

