firstpos(n')$ $\cup firstpos(n'')$	$lastpos(n')$ $\cup lastpos(n'')$	$nullable(n')$ or $nullable(n'')$
	if $nullable(n')$ then $firstpos(n')$ $\cup firstpos(n'')$ else $firstpos(n')$	if $nullable(n'')$ then $lastpos(n')$ $\cup lastpos(n'')$ else $firstpos(n'')$	$nullable(n')$ and $nullable(n'')$
	$firstpos(n')$	$lastpos(n')$	true

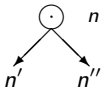


## Algoritm ER $\rightarrow$ AFD

Apoi construim mulțimile  $followpos(i)$ , astfel:

- pentru orice poziție  $i$ , inițializez  $followpos(i) = \emptyset$ ;

- pentru orice nod  
la  $followpos(i)$ ;



și orice  $i \in lastpos(n')$ , reunesc  $firstpos(n'')$

- pentru orice nod  
la  $followpos(i)$ .



și orice  $i \in lastpos(n')$ , reunesc  $firstpos(n')$

# Algoritm ER $\rightarrow$ AFD

## Observație:

- $firstpos(n)/lastpos(n)$  este mulțimea pozițiilor pe care le poate avea primul simbol/ultimul simbol al unui șir generat de subarborele cu rădăcina  $n$ ;  
     $nullable(n)=true$  d.d. subarborele cu rădăcina  $n$  poate genera șirul vid.  
     $followpos(i)$  este mulțimea pozițiilor pe care le pot avea simbolurile ce urmează simbolului de pe poziția  $i$  într-un șir generat de arborele ER;
- toate mulțimile  $firstpos$ ,  $lastpos$ ,  $nullable$ ,  $followpos$  se pot calcula printr-o singură parcurgere în adâncime a arborelui (la pașii de revenire).

## Algoritm ER $\rightarrow$ AFD

### Pasul 3:

Se construiește  $A = \langle Q, V, \delta, q_0, F \rangle$  AFD a.î.  $L(A) = \bar{e}$ , astfel:

- elementele lui  $Q$  sunt mulțimi de poziții;
- $q_0 = \text{firpos}(\text{rădăcină})$ ;
- $Q, F$  și  $\delta$  se construiesc astfel:

```
Q ← {firstpos(rădăcină)}, firstpos(rădăcină) stare nemarcată;  
if [firstpos(rădăcină) conține poziția lui #] then F ← {firstpos(rădăcină)}  
    else F ← ∅;  
while [există T ∈ Q nemarcată] do begin  
    marchează T;  
    for [fiecare a ∈ V] do begin  
        U ← {j ∈ followpos(i) : i ∈ T, i etichetat cu a};  
        if U ∉ Q then begin  
            Q ← Q ∪ {U}; U stare nemarcată;  
            if [U conține poziția lui #] then F ← F ∪ {U}  
        end;  
        δ(T, a) ← U  
    end  
end
```

# Algoritm ER $\longrightarrow$ AFD

## Exemplu:

Să aplicăm algoritmul pentru ER *idf* din exemplul anterior.

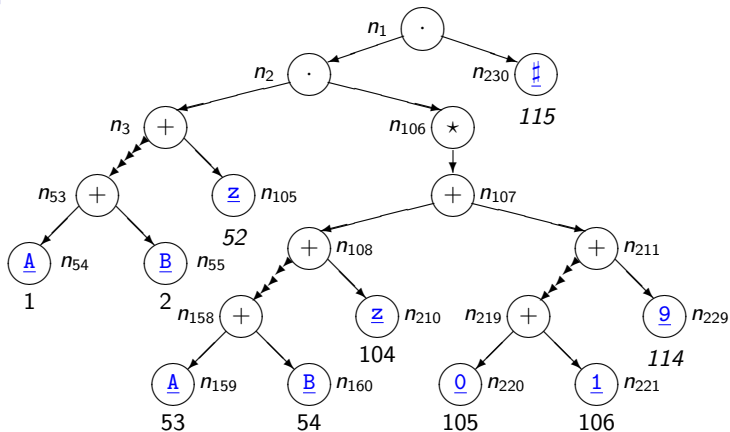
Expandând ER îmbricate și considerând ER extinsă, aceasta se scrie:

$$idf = ((\underline{A} + \dots + \underline{Z} + \underline{a} + \dots + \underline{z}) ((\underline{A} + \dots + \underline{Z} + \underline{a} + \dots + \underline{z}) + (\underline{0} + \dots + \underline{9}))^*) \#$$

Construim arborele asociat ER extinse, numim nodurile și numerotăm pozițiile:

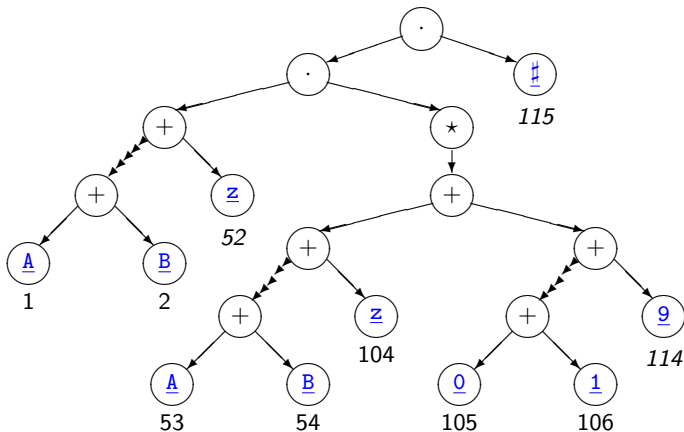
# Algorithm ER $\rightarrow$ AFD

Exemplu:



## Algoritmo ER $\rightarrow$ AFD

Exemplu:

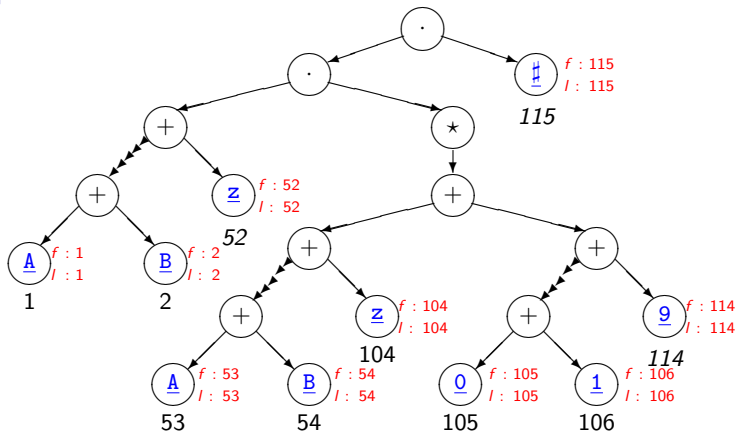


Calculăm mulțimile *firstpos* (*f*), *lastpos* (*l*), *nullable* (vom nota doar nodurile pentru care *nullable* = *true*, scriind *n* : *t*), *followpos*, făcând o parcurgere în adâncime a arborelui (când lucrăm direct pe desen nu mai sunt necesari identificatorii vârfurilor și i-am omis).

Vom scrie cu roșu valorile nou adăugate la fiecare pas.

# Algorithm ER $\rightarrow$ AFD

Exemplu:



*followpos*(1):

...

*followpos*(52):

*followpos*(53):

...

*followpos*(104):

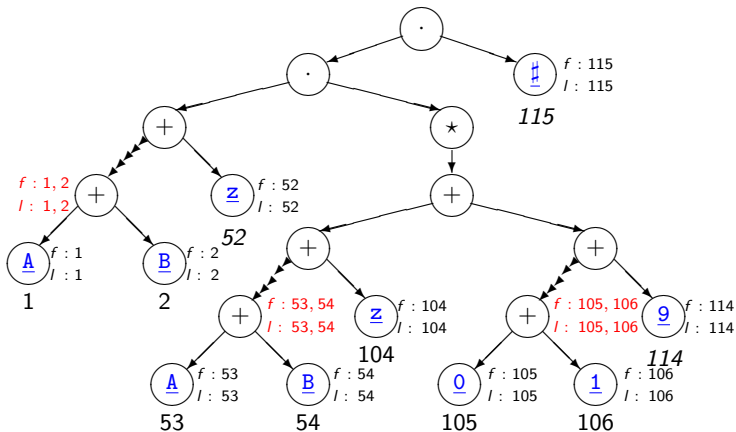
*followpos*(105):

...

*followpos*(114):

# Algorithm ER $\rightarrow$ AFD

Exemplu:



*followpos*(1):

...

*followpos*(52):

*followpos*(53):

...

*followpos*(104):

*followpos*(105):

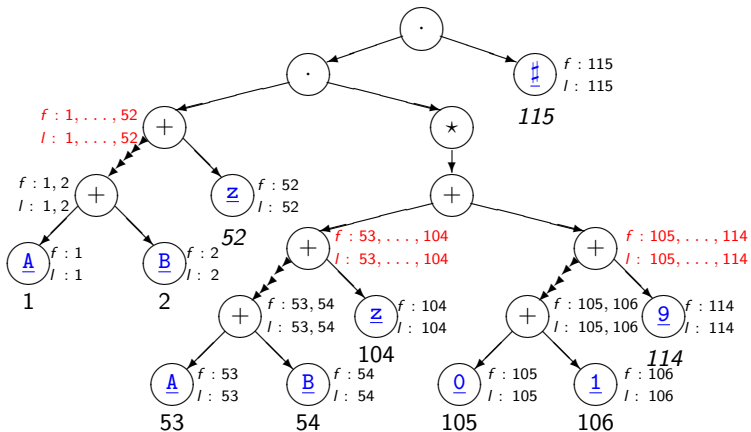
...

*followpos*(114):



# Algorithm ER $\rightarrow$ AFD

Exemplu:



$followpos(1):$

...

$followpos(52):$

$followpos(53):$

...

$followpos(104):$

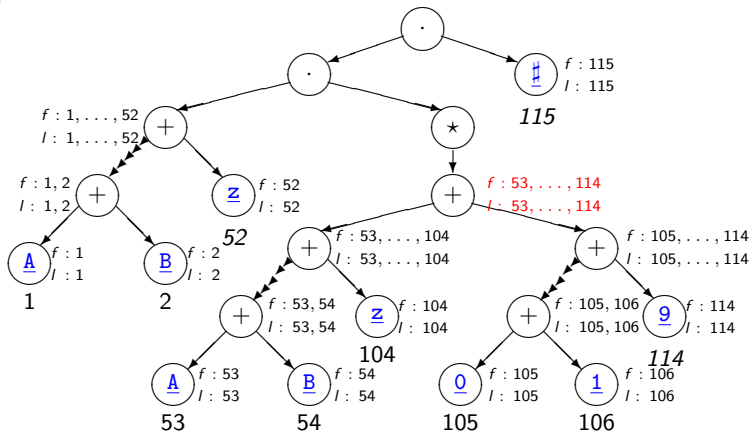
$followpos(105):$

...

$followpos(114):$

# Algorithm ER $\rightarrow$ AFD

Exemplu:



$followpos(1):$

...

$followpos(52):$

$followpos(53):$

...

$followpos(104):$

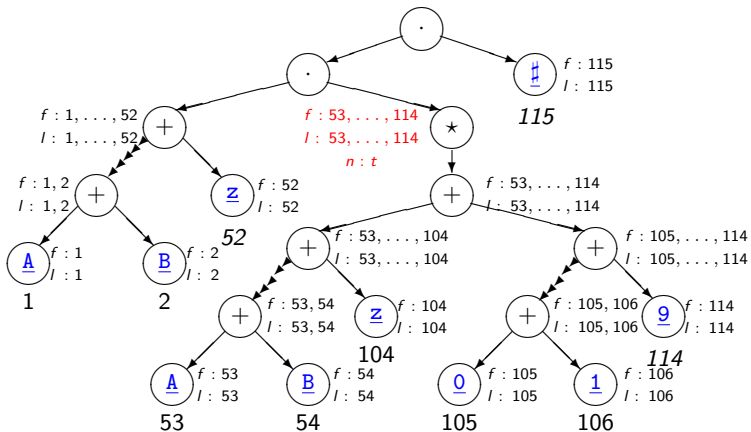
$followpos(105):$

...

$followpos(114):$

# Algorithm ER $\rightarrow$ AFD

Exemplu:



$followpos(1):$

...

$followpos(52):$

$followpos(53):$  53, ..., 114

...

$followpos(104):$  53, ..., 114

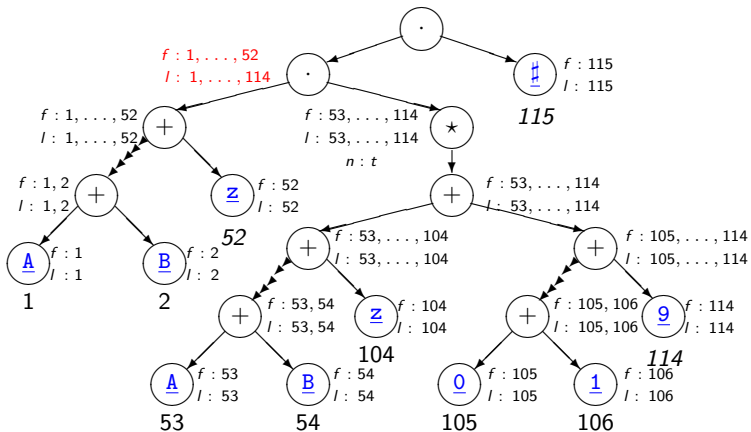
$followpos(105):$  53, ..., 114

...

$followpos(114):$  53, ..., 114

# Algorithm ER $\rightarrow$ AFD

Exemplu:



$followpos(1): 53, \dots, 114$      $followpos(53): 53, \dots, 114$      $followpos(105): 53, \dots, 114$

...

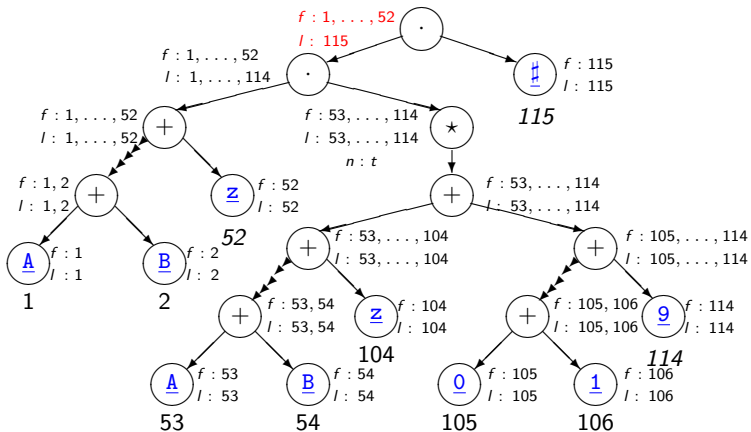
...

...

$followpos(52): 53, \dots, 114$      $followpos(104): 53, \dots, 114$      $followpos(114): 53, \dots, 114$

## Algoritmo ER $\rightarrow$ AFD

Exemplu:



*followpos*(1): 53, ..., 114, 115    *followpos*(53): 53, ..., 114, 115    *followpos*(105): 53, ..., 114, 115

• • •

• • •

• • •

*followpos*(52): 53, ..., 114, 115    *followpos*(104): 53, ..., 114, 115    *followpos*(114): 53, ..., 114, 115

## Algoritm ER $\longrightarrow$ AFD

### Exemplu:

Din calculele de mai sus este suficient să reținem că

$firstpos(rădăcină) = \{1, \dots, 52\}$

$followpos(i) = 53, \dots, 115$  pentru orice  $i$

poziția lui # este 115

Construim automatul, mai exact  $q_0, Q, F, \delta$ :

- avem  $q_0 = \{1, \dots, 52\}$
- inițializăm  $Q, F, \delta$  cu  $\emptyset$ , apoi avem succesiv (am scris cu roșu stările și tranzițiile noi adăugate, am subliniat stările marcate, iar la ultima stare marcată am consemnat deasupra simbolurile corespunzătoare pozițiilor):

$Q$  :

$F$  :

$\delta$  :

## Algoritm ER $\longrightarrow$ AFD

### Exemplu:

Din calculele de mai sus este suficient să reținem că

$firstpos(r\grave{a}d\acute{a}cin\grave{a}) = \{1, \dots, 52\}$

$followpos(i) = 53, \dots, 115$  pentru orice  $i$

poziția lui # este 115

Construim automatul, mai exact  $q_0$ ,  $Q$ ,  $F$ ,  $\delta$ :

- avem  $q_0 = \{1, \dots, 52\}$
- inițializăm  $Q$ ,  $F$ ,  $\delta$  cu  $\emptyset$ , apoi avem succesiv (am scris cu roșu stările și tranzițiile noi adăugate, am subliniat stările marcate, iar la ultima stare marcată am consemnat deasupra simbolurile corespunzătoare pozițiilor):

$Q : \{1, \dots, 52\}$

$F :$

$\delta :$

## Algoritm ER $\longrightarrow$ AFD

### Exemplu:

Din calculele de mai sus este suficient să reținem că

$firstpos(rădăcină) = \{1, \dots, 52\}$

$followpos(i) = 53, \dots, 115$  pentru orice  $i$

poziția lui # este 115

Construim automatul, mai exact  $q_0, Q, F, \delta$ :

- avem  $q_0 = \{1, \dots, 52\}$
- inițializăm  $Q, F, \delta$  cu  $\emptyset$ , apoi avem succesiv (am scris cu roșu stările și tranzițiile noi adăugate, am subliniat stările marcate, iar la ultima stare marcată am consemnat deasupra simbolurile corespunzătoare pozițiilor):

A ... Z  
 $Q : \{ \underline{1}, \dots, \underline{52} \}, \{ 53, \dots, 115 \}$   
 $F : \{ 53, \dots, 115 \}$   
 $\delta : \{ \underline{1}, \dots, \underline{52} \} \longrightarrow \{ 53, \dots, 115 \}$   
          A, ..., Z



## Algoritm ER $\longrightarrow$ AFD

### Exemplu:

Din calculele de mai sus este suficient să reținem că

$firstpos(\text{rădăcină}) = \{1, \dots, 52\}$

$followpos(i) = 53, \dots, 115$  pentru orice  $i$

poziția lui # este 115

Construim automatul, mai exact  $q_0$ ,  $Q$ ,  $F$ ,  $\delta$ :

- avem  $q_0 = \{1, \dots, 52\}$
- inițializăm  $Q$ ,  $F$ ,  $\delta$  cu  $\emptyset$ , apoi avem succesiv (am scris cu roșu stările și tranzițiile noi adăugate, am subliniat stările marcate, iar la ultima stare marcată am consemnat deasupra simbolurile corespunzătoare pozițiilor):

$Q : \{1, \dots, 52\}, \underline{\text{A}} \dots \underline{\text{z}} \quad \underline{\text{0}} \dots \underline{\text{9}} \quad \underline{\text{\#}}$   
 $\underline{\{53, \dots, 104, 105, \dots, 114, 115\}}$

$F : \{53, \dots, 115\}$

$\delta : \{1, \dots, 52\} \longrightarrow \{53, \dots, 115\}, \text{\textcolor{red}{\{53, \dots, 115\}}} \longrightarrow \text{\textcolor{red}{\{53, \dots, 115\}}}$   
 $\underline{\text{A}}, \dots, \underline{\text{z}} \qquad \underline{\text{A}}, \dots, \underline{\text{z}}, \underline{\text{0}}, \dots, \underline{\text{9}}$

## Algoritm ER $\longrightarrow$ AFD

### Exemplu:

Numerotând starile:  $\{1, \dots, 52\} = 0$ ,  $\{53, \dots, 115\} = 1$ , avem automatul definit prin:

$$V = \{\underline{A}, \dots, \underline{z}, \underline{0}, \dots, \underline{9}\}$$

$$Q = \{0, 1\}$$

$$q_0 = 0$$

$$F = \{1\}$$

$\delta$	<u>A</u>	...	<u>z</u>	<u>0</u>	...	<u>9</u>
0	1	...	1	1	...	1
1	1	...	1	1	...	1

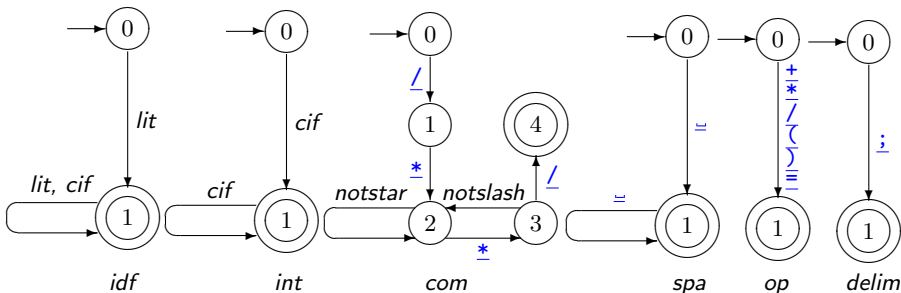
adică tocmai automatul desenat mai înainte.

Obs: pentru a ușura lucrul cu grupuri numeroase de simboluri ale alfabetului folosite la fel (și a evita notațiile cu "..."), se pot introduce **clase** de simboluri (în exemplele noastre: *lit*, *cif*, *notstar*, *notslash*).

# Analiza lexicală

## Exemplu:

Procedând în același fel, obținem următoarele AFD echivalente cu respectiv cele 6 ER ce defineau clasele de tokeni în exemplele anterioare (exercițiu):



(prin "lit", "cif" și "lit, cif" am notat și acum totalitatea simbolurilor din  $V$  descrise de respectiv aceste ER).

# Analiza lexicală

Uneori AFD obținute prin diverse procedee pot fi micșorate, în sensul reducerii numărului de stări (nu e cazul celor 6 AFD de mai înainte).

În acest scop se poate aplica un algoritm de minimizare AFD (a se vedea cursul de limbaje formale, anul I).

# Analiza lexicală

Odată construit câte un AFD pentru fiecare tip de token, putem face analiza lexicală folosind în paralel aceste AFD, sau putem folosi un AFD unic (construit din AFD respective), pentru care am asociat câte un tip de token fiecărei stări finale.

AFD unic se poate construi întotdeauna, în baza următoarelor rezultate și algoritmi:

## Propoziție:

*Clasa limbajelor regulate este închisă la reuniunea finită, mai exact:  
Dacă  $L_1, \dots, L_n \subseteq V^*$  sunt limbaje regulate ( $n \in \mathbf{N}$ ,  $n \geq 1$ ), atunci  
 $L_1 \cup \dots \cup L_n$  este limbaj regulat.*

Demonstrația se face considerând câte un instrument de limbaje regulate (AFD, AFN, AFNL, ER, GR) pentru fiecare dintre  $L_1, \dots, L_n$  și construind un alt instrument de limbaje regulate (AFD, AFN, AFNL, ER, GR) despre care se demonstrează că descrie  $L_1 \cup \dots \cup L_n$  (a se vedea cursul de limbaje formale, anul I).

În funcție de alegerile făcute, se obțin diverse demonstrații.

De exemplu din  $n$  AFD-uri se construiește un alt AFD; sau din  $n$  AFD-uri se construiește un AFNL (e mai ușor), etc.

# Analiza lexicală

Ne va fi utilă următoarea variantă a propoziției anterioare:

## Propoziție:

*Fiind date AFD totale  $A_i = \langle Q_i, V, \delta_i, q_0^i, F_i \rangle$ ,  $i = 1, \dots, n$ , a.î.*

*$L(A_i) \cap L(A_j) = \emptyset$  pentru orice  $i \neq j$ , există un AFD  $A = \langle Q, V, \delta, q_0, F \rangle$  și o funcție  $\text{tok} : F \rightarrow \{1, \dots, n\}$  a.î.:*

- $L(A) = L(A_1) \cup \dots \cup L(A_n)$ ;
- *pentru orice  $\alpha \in V^*$  și  $i \in \{1, \dots, n\}$  avem  $\alpha \in L(A_i)$  d.d.  $\delta(q_0, \alpha) \in F$  și  $\text{tok}(\delta(q_0, \alpha)) = i$ .*

Algoritmul de construcție a lui  $A$  este următorul (demonstrațiile sunt simple - exercițiu):

## Algoritm reuniune disjunctă AFD

```
Q ← {⟨q01, ..., q0n⟩}, ⟨q01, ..., q0n⟩ stare nemarcată;  
q0 ← ⟨q01, ..., q0n⟩;  
if [există i ∈ {1, ..., n} a.î. q0i ∈ Fi] then begin F ← {q0}; tok(q0) ← i end  
                                else F ← ∅;  
while [există q = ⟨q1, ..., qn⟩ ∈ Q nemarcată] do begin  
    marchează q;  
    for [fiecare a ∈ V] do begin  
        r ← ⟨δ1(q1, a), ..., δn(qn, a)⟩ not≡ ⟨r1, ..., rn⟩;  
        if r ∉ Q then begin  
            Q ← Q ∪ {r}; r stare nemarcată;  
            if [există i ∈ {1, ..., n} a.î. ri ∈ Fi] then begin  
                F ← F ∪ {r}; tok(r) ← i  
            end  
        end;  
        δ(q, a) ← r  
    end  
end  
end
```

Obs: faptul că  $L(A_i) \cap L(A_j) = \emptyset$  pentru orice  $i \neq j$  garantează că mai sus întotdeauna va exista cel mult un  $i \in \{1, \dots, n\}$  a.î.  $q_0^i \in F_i$ , respectiv  $r^i \in F_i$ , și astfel funcția *tok* este bine definită.



# Algoritm reuniune disjunctă AFD

## Exemplu:

Să aplicăm algoritmul celor 6 AFD ce defineau clase de tokeni în exemplele anterioare. Pentru aceasta le completăm mai întâi la niște AFD totale:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u><i>⌊ * V \ {⌊, *}⌋</i></u>	$\delta$   <i>= altesb</i>	$\delta$   <u><i>+ *, /, (, ), = altesb</i></u>	$\delta$   <i> ; altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
F: 1	F: 1	3   4 2 2	F: 1	F: 1	F: 1
		4   5 5 5			
		5   5 5 5			
		F: 4			

Am scris automatele mai compact: din fiecare tabel se deduce mulțimea stărilor și funcția de tranziție, starea inițială este 0, cele finale sunt scrise dedesubt, "*lit*", "*cif*", "*notstar*", "*notslash*" desemnează mulțimile respective de simboluri, iar "*altesb*" desemnează în cazul fiecărui automat mulțimea simbolurilor din  $V$  neacoperite de celelalte coloane. De asemenea, în dreptul numelui fiecărui automat am scris între paranteze numărul clasei de tokeni pe care o definește.

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u><i>/ *</i></u> $V \setminus \{/, *\}$	$\delta$   <u><i>=</i></u> <i>altesb</i>	$\delta$   <u><i>+, *, /, (, ), =</i></u> <i>altesb</i>	$\delta$   <u><i>;</i></u> <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F: 1</i>	<i>F: 1</i>	3   4 2 2	<i>F: 1</i>	<i>F: 1</i>	<i>F: 1</i>
		4   5 5 5			
		5   5 5 5			
		<i>F: 4</i>			

Avem (am numerotat stările noului AFD începând de la 0, am scris cu roșu stările și tranzițiile noi adăugate, am subliniat stările marcate și am scris între [ ] valoarea *tok* pentru fiecare stare finală):

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit</i> <i>cif</i> <i>altesb</i>	$\delta$   <i>cif</i> <i>altesb</i>	$\delta$   <u>/</u> <u>*</u> $V \setminus \{/, *\}$	$\delta$   <u>=</u> <i>altesb</i>	$\delta$   <u>+, *, /, (, ), =</u> <i>altesb</i>	$\delta$   <u>;</u> <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F</i> : 1	<i>F</i> : 1	3   4 2 2	<i>F</i> : 1	<i>F</i> : 1	<i>F</i> : 1
		4   5 5 5			
		5   5 5 5			
		<i>F</i> : 4			

$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0$

$q_0 = 0;$

$F :$

$\delta :$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   $\frac{\textcolor{blue}{/} \textcolor{blue}{*}}{V \setminus \{\textcolor{blue}{/}, \textcolor{blue}{*}\}}$	$\delta$   $\textcolor{blue}{=} \textcolor{blue}{\neq}$	$\delta$   $\textcolor{blue}{+}, \textcolor{blue}{-}, \textcolor{blue}{/}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}$	$\delta$   $\textcolor{blue}{;}$
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F</i> : 1	<i>F</i> : 1	3   4 2 2	<i>F</i> : 1	<i>F</i> : 1	<i>F</i> : 1
		4   5 5 5			
		5   5 5 5			
		<i>F</i> : 4			

$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2,$   
 $\langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5,$   
 $\langle 2, 2, 5, 2, 2, 1 \rangle = 6$

$q_0 = 0;$

*F* : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]

$\delta : 0 \xrightarrow{\textcolor{red}{lit}} 1, 0 \xrightarrow{\textcolor{red}{cif}} 2, 0 \xrightarrow{\textcolor{blue}{/}} 3, 0 \xrightarrow{\textcolor{blue}{*}} 4, 0 \xrightarrow{\textcolor{blue}{=}} 5, 0 \xrightarrow{\textcolor{blue}{+}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}} 4, 0 \xrightarrow{\textcolor{blue}{;}} 6$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u><i>/ *</i></u> <i>V \ {/,*}</i>	$\delta$   <i>= altesb</i>	$\delta$   <u><i>+,*,/,(,),=</i></u> <i>altesb</i>	$\delta$   <u><i>;</i></u> <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F: 1</i>	<i>F: 1</i>	3   4 2 2	<i>F: 1</i>	<i>F: 1</i>	<i>F: 1</i>
		4   5 5 5			
		5   5 5 5			
		<i>F: 4</i>			

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{\text{lit}} 1, 0 \xrightarrow{\text{cif}} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+,*,/,(,),=} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{\text{lit}, \text{cif}} 1, 1 \xrightarrow{V \setminus (\text{lit} \cup \text{cif})} 7$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   $\frac{\textcolor{blue}{/}}{\textcolor{blue}{*}} \textcolor{blue}{V} \setminus \{\textcolor{blue}{/}, \textcolor{blue}{*}\}$	$\delta$   $\textcolor{blue}{=} \textcolor{blue}{altesb}$	$\delta$   $\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{/}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}$ altesb	$\delta$   $\textcolor{blue}{;} \textcolor{blue}{altesb}$
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F</i> : 1	<i>F</i> : 1	3   4 2 2	<i>F</i> : 1	<i>F</i> : 1	<i>F</i> : 1
		4   5 5 5			
		5   5 5 5			
		<i>F</i> : 4			

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{\textcolor{blue}{lit}} 1, 0 \xrightarrow{\textcolor{blue}{cif}} 2, 0 \xrightarrow{\textcolor{blue}{/}} 3, 0 \xrightarrow{\textcolor{blue}{=}} 5, 0 \xrightarrow{\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}} 4, 0 \xrightarrow{\textcolor{blue}{;}} 6, \\ 1 \xrightarrow{\textcolor{blue}{lit}, \textcolor{blue}{cif}} 1, 1 \xrightarrow{\textcolor{blue}{V} \setminus (\textcolor{blue}{lit} \cup \textcolor{blue}{cif})} 7, \textcolor{red}{2} \xrightarrow{\textcolor{red}{cif}} \textcolor{red}{2}, \textcolor{red}{2} \xrightarrow{\textcolor{red}{V} \setminus \textcolor{red}{cif}} \textcolor{red}{7}$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$	$\delta$	$\delta$	$\delta$	$\delta$	$\delta$
$\begin{array}{c ccc} & lit & cif & altesb \\ 0 & 1 & 2 & 2 \\ 1 & 1 & 1 & 2 \\ 2 & 2 & 2 & 2 \\ F: 1 \end{array}$	$\begin{array}{c ccc} & cif & altesb \\ 0 & 1 & 2 \\ 1 & 1 & 2 \\ 2 & 2 & 2 \\ F: 1 \end{array}$	$\begin{array}{c ccc} & \underline{L} & \underline{*} & V \setminus \{\underline{L}, \underline{*}\} \\ 0 & 1 & 5 & 5 \\ 1 & 5 & 2 & 5 \\ 2 & 2 & 3 & 2 \\ 3 & 4 & 2 & 2 \\ 4 & 5 & 5 & 5 \\ 5 & 5 & 5 & 5 \\ F: 4 \end{array}$	$\begin{array}{c ccc} & = & altesb \\ 0 & 1 & 2 \\ 1 & 1 & 2 \\ 2 & 2 & 2 \\ F: 1 \end{array}$	$\begin{array}{c ccc} & +, *, /, (, ), = & altesb \\ 0 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 2 \\ F: 1 \end{array}$	$\begin{array}{c ccc} & ; & altesb \\ 0 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 2 \\ F: 1 \end{array}$

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{lit} 1, 0 \xrightarrow{cif} 2, 0 \xrightarrow{\underline{L}} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, /, (, ), =} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{lit, cif} 1, 1 \xrightarrow{V \setminus (lit \cup cif)} 7, 2 \xrightarrow{cif} 2, 2 \xrightarrow{V \setminus cif} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$	$\delta$	$\delta$	$\delta$	$\delta$	$\delta$
$\begin{array}{c ccc} \text{lit} & \text{cif} & \text{altesb} & \\ \hline 0 & 1 & 2 & 2 \\ 1 & 1 & 1 & 2 \\ 2 & 2 & 2 & 2 \end{array}$	$\begin{array}{c ccc} \text{cif} & \text{altesb} & & \\ \hline 0 & 1 & 2 & \\ 1 & 1 & 2 & \\ 2 & 2 & 2 & \end{array}$	$\begin{array}{c ccc} \text{com} & V \setminus \{/, *\} & & \\ \hline 0 & 1 & 5 & 5 \\ 1 & 5 & 2 & 5 \\ 2 & 2 & 3 & 2 \\ 3 & 4 & 2 & 2 \\ 4 & 5 & 5 & 5 \\ 5 & 5 & 5 & 5 \end{array}$	$\begin{array}{c ccc} \text{spa} & \text{altesb} & & \\ \hline 0 & 1 & 2 & \\ 1 & 1 & 2 & \\ 2 & 2 & 2 & \end{array}$	$\begin{array}{c ccc} \text{op} & +, *, /, (, ), = & \text{altesb} & \\ \hline 0 & 1 & 2 & \\ 1 & 2 & 2 & \\ 2 & 2 & 2 & \end{array}$	$\begin{array}{c ccc} \text{delim} & ; & \text{altesb} & \\ \hline 0 & 1 & 2 & \\ 1 & 2 & 2 & \\ 2 & 2 & 2 & \end{array}$
$F: 1$	$F: 1$	$F: 4$	$F: 1$	$F: 1$	$F: 1$

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{\text{lit}} 1, 0 \xrightarrow{\text{cif}} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, (, ), =} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{\text{lit}, \text{cif}} 1, 1 \xrightarrow{V \setminus (\text{lit} \cup \text{cif})} 7, 2 \xrightarrow{\text{cif}} 2, 2 \xrightarrow{V \setminus \text{cif}} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7, 4 \xrightarrow{V} 7$$



# Algoritm reuniune disjunctă AFD

Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u><i>/ * V \ { / , * }</i></u>	$\delta$   <u><i>= altesb</i></u>	$\delta$   <u><i>+ , * , / , ( , ) , = altesb</i></u>	$\delta$   <u><i> ; altesb</i></u>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F</i> : 1	<i>F</i> : 1	3   4 2 2	<i>F</i> : 1	<i>F</i> : 1	<i>F</i> : 1
		4   5 5 5			
		5   5 5 5			
		<i>F</i> : 4			

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{lit} 1, 0 \xrightarrow{cif} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, (, ), =} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{lit, cif} 1, 1 \xrightarrow{V \setminus (lit \cup cif)} 7, 2 \xrightarrow{cif} 2, 2 \xrightarrow{V \setminus cif} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7, 4 \xrightarrow{V} 7, \\ 5 \xrightarrow{=} 5, 5 \xrightarrow{V \setminus \{=\}} 7$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <i>/ * V \ {/, *}</i>	$\delta$   <i>= altesb</i>	$\delta$   <i>+, *, /, (, ), = altesb</i>	$\delta$   <i>; altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
F: 1	F: 1	3   4 2 2	F: 1	F: 1	F: 1
		4   5 5 5			
		5   5 5 5			
		F: 4			

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{\text{lit}} 1, 0 \xrightarrow{\text{cif}} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, /, (, ), =} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{\text{lit}, \text{cif}} 1, 1 \xrightarrow{V \setminus (\text{lit} \cup \text{cif})} 7, 2 \xrightarrow{\text{cif}} 2, 2 \xrightarrow{V \setminus \text{cif}} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7, 4 \xrightarrow{V} 7, \\ 5 \xrightarrow{=} 5, 5 \xrightarrow{V \setminus \{=\}} 7, 6 \xrightarrow{V} 7$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u>/</u> <u>*</u> $V \setminus \{/, *\}$	$\delta$   <u>=</u> <i>altesb</i>	$\delta$   <u>+, *, /, (, ), =</u> <i>altesb</i>	$\delta$   <u>;</u> <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
F: 1	F: 1	3   4 2 2	F: 1	F: 1	F: 1
		4   5 5 5			
		5   5 5 5			
		F: 4			

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\delta : 0 \xrightarrow{\text{lit}} 1, 0 \xrightarrow{\text{cif}} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, /, (, ), =} 4, 0 \xrightarrow{;} 6, \\ 1 \xrightarrow{\text{lit}, \text{cif}} 1, 1 \xrightarrow{V \setminus (\text{lit} \cup \text{cif})} 7, 2 \xrightarrow{\text{cif}} 2, 2 \xrightarrow{V \setminus \text{cif}} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7, 4 \xrightarrow{V} 7, \\ 5 \xrightarrow{=} 5, 5 \xrightarrow{V \setminus \{=\}} 7, 6 \xrightarrow{V} 7, 7 \xrightarrow{V} 7$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   <u><i>/ *</i></u> <i>V \ {/ , *}</i>	$\delta$   <u><i>=</i></u> <i>altesb</i>	$\delta$   <u><i>+, *, /, (, ), =</i></u> <i>altesb</i>	$\delta$   <u><i>;</i></u> <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
<i>F: 1</i>	<i>F: 1</i>	3   4 2 2	<i>F: 1</i>	<i>F: 1</i>	<i>F: 1</i>
		4   5 5 5			
		5   5 5 5			
		<i>F: 4</i>			

$$\begin{aligned}
 Q : \langle 0, 0, 0, 0, 0, 0 \rangle &= 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\
 \langle 2, 2, 1, 2, 1, 2 \rangle &= 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\
 \langle 2, 2, 5, 2, 2, 1 \rangle &= 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8, \\
 \langle 2, 2, 3, 2, 2, 2 \rangle &= 9,
 \end{aligned}$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6]$$

$$\begin{aligned}
 \delta : 0 &\xrightarrow{\text{lit}} 1, 0 \xrightarrow{\text{cif}} 2, 0 \xrightarrow{/} 3, 0 \xrightarrow{=} 5, 0 \xrightarrow{+, *, /, (, ), =} 4, 0 \xrightarrow{;} 6, \\
 1 &\xrightarrow{\text{lit, cif}} 1, 1 \xrightarrow{V \setminus (\text{lit} \cup \text{cif})} 7, 2 \xrightarrow{\text{cif}} 2, 2 \xrightarrow{V \setminus \text{cif}} 7, 3 \xrightarrow{*} 8, 3 \xrightarrow{V \setminus \{*\}} 7, 4 \xrightarrow{V} 7, \\
 5 &\xrightarrow{=} 5, 5 \xrightarrow{V \setminus \{=\}} 7, 6 \xrightarrow{V} 7, 7 \xrightarrow{V} 7, 8 \xrightarrow{*} 9, 8 \xrightarrow{V \setminus \{*\}} 8
 \end{aligned}$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   $\begin{array}{c} \underline{\textcolor{blue}{/}} \textcolor{blue}{*} \quad V \setminus \{\underline{\textcolor{blue}{/}}, \textcolor{blue}{*}\} \\ 0 \mid 1 \ 5 \quad 5 \\ 1 \mid 5 \ 2 \quad 5 \\ 2 \mid 2 \ 3 \quad 2 \\ 3 \mid 4 \ 2 \quad 2 \\ 4 \mid 5 \ 5 \quad 5 \\ 5 \mid 5 \ 5 \quad 5 \end{array}$	$\delta$   $\textcolor{blue}{=} \text{ altesb}$	$\delta$   $\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=} \text{ altesb}$	$\delta$   $\textcolor{blue}{;} \text{ altesb}$
0   1 2 2 1   1 1 2 2   2 2 2 F: 1	0   1 2 1   1 2 2   2 2 F: 1	0   1 5 5 1   5 2 5 2   2 3 2 3   4 2 2 4   5 5 5 5   5 5 5 F: 4	0   1 2 1   1 2 2   2 2 F: 1	0   1 2 1   2 2 2   2 2 F: 1	0   1 2 1   2 2 2   2 2 F: 1

$$Q : \langle 0, 0, 0, 0, 0, 0 \rangle = 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\ \langle 2, 2, 1, 2, 1, 2 \rangle = 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\ \langle 2, 2, 5, 2, 2, 1 \rangle = 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8, \\ \langle 2, 2, 3, 2, 2, 2 \rangle = 9, \textcolor{red}{\langle 2, 2, 4, 2, 2, 2 \rangle = 10}$$

$$q_0 = 0;$$

$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6], \textcolor{red}{10 [3]}$$

$$\delta : 0 \xrightarrow{\textcolor{blue}{/}} 1, 0 \xrightarrow{\textcolor{blue}{*}} 2, 0 \xrightarrow{\textcolor{blue}{(}} 3, 0 \xrightarrow{\textcolor{blue}{=}} 5, 0 \xrightarrow{\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}} 4, 0 \xrightarrow{\textcolor{blue}{;}} 6, \\ 1 \xrightarrow{\textcolor{blue}{/}} 1, 1 \xrightarrow{V \setminus \{\textcolor{blue}{/}\}} 7, 2 \xrightarrow{\textcolor{blue}{*}} 2, 2 \xrightarrow{V \setminus \{\textcolor{blue}{*}\}} 7, 3 \xrightarrow{\textcolor{blue}{(}} 8, 3 \xrightarrow{V \setminus \{\textcolor{blue}{(}\}} 7, 4 \xrightarrow{V} 7, \\ 5 \xrightarrow{\textcolor{blue}{=}} 5, 5 \xrightarrow{V \setminus \{\textcolor{blue}{=}\}} 7, 6 \xrightarrow{V} 7, 7 \xrightarrow{V} 7, 8 \xrightarrow{\textcolor{blue}{+}} 9, 8 \xrightarrow{V \setminus \{\textcolor{blue}{+}\}} 8, \\ \textcolor{red}{9 \xrightarrow{\textcolor{blue}{/}} 10}, \textcolor{red}{9 \xrightarrow{V \setminus \{\textcolor{blue}{/}\}} 8}$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

<i>lit</i> (1)	<i>int</i> (2)	<i>com</i> (3)	<i>spa</i> (4)	<i>op</i> (5)	<i>delim</i> (6)
$\delta$   <i>lit cif altesb</i>	$\delta$   <i>cif altesb</i>	$\delta$   $\frac{\textcolor{blue}{/}}{\textcolor{blue}{*}} \quad V \setminus \{\textcolor{blue}{/}, \textcolor{blue}{*}\}$	$\delta$   $\textcolor{blue}{=} \quad \textit{altesb}$	$\delta$   $\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{/}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}$ <i>altesb</i>	$\delta$   $\textcolor{blue}{;}$ <i>altesb</i>
0   1 2 2	0   1 2	0   1 5 5	0   1 2	0   1 2	0   1 2
1   1 1 2	1   1 2	1   5 2 5	1   1 2	1   2 2	1   2 2
2   2 2 2	2   2 2	2   2 3 2	2   2 2	2   2 2	2   2 2
F: 1	F: 1	3   4 2 2	F: 1	F: 1	F: 1
		4   5 5 5			
		5   5 5 5			
		F: 4			

$$\begin{aligned}
 Q : \langle 0, 0, 0, 0, 0, 0 \rangle &= 0, \langle 1, 2, 5, 2, 2, 2 \rangle = 1, \langle 2, 1, 5, 2, 2, 2 \rangle = 2, \\
 \langle 2, 2, 1, 2, 1, 2 \rangle &= 3, \langle 2, 2, 5, 2, 1, 2 \rangle = 4, \langle 2, 2, 5, 1, 2, 2 \rangle = 5, \\
 \langle 2, 2, 5, 2, 2, 1 \rangle &= 6, \langle 2, 2, 5, 2, 2, 2 \rangle = 7, \langle 2, 2, 2, 2, 2, 2 \rangle = 8 \\
 \langle 2, 2, 3, 2, 2, 2 \rangle &= 9, \langle 2, 2, 4, 2, 2, 2 \rangle = 10
 \end{aligned}$$

$$q_0 = 0;$$

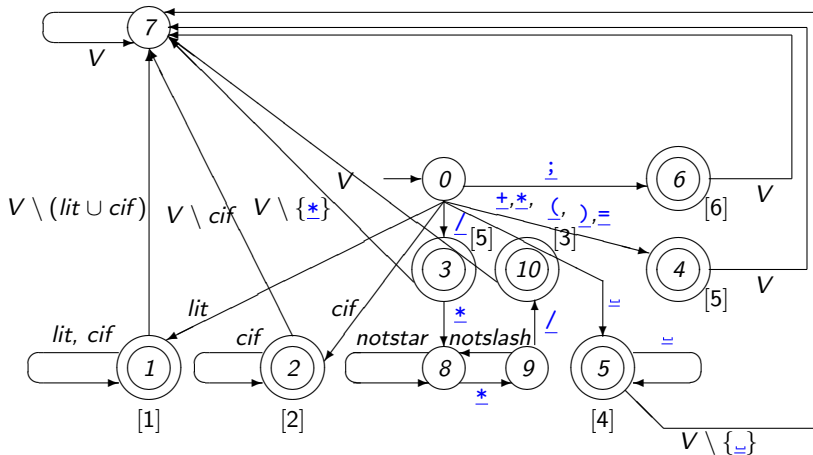
$$F : 1 [1], 2 [2], 3 [5], 4 [5], 5 [4], 6 [6], 10 [3]$$

$$\begin{aligned}
 \delta : 0 &\xrightarrow{\textit{lit}} 1, 0 \xrightarrow{\textit{cif}} 2, 0 \xrightarrow{\textcolor{blue}{/}} 3, 0 \xrightarrow{\textcolor{blue}{=}} 5, 0 \xrightarrow{\textcolor{blue}{+}, \textcolor{blue}{*}, \textcolor{blue}{(}, \textcolor{blue}{)}, \textcolor{blue}{=}} 4, 0 \xrightarrow{\textcolor{blue}{;}} 6, \\
 1 &\xrightarrow{\textit{lit}, \textit{cif}} 1, 1 \xrightarrow{V \setminus (\textit{lit} \cup \textit{cif})} 7, 2 \xrightarrow{\textit{cif}} 2, 2 \xrightarrow{V \setminus \textit{cif}} 7, 3 \xrightarrow{\textcolor{blue}{*}} 8, 3 \xrightarrow{V \setminus \{\textcolor{blue}{*}\}} 7, 4 \xrightarrow{V} 7, \\
 5 &\xrightarrow{\textcolor{blue}{=}} 5, 5 \xrightarrow{V \setminus \{\textcolor{blue}{=}\}} 7, 6 \xrightarrow{V} 7, 7 \xrightarrow{V} 7, 8 \xrightarrow{\textcolor{blue}{*}} 9, 8 \xrightarrow{V \setminus \{\textcolor{blue}{*}\}} 8, \\
 9 &\xrightarrow{\textcolor{blue}{/}} 10, 9 \xrightarrow{V \setminus \{\textcolor{blue}{/}\}} 8, \textcolor{red}{10} \xrightarrow{V} \textcolor{red}{7}
 \end{aligned}$$

# Algoritm reuniune disjunctă AFD

## Exemplu:

Desenat, automatul rezultat arată astfel (în dreptul fiecărei stări finale am scris între [ ] clasele de tokeni definite de ea):



# Analiza lexicală

Observăm că starea 7 nu este utilizată la recunoașterea cuvintelor - ea doar captează tranzițiile care nu ar trebui definite, pentru a avea un AFD total. Am putea elimina această stare și tranzițiile ce intră/ies din ea, obținând un AFD echivalent (inclusiv în ceea ce privește identificarea tipurilor de token), cu mai puține stări, dar parțial. Se folosesc următoarele rezultate/algoritmi (demonstrațiile sunt simple - exercițiu):



# Analiza lexicală

## Definiție:

Fie  $A = \langle Q, V, \delta, q_0, F \rangle$  un AFD.

O stare  $q \in Q$  este:

- **accesibilă**, dacă există  $\alpha \in V^*$  a.î.  $\delta(q_0, \alpha) = q$ ;
- **coaccesibilă**, dacă există  $\alpha \in V^*$  a.î.  $\delta(q, \alpha) \in F$ .

## Observație:

- doar stările ce sunt și accesibile și coaccesibile sunt utilizate efectiv la recunoașterea cuvintelor;
- starea inițială este accesibilă, stările finale sunt coaccesibile (aplicăm definiția pentru  $\alpha = \lambda$ );
- limbajul recunoscut de automat este nevid  $\Leftrightarrow$  există stări ce sunt și accesibile și coaccesibile  $\Leftrightarrow$  există cel puțin o stare finală accesibilă  $\Leftrightarrow$  starea inițială este coaccesibilă.

# Analiza lexicală

## Propoziție:

Fie  $A = \langle Q, V, \delta, q_0, F \rangle$  un AFD și  $\text{tok} : F \longrightarrow \{1, \dots, n\}$ ,  $n \in \mathbf{N}$ ,  $n \geq 1$ , o funcție.

Construim AFD  $A' = \langle Q', V, \delta', q'_0, F' \rangle$  c și funcția  $\text{tok}' : F \longrightarrow \{1, \dots, n\}$  astfel:

- $Q' = \{q_0\} \cup \{q \in Q : q \text{ este și accesibilă și coaccesibilă în } A\}$ ;
- $q'_0 = q_0$ ;
- pentru orice  $q \in Q'$  și  $a \in V$ , dacă  $\delta(q, a)$  este definit în  $A$  și  $\delta(q, a)$  este și accesibilă și coaccesibilă în  $A$ , atunci  $\delta'(q, a) = \delta(q, a)$ , altfel  $\delta'(q, a)$  este nedefinit;
- $F' = \{q \in F : q \text{ este accesibilă în } A\}$ .
- $\text{tok}' = \text{tok}|_{F'}$  (dacă  $F' = \emptyset$  atunci  $\text{tok}'$  este funcția banală).

Atunci:

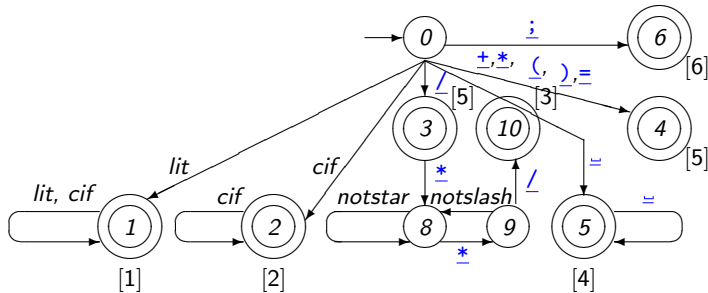
- pentru orice  $\alpha \in V^*$  avem  $\delta(q_0, \alpha) \in F$  d.d.  $\delta'(q'_0, \alpha) \in F'$ , iar în acest caz avem  $\text{tok}(\delta(q_0, \alpha)) = \text{tok}'(\delta'(q'_0, \alpha))$ ;
- în particular,  $L(A) = L(A')$ .

Deci, pentru recunoașterea tokenilor putem folosi automatul restrâns.

# Analiza lexicală

## Exemplu:

Aplicând algoritmul din propoziția precedentă automatului ce definea clasele de tokeni în exemplul precedent, obținem:



Acesta este automatul pe care îl vom folosi la recunoașterea tokenilor.

# Analiza lexicală

Analiza lexicală (bazată pe automatul de mai înainte) se efectuează astfel:

- pornind cu automatul din starea inițială, luăm pe rând caracterele din codul sursă și facem tranziții cu ele, până se termină codul sursă sau automatul se blochează (ajunge într-o stare din care cu caracterul curent nu mai există tranziție);
- dacă starea în care s-a ajuns (prin terminarea codului sursă sau blocarea automatului) este finală, atunci s-a mai identificat un token; tipul său este dat de starea finală respectivă (funcția *tok* a acesteia), iar valoarea tokenului este șirul consumat la recunoașterea lui; după emiterea tokenului, dacă nu s-a terminat codul sursă, restăm automatul (îl punem din nou în starea inițială) și facem tranziții cu următoarele caractere, începând cu caracterul cu care s-a blocat mai devreme (din starea inițială s-ar putea să existe tranziție cu el) până când codul sursă se termină sau automatul se blochează din nou;
- dacă starea respectivă nu este finală, atunci am întâlnit o eroare lexicală; în acest moment putem încheia scanarea semnalând eroare, sau putem încerca o recuperare a ei;

# Analiza lexicală

- pentru recuperarea erorii vom parcurge în sens invers drumul de stări urmat în automat la scanarea tokenului curent, introducând înapoi în input caracterele cu care s-au făcut tranzițiile;
- dacă vom întâlni o stare finală, ea va fi ultima stare finală prin care am trecut la scanarea tokenului curent, și atunci vom emite tokenul cu tipul dat de acea stare finală și valoarea dată de șirul consumat până la ea, după care vom reseta automatul și cu caracterul următor (ultimul reintrodus în input) vom face iar tranziții;
- dacă vom reveni în starea inițială fără să fi întâlnit nici o stare finală, atunci eroarea lexicală este atât de mare încât nu conține nici un token ca prefix; în acest caz vom opri scanarea, semnalând eroarea lexicală.

# Analiza lexicală

## Observație:

Faptul că tranzițiile se fac până automatul se blochează (sau se termină codul sursă) ne garantează, în cazul când ne-am oprit într-o stare finală, că aceasta este ultima stare finală prin care am trecut la scanarea tokenului curent. Astfel, fiecare token identificat este cel mai lung care se poate forma de la poziția de unde a început scanarea lui, deci se obține o interpretare orientată spre dreapta.

# Analiza lexicală

Putem formaliza algoritmul de analiză lexicală a unui șir (cod sursă)  $w \in V^*$  bazat pe AFD  $A = \langle Q, V, \delta, q_0, F \rangle$  și funcția  $tok : F \longrightarrow \{1, \dots, n\}$  astfel:

- Lucrăm cu configurații de forma  $(\alpha; s; \beta; \gamma)$  unde:
  - $\alpha$  este o stivă cu vârful la stânga, reprezentând partea rămasă din codul sursă scanat;
  - $s$  poate fi  $c$  (continuare),  $b$  (blocaj),  $s$  (succes),  $e$  (eroare);
  - $\beta$  este o stivă cu vârful la dreapta, conținând traseul (stări intercalate cu caractere) parcurs în automat la scanarea tokenului curent;
  - $\gamma$  este o stivă cu vârful la dreapta, conținând partea generată din interpretare;
- Configurația inițială este  $(w; c; q_0; \lambda)$ ;
- Configurație finală este una de forma:
  - $(\lambda; s; q_0; \gamma)$ ; este un caz de succes, în acest caz  $\gamma$  este interpretarea lui  $w$ ;
  - $(\alpha; e; q_0; \gamma)$ ; este un caz de eșec (eroare lexicală nerecuperabilă);

## Analiza lexicală

• Trecerea de la o configurație la alta se face în baza următoarelor reguli (în fiecare moment se poate aplica cel mult una dintre ele):

1.  $(a\alpha; c; \beta q; \gamma) \vdash (\alpha; c; \beta qar; \gamma)$ , dacă  $\delta(q, a) = r$ ;
2.  $(a\alpha; c; \beta q; \gamma) \vdash (a\alpha; b; \beta q; \gamma)$ , dacă  $\delta(q, a)$  nu e definită;
3.  $(\lambda; c; \beta q; \gamma) \vdash (\lambda; b; \beta q; \gamma)$ ;
4.  $(\alpha; b; \beta q; \gamma) \vdash (\alpha; c; q_0; \gamma\langle\varphi, tok(q)\rangle)$ , dacă  $q \in F$  și  $\alpha \neq \lambda$ ;  $\varphi$  este cuvântul din  $V^*$  care se obține din  $\beta q$  prin eliminarea stărilor intercalate între simbolurile din  $V$ ;
5.  $(\lambda; b; \beta q; \gamma) \vdash (\lambda; s; q_0; \gamma\langle\varphi, tok(q)\rangle)$ , dacă  $q \in F$ ;  $\varphi$  este cuvântul din  $V^*$  care se obține din  $\beta q$  prin eliminarea stărilor intercalate între simbolurile din  $V$ ;
6.  $(\alpha; b; \beta raq; \gamma) \vdash (a\alpha; b; \beta r; \gamma)$ , dacă  $q \notin F$ ;
7.  $(\alpha; b; q_0; \gamma) \vdash (\alpha; e; q_0; \gamma)$ , dacă  $q_0 \notin F$ ;

Obs: ultimele două reguli sunt considerate doar dacă vrem să încercăm recuperarea erorilor; altfel, în locul lor vom considera regula:

- 6'.  $(\alpha; b; \beta q; \gamma) \vdash (\alpha; e; \beta q; \gamma)$ , dacă  $q \notin F$ .



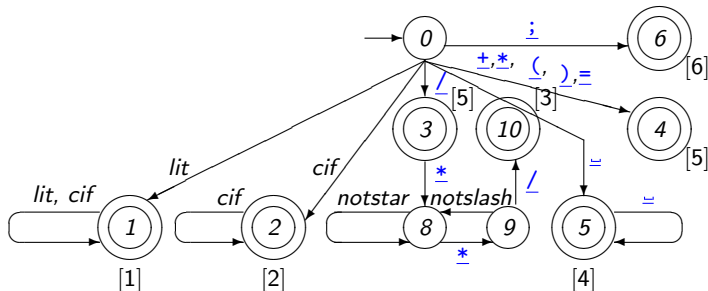
## Analiza lexicală

Să aplicăm algoritmul de analiză lexicală în cazul când AFD, *tok* și *w* sunt cele din exemplele anterioare:

```
(abc=(12+x30)+1; _y_=abc;; c; 0; λ)
```

## Exemplu:

## Analiza lexicală



$\vdash^2 (\underline{=(12+x30)+1}; \underline{\_y\_}=abc;; b; 0\underline{a}1\underline{b}1\underline{c}1; \lambda)$

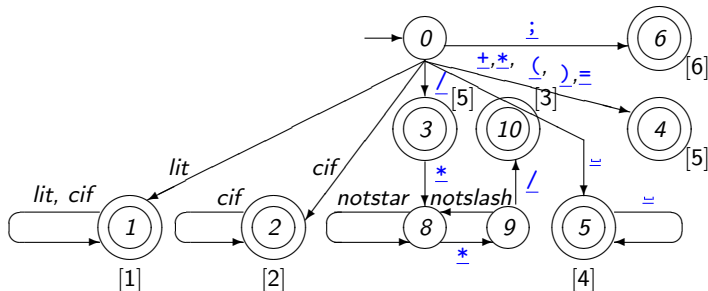
$\vdash^4 (\underline{=(12+x30)+1}; \underline{\_y\_}=abc;; c; 0; \langle \underline{abc}, 1 \rangle)$

$\vdash^1 (\underline{(12+x30)+1}; \underline{\_y\_}=abc;; c; 0\underline{=4}; \langle \underline{abc}, 1 \rangle)$

$\vdash^2 (\underline{(12+x30)+1}; \underline{\_y\_}=abc;; b; 0\underline{=4}; \langle \underline{abc}, 1 \rangle)$

Exemplu:

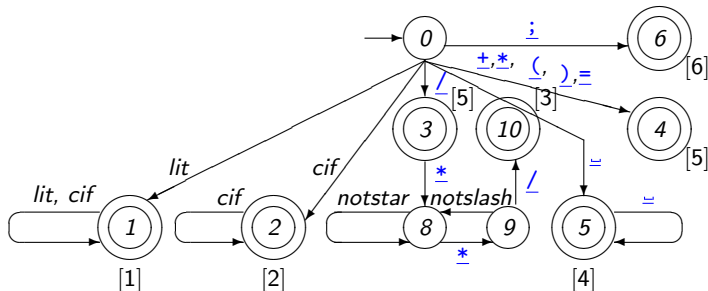
## Analiza lexicală



$$\begin{aligned} & \vdash^4 ((12+x30)+1; \text{cif}=\text{abc}; ; c; 0; \langle \text{abc}, 1 \rangle \langle \text{=}, 5 \rangle) \\ & \vdash^1 ((12+x30)+1; \text{cif}=\text{abc}; ; c; 0 \langle 4; \langle \text{abc}, 1 \rangle \langle \text{=}, 5 \rangle) \\ & \vdash^2 ((12+x30)+1; \text{cif}=\text{abc}; ; b; 0 \langle 4; \langle \text{abc}, 1 \rangle \langle \text{=}, 5 \rangle) \\ & \vdash^4 ((12+x30)+1; \text{cif}=\text{abc}; ; c; 0; \langle \text{abc}, 1 \rangle \langle \text{=}, 5 \rangle \langle \langle, 5 \rangle) \end{aligned}$$

## Exemplu:

## Analiza lexicală



<sup>1</sup>  
 $\vdash (\underline{2+x30})+1; \underline{\_y\_}=abc;; c; 0\underline{1}2; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle)$

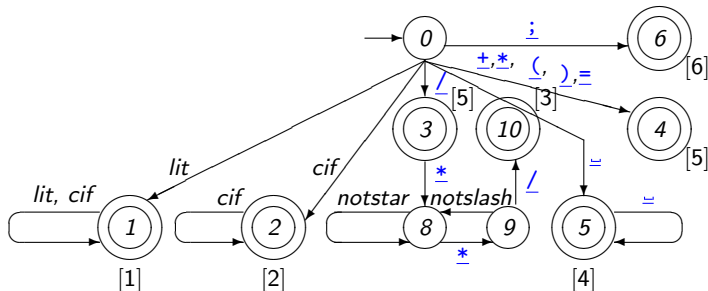
<sup>1</sup>  
 $\vdash (\underline{+x30})+1; \underline{\_y\_}=abc;; c; 0\underline{1}2\underline{2}2; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle)$

<sup>2</sup>  
 $\vdash (\underline{+x30})+1; \underline{\_y\_}=abc;; b; 0\underline{1}2\underline{2}2; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle)$

<sup>4</sup>  
 $\vdash (\underline{+x30})+1; \underline{\_y\_}=abc;; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle)$

## Exemplu:

## Analiza lexicală



<sup>1</sup>  
 $\vdash (\underline{x30})+1; \underline{y}=abc;; c; 0+4; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle)$

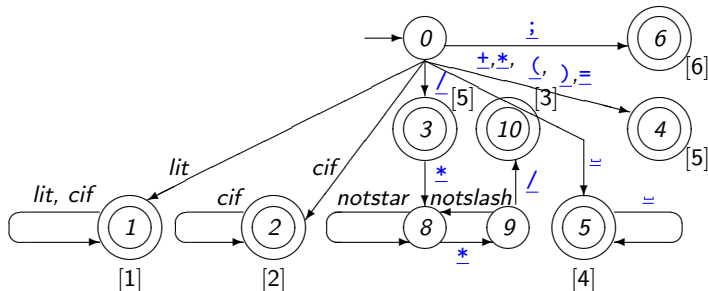
<sup>2</sup>  
 $\vdash (\underline{x30})+1; \underline{y}=abc;; b; 0+4; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle)$

<sup>4</sup>  
 $\vdash (\underline{x30})+1; \underline{y}=abc;; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle)$

<sup>1</sup>  
 $\vdash (\underline{30})+1; \underline{y}=abc;; c; 0\underline{x}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle)$

## Exemplu:

## Analiza lexicală



<sup>1</sup>  
 $\vdash (\underline{0})+1; \underline{y}=abc;; c; 0\underline{x}131; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle)$

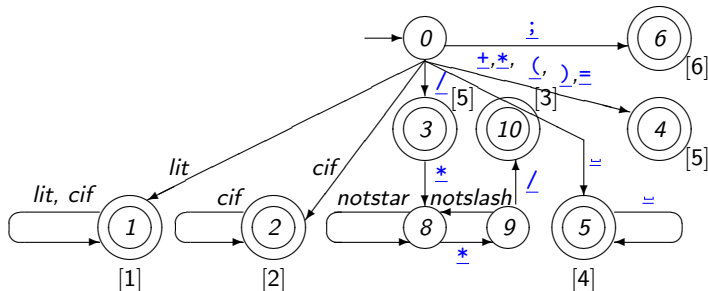
<sup>1</sup>  
 $\vdash (\underline{)}+1; \underline{y}=abc;; c; 0\underline{x}131\underline{0}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle)$

<sup>2</sup>  
 $\vdash (\underline{)}+1; \underline{y}=abc;; b; 0\underline{x}131\underline{0}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle)$

<sup>4</sup>  
 $\vdash (\underline{)}+1; \underline{y}=abc;; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x}30, 1 \rangle)$

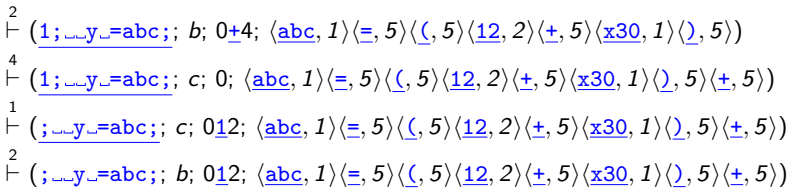
## Exemplu:

## Analiza lexicală



$$\begin{aligned}
 & \stackrel{1}{\vdash} (\underline{+1}; \underline{\text{y}} = \underline{\text{abc}}; ; \text{ c; } 0 \underline{4}; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle) \\
 & \stackrel{2}{\vdash} (\underline{+1}; \underline{\text{y}} = \underline{\text{abc}}; ; \text{ b; } 0 \underline{4}; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle) \\
 & \stackrel{4}{\vdash} (\underline{+1}; \underline{\text{y}} = \underline{\text{abc}}; ; \text{ c; } 0; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle) \\
 & \stackrel{1}{\vdash} (\underline{1}; \underline{\text{y}} = \underline{\text{abc}}; ; \text{ c; } 0 \underline{+4}; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle)
 \end{aligned}$$

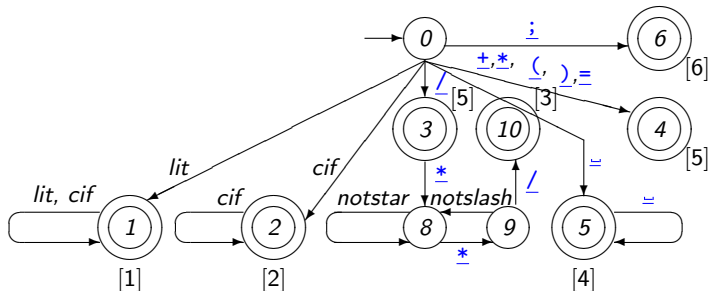
## Analiza lexicală





## Exemplu:

## Analiza lexicală



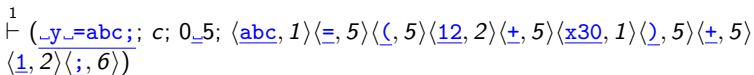
<sup>4</sup>  
 $\vdash (\underline{\text{;y=abc;}}; c; 0; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle)$

<sup>1</sup>  
 $\vdash (\underline{\text{;y=abc;}}; c; 0; \underline{6}; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle)$

<sup>2</sup>  
 $\vdash (\underline{\text{;y=abc;}}; b; 0; \underline{6}; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle)$

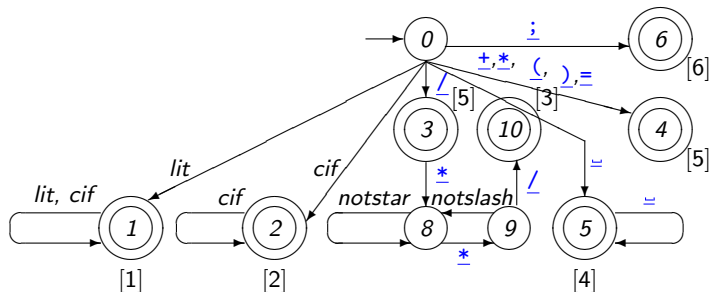
<sup>4</sup>  
 $\vdash (\underline{\text{;y=abc;}}; c; 0; \langle \underline{\text{abc}}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{\text{x30}}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle)$

## Analiza lexicală


$$\vdash^1 (\underline{y} = \underline{abc};; c; 0 \underline{5} \underline{5}; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{()}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{()}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle)$$
$$\vdash^2 (\underline{y} = \underline{abc};; b; 0 \underline{5} \underline{5}; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{()}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{()}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle)$$
$$\vdash^4 (\underline{y} = \underline{abc}; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{()}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{()}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{=}, 4 \rangle)$$

## Exemplu:

## Analiza lexicală



<sup>1</sup>  
 $\vdash (\underline{u}=\underline{abc};; c; 0\underline{y}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{u}, 4 \rangle)$

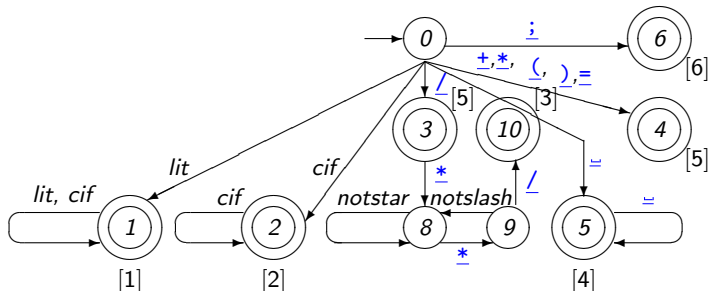
<sup>2</sup>  
 $\vdash (\underline{u}=\underline{abc};; b; 0\underline{y}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{u}, 4 \rangle)$

<sup>4</sup>  
 $\vdash (\underline{u}=\underline{abc};; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{u}, 4 \rangle \langle \underline{y}, 1 \rangle)$

<sup>1</sup>  
 $\vdash (\underline{u}=\underline{abc};; c; 0\underline{u}5; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{u}, 4 \rangle \langle \underline{y}, 1 \rangle)$

## Exemplu:

## Analiza lexicală



<sup>2</sup>  
 $\vdash (=abc; ; b; 0\_5; \langle abc, 1 \rangle \langle \_, 5 \rangle \langle (, 5 \rangle \langle 12, 2 \rangle \langle +, 5 \rangle \langle x30, 1 \rangle \langle ), 5 \rangle \langle +, 5 \rangle \langle 1, 2 \rangle \langle ;, 6 \rangle \langle \_, 4 \rangle \langle y, 1 \rangle)$

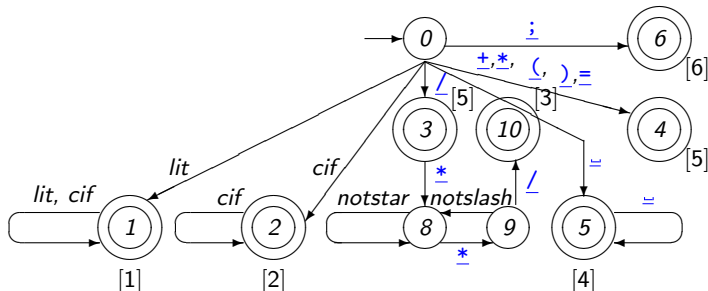
<sup>4</sup>  
 $\vdash (=abc; ; c; 0; \langle abc, 1 \rangle \langle \_, 5 \rangle \langle (, 5 \rangle \langle 12, 2 \rangle \langle +, 5 \rangle \langle x30, 1 \rangle \langle ), 5 \rangle \langle +, 5 \rangle \langle 1, 2 \rangle \langle ;, 6 \rangle \langle \_, 4 \rangle \langle y, 1 \rangle \langle \_, 4 \rangle)$

<sup>1</sup>  
 $\vdash (abc; ; c; 0=4; \langle abc, 1 \rangle \langle \_, 5 \rangle \langle (, 5 \rangle \langle 12, 2 \rangle \langle +, 5 \rangle \langle x30, 1 \rangle \langle ), 5 \rangle \langle +, 5 \rangle \langle 1, 2 \rangle \langle ;, 6 \rangle \langle \_, 4 \rangle \langle y, 1 \rangle \langle \_, 4 \rangle)$

<sup>2</sup>  
 $\vdash (abc; ; b; 0=4; \langle abc, 1 \rangle \langle \_, 5 \rangle \langle (, 5 \rangle \langle 12, 2 \rangle \langle +, 5 \rangle \langle x30, 1 \rangle \langle ), 5 \rangle \langle +, 5 \rangle \langle 1, 2 \rangle \langle ;, 6 \rangle \langle \_, 4 \rangle \langle y, 1 \rangle \langle \_, 4 \rangle)$

## Exemplu:

## Analiza lexicală



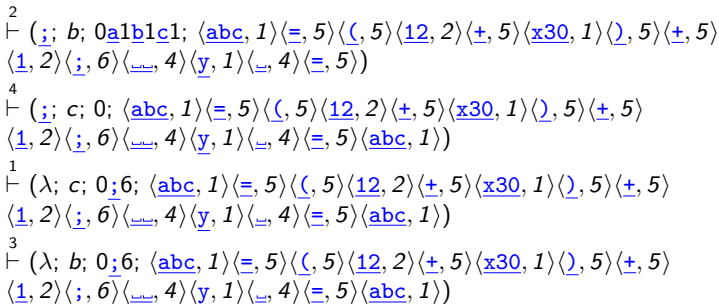
4  
 $\vdash (\underline{abc};; c; 0; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{=}, 4 \rangle \langle \underline{y}, 1 \rangle \langle \underline{=}, 4 \rangle \langle \underline{=}, 5 \rangle)$

1  
 $\vdash (\underline{bc};; c; 0\underline{a}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{=}, 4 \rangle \langle \underline{y}, 1 \rangle \langle \underline{=}, 4 \rangle \langle \underline{=}, 5 \rangle)$

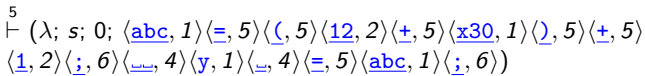
1  
 $\vdash (\underline{c};; c; 0\underline{a}1\underline{b}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{=}, 4 \rangle \langle \underline{y}, 1 \rangle \langle \underline{=}, 4 \rangle \langle \underline{=}, 5 \rangle)$

1  
 $\vdash (\underline{;}; c; 0\underline{a}1\underline{b}1\underline{c}1; \langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{=}, 4 \rangle \langle \underline{y}, 1 \rangle \langle \underline{=}, 4 \rangle \langle \underline{=}, 5 \rangle)$

## Analiza lexicală



## Analiza lexicală



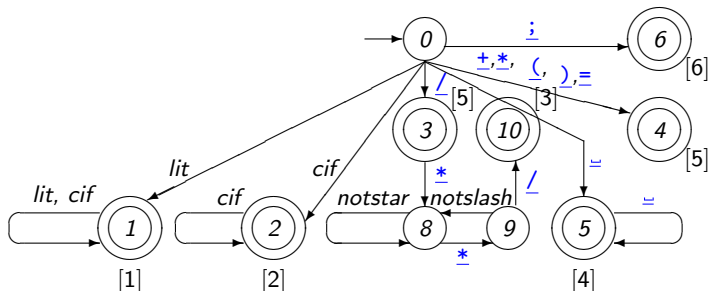
Exemplu:

$$\vdash^1 (\underline{2/*x12*}; c; 0\underline{a}1\underline{b}1\underline{1}; \lambda)$$



Exemplu:

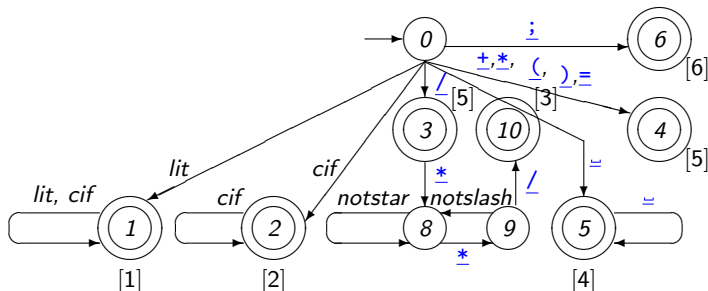
## Analiza lexicală



<sup>1</sup>  
 $\vdash (\underline{/*x12*}; c; 0\underline{a1b11121}; \lambda)$   
<sup>2</sup>  
 $\vdash (\underline{/*x12*}; b; 0\underline{a1b11121}; \lambda)$   
<sup>4</sup>  
 $\vdash (\underline{/*x12*}; c; 0; \langle \underline{ab12}, 1 \rangle)$   
<sup>1</sup>  
 $\vdash (\underline{*x12*}; c; 0\underline{/3}; \langle \underline{ab12}, 1 \rangle)$

## Exemplu:

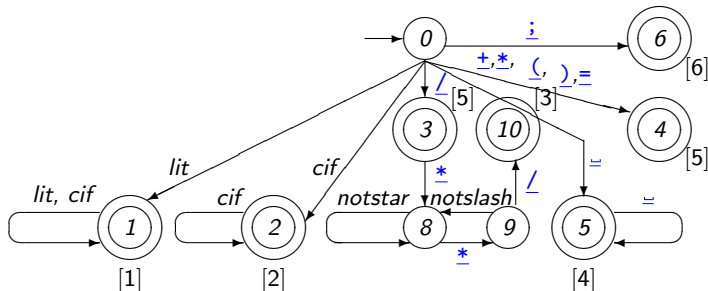
## Analiza lexicală



$\vdash^1 (\underline{x12*}; c; 0/\underline{3*8}; \langle \underline{ab12}, 1 \rangle)$   
 $\vdash^1 (\underline{12*}; c; 0/\underline{3*8x8}; \langle \underline{ab12}, 1 \rangle)$   
 $\vdash^1 (\underline{2*}; c; 0/\underline{3*8x818}; \langle \underline{ab12}, 1 \rangle)$   
 $\vdash^1 (\underline{*}; c; 0/\underline{3*8x81828}; \langle \underline{ab12}, 1 \rangle)$

## Exemplu:

## Analiza lexicală



<sup>1</sup>  
 $\vdash (\lambda; c; 0/3*8x81828*9; \langle \underline{ab12}, 1 \rangle)$

<sup>3</sup>  
 $\vdash (\lambda; b; 0/3*8x81828*9; \langle \underline{ab12}, 1 \rangle)$

(în acest moment, dacă nu am dori recuperarea erorilor, am efectua

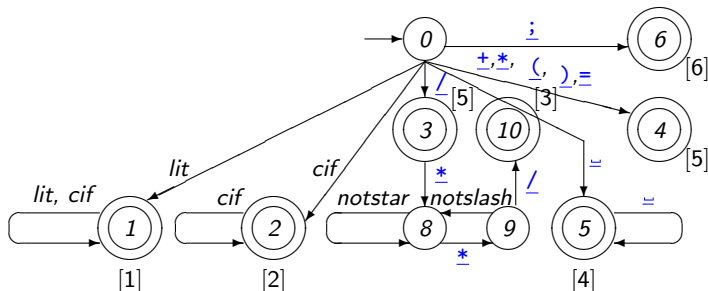
<sup>6'</sup>  
 $\vdash (\lambda; e; 0/3*8x81828*9; \langle \underline{ab12}, 1 \rangle)$  și am încheia analiza)

<sup>6</sup>  
 $\vdash (*; b; 0/3*8x81828; \langle \underline{ab12}, 1 \rangle)$

<sup>6</sup>  
 $\vdash (2*; b; 0/3*8x818; \langle \underline{ab12}, 1 \rangle)$

## Exemplu:

## Analiza lexicală



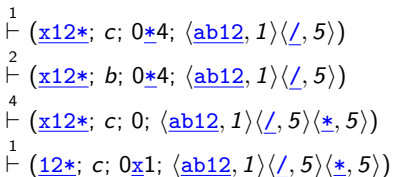
$\vdash^6 (\underline{12*}; b; 0/\underline{3*8}\underline{x8}; \langle \underline{ab12}, 1 \rangle)$

$\vdash^6 (\underline{x12*}; b; 0/\underline{3*8}; \langle \underline{ab12}, 1 \rangle)$

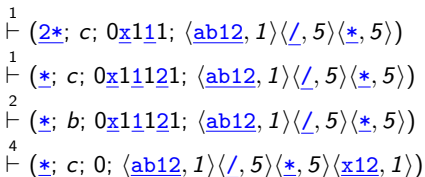
$\vdash^6 (\underline{*x12*}; b; 0/\underline{3}; \langle \underline{ab12}, 1 \rangle)$

$\vdash^4 (\underline{*x12*}; c; 0; \langle \underline{ab12}, 1 \rangle \langle \underline{/}, 5 \rangle)$

## Analiza lexicală

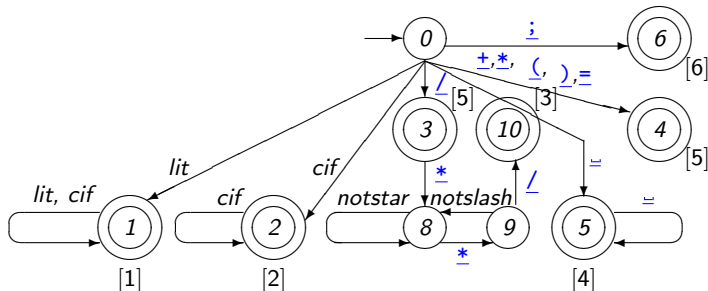


## Analiza lexicală



## Exemplu:

## Analiza lexicală



$$\begin{aligned}
 & \stackrel{1}{\vdash} (\lambda; c; 0\underline{*}4; \langle \underline{ab12}, 1 \rangle \langle \underline{/}, 5 \rangle \langle \underline{*}, 5 \rangle \langle \underline{x12}, 1 \rangle) \\
 & \stackrel{3}{\vdash} (\lambda; b; 0\underline{*}4; \langle \underline{ab12}, 1 \rangle \langle \underline{/}, 5 \rangle \langle \underline{*}, 5 \rangle \langle \underline{x12}, 1 \rangle) \\
 & \stackrel{5}{\vdash} (\lambda; s; 0; \langle \underline{ab12}, 1 \rangle \langle \underline{/}, 5 \rangle \langle \underline{*}, 5 \rangle \langle \underline{x12}, 1 \rangle \langle \underline{*}, 5 \rangle)
 \end{aligned}$$

# Analiza lexicală

## Observație:

*În ambele exemple de mai sus, dacă am fi lucrat cu automatul total, la fiecare token am fi parcurs codul sursă până la capăt (deoarece automatul nu se blochează), apoi ne-am fi întors până la starea finală corespunzătoare. Astfel, scanarea ar fi fost mai puțin eficientă.*



# Aspecte de implementare

Aspecte de implementare a unui analizor lexical (scanner):

- Structura scannerului:

Analizorul lexical poate fi implementat în mai multe feluri.

O idee este să-l construim a.î. să se poată integra în schema de compilator prezentată la început.

Astfel, el va fi implementat ca o componentă care poate fi invocată din exterior și la fiecare invocare să furnizeze un nou token.

De exemplu, el poate fi un obiect (în sensul programării orientate pe obiecte) având o metodă publică "gettoken()", care la fiecare apel va furniza un nou token (tokenul curent); în acest scop, obiectul scanner va reține prin membrii săi dată poziția în codul sursă până la care a ajuns cu scanarea.

## Aspecte de implementare

Conceptul de token poate fi și el implementat ca un obiect (sau structură) cu membri (cel puțin) pentru tipul și valoarea tokenului. La fiecare apel, metoda "gettoken()" va returna un obiect token.

Tipul poate fi un întreg (se fac niște convenții, ex: 0 = identificator, 1 = întreg. etc.) sau un element al unui tip enumerare.

Valoarea este un string.

Pentru a lucra unitar, "gettoken()" va semnala și apariția unei erori tot prin returnarea unui token, special - de ex. tipul să fie -1 iar valoarea să fie un întreg, desemnând poziția în cod unde a fost întâlnită eroarea (a se vedea și comentariile următoare referitoare la tabela de stringuri, unde vom arăta că valoarea tokenului poate fi returnată întotdeauna ca un întreg).

În cele ce urmează, dacă nu vom preciza altfel, aceasta este structura de scanner pe care o vom avea în vedere.

## Aspecte de implementare

- Tabela de stringuri:

Celelalte componente ale compilatorului nu au nevoie de valorile tokenilor (ca stringuri) ci doar de o informație mai restrânsă, din care să se poată determina dacă doi tokeni diferiți au sau nu aceeași valoare.

Astfel, putem face următoarea optimizare: obiectul scanner va gestiona o tabelă de stringuri, care va fi un membru privat al său (deoarece doar el are nevoie de ea) și în care va reține valorile unice ale tokenilor întâlniți la scanarea codului curent.

Fiecare apel al lui "gettoken()", după ce va identifica noul token, va căuta valoarea lui în tabela de stringuri; dacă nu o va găsi, o va insera; în ambele cazuri, poziția respectivă din tabel (unde a fost găsită/inserată) va fi furnizată ca valoare a tokenului (sub forma unui întreg, pointer, referință, etc.). Astfel, obiectul token returnat poate fi format întotdeauna din doi membri întregi - unul pentru tip, altul pentru valoare (dacă nu vom menționa altfel, așa vom considera în cele ce urmează).

Astfel, rezultă o economie de memorie în cazul manevrării simultane a mai multor tokeni cu aceeași valoare.

# Aspecte de implementare

În continuare, putem face diverse convenții/optimizări, de exemplu: fiecare operator să aparțină unui tip distinct de token - în acest caz în tabela de stringuri nu mai trebuie să reținem și operatorii (tipul tokenului va indica în mod unic valoarea), sau toți operatorii să aparțină unui același tip token - în acest caz vom reține în tabela de stringuri și operatorii.

La fel în cazul cuvintelor cheie.

# Aspecte de implementare

- Cuvinte cheie:

În diversele limbaje de programare întâlnim cuvinte cheie. Acestea au structura unui identificador, dar sunt tokeni de alt tip, așa că pentru recunoașterea lor nu putem folosi, teoretic, starea finală a identificatorilor.

De exemplu, putem conveni ca fiecare cuvânt cheie să definească un tip distinct de token (cu o unică valoare posibilă) și atunci să necesite o stare finală distinctă.

Existența cuvintelor cheie poate fi modelată în cadrul teoriei generale folosite până acum, în baza următoarelor rezultate:

- *Orice limbaj format dintr-un singur cuvânt este regulat.*

De exemplu limbajul  $\{\alpha\}$ , unde  $\alpha = a_1 \dots a_n$ ,  $a_i \in V^*$ , poate fi descris de ER  $a_1 \cdot \dots \cdot a_n$  sau recunoscut de un AFD lanț.

- *Orice limbaj finit este regulat.*

Într-adevăr, el este o reuniune finită de limbaje regulate; putem construi un AFD pentru el și care în plus să asocieze un tip distinct fiecărui cuvânt component (printr-o funcție "tok") cu algoritmul de reuninune disjunctă AFD prezentat mai înainte.

## Aspecte de implementare

- o *Clasa limbajelor regulate este închisă la diferență (adică dacă  $L_1$  și  $L_2$  sunt limbaje regulate, atunci  $L_1 \setminus L_2$  este un limbaj regulat).*

Pentru demonstrație se pornește cu două AFD totale pentru cele două limbaje și se construiește un AFD total pentru diferența lor, ca în algoritmul de reuniune disjunctă AFD, cu următoarele diferențe:

- se lucrează doar cu două automate, nu cu  $n$ ; deci stările noului AFD vor fi de forma  $\langle q^1, q^2 \rangle$ , cu  $q^1 \in Q_1$ ,  $q^2 \in Q_2$ ;
- nu mai avem de a face cu funcția "tok" (deoarece cuvintele recunoscute de noul AFD sunt printre cele recunoscute de primul AFD sursă - deci știm care este);
- În noul AFD vom avea  $\langle q^1, q^2 \rangle \in F$  d.d.  $q^1 \in F_1$  și  $q^2 \notin F_2$ ;  
(demonstrațiile sunt simple - exercițiu).

## Aspecte de implementare

○ Dacă notăm cu  $M$  limbajul descris de ER "idf" din exemplele anterioare și fixăm o submulțime finită  $C \subseteq M$ , atunci putem considera  $M \setminus C$  mulțimea identificatorilor și  $C$  mulțimea cuvintelor cheie; ambele sunt limbaje regulate, sunt disjuncte, iar prin tehnicile menționate mai sus putem construi câte un AFD pentru fiecare (care este sau se poate completa la unul total); cu aceste AFD și cu cele care definesc celelalte clase de token considerate (aceste clase sunt disjuncte de  $M$ , deci și de  $M \setminus C$  și  $C$ ) putem continua în maniera cunoscută, aplicând algoritmul de reuniune disjunctă AFD, pentru a obține un AFD capabil să recunoască tokenii și care va asocia un tip identificatorilor și cate un nou tip fiecărui cuvânt cheie.

## Aspecte de implementare

Din punctul de vedere al implementării putem evita însă complicarea automatului prin adăugarea de noi stări finale și putem folosi la recunoașterea cuvintelor cheie starea finală a identificatorilor, în felul următor:

dacă la apelul curent al lui "gettoken()" scanarea s-a oprit (prin blocarea automatului sau terminarea codului sursă) în starea finală a identificatorilor, șirul consumat (valoarea noului token) se caută mai întâi într-o tabelă de cuvinte cheie, care este fixată; dacă nu este găsit, se va returna un token identicator; dacă este găsit, se va returna un token de alt tip, corespunzător cuvântului cheie.



# Aspecte de implementare

- Spații:

În diversele limbaje de programare se folosesc și caractere albe (blank, tab, cap de linie, etc.) sau comentarii.

Prezența lor într-un cod sursă ajută la o mai bună înțelegere a codului de către om, dar nu au efect la execuție.

Pentru a clarifica ideile, vom numi **spațiu** o succesiune maximală de caractere albe și comentarii. În general, într-un program putem insera un spațiu între orice doi tokeni de alt tip, fără a afecta semantica programului (efectul la execuție). Spațiile se pot defini și identifica la fel ca și ceilalți tokeni - ele pot fi definite cu o ER iar pentru recunoașterea lor vom mai adăuga o stare finală în AFD.

Pentru compilator spațiile au ca singur scop să ajute scannerul să identifice ceilalți tokeni (care au efect la execuție). Odată identificați și livrați mai departe, aceștia sunt delimitați logic prin structurile de program ale compilatorului, deci spațiile nu mai sunt necesare - restul compilatorului nu mai are nevoie de spații.

## Aspecte de implementare

De aceea, scannerul va fi construit a.î. să recunoască spațiile la fel ca pe ceilalți tokeni, dar să nu le livreze mai departe. Mai exact, dacă la un apel al funcției "gettoken()" automatul s-a blocat în starea finală a spațiilor, să nu se returneze tokenul spațiu, să se rămână în apelul curent, să se scaneze și să se returneze următorul token (având în vedere maximalitatea spațiilor, următorul token nu va fi spațiu).

Pentru a trata unitar și cazul când spațiul a fost ultimul token (deci nu mai există alt token după el care să se returneze) se poate inventa un token - sfârșit de fișier (*EOF*), care să se returneze dacă partea răasă din cod este  $\lambda$ , indiferent dacă înainte a fost sau nu un token spațiu, iar restul compilatorului să fie scris a.î. să ignore acest ultim token furnizat.

## Aspecte de implementare

De exemplu la analizarea codului  $w = \underline{abc} = (12 + x30) + 1; \underline{\underline{y}} = \underline{abc};$  din exemplele precedente, scannerul va furniza interpretarea:

$\langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{y}, 1 \rangle$   
 $\langle \underline{=}, 5 \rangle \langle \underline{abc}, 1 \rangle \langle \underline{;}, 6 \rangle$

sau interpretarea:

$\langle \underline{abc}, 1 \rangle \langle \underline{=}, 5 \rangle \langle \underline{(}, 5 \rangle \langle \underline{12}, 2 \rangle \langle \underline{+}, 5 \rangle \langle \underline{x30}, 1 \rangle \langle \underline{)}, 5 \rangle \langle \underline{+}, 5 \rangle \langle \underline{1}, 2 \rangle \langle \underline{;}, 6 \rangle \langle \underline{y}, 1 \rangle$   
 $\langle \underline{=}, 5 \rangle \langle \underline{abc}, 1 \rangle \langle \underline{;}, 6 \rangle EOF$

# Analiza lexicală

Exerciții seminar:

1. Aplicați algoritmul ER  $\rightarrow$  AFD pentru ER  $(a+b)^*abb$
2. Aplicați algoritmul ER  $\rightarrow$  AFD pentru ER "com" din exemplele date.
3. Considerăm clasele de tokeni:
  1. identificatori (ca în exemplele date)
  2. cuvântul cheie: while
  3. cuvântul cheie: for
  4. operatori: +
  5. paranteză deschisă: (
  6. paranteză închisă: )
  7. delimitator: ;

Specificați aceste clase de tokeni prin expresii regulate.

Construiți AFD unic și funcția *tok* corespunzătoare pentru recunoașterea tokenilor.

4. Faceți demonstrațiile lăsate ca exercițiu în lecție.

# Analiza lexicală

## Teme laborator:

1. Scrieți un obiect scanner pentru un limbaj de programare uzual, respectând aspectele de implementare menționate în lecție. Integrați acest scanner într-un program care primește la intrare un cod (fișier) sursă și afișază lista tokenilor; pentru fiecare token va afișa tipul și valoarea ca stringuri.

Programul principal va conține un ciclu care la fiecare iterație apelează metoda "gettoken()" a scannerului și afișază tokenul returnat. Acest token are doi membri întregi. Stringul corespunzător membrului valoare se poate obține cu o metodă publică a scannerului, care primește întregul ca parametru și returnează o copie a valorii din tabela de stringuri. Stringul corespunzător membrului tip se poate obține cu o altă metodă publică a scannerului ce primește întregul respectiv ca parametru.

2. Scrieți un program care primește la intrare niște ER și o listă de cuvinte cheie și generează codul unui obiect scanner corespunzător (care să se poată integra de exemplu într-un program ca la tema 1). Se presupune predefinit tipul "identificator".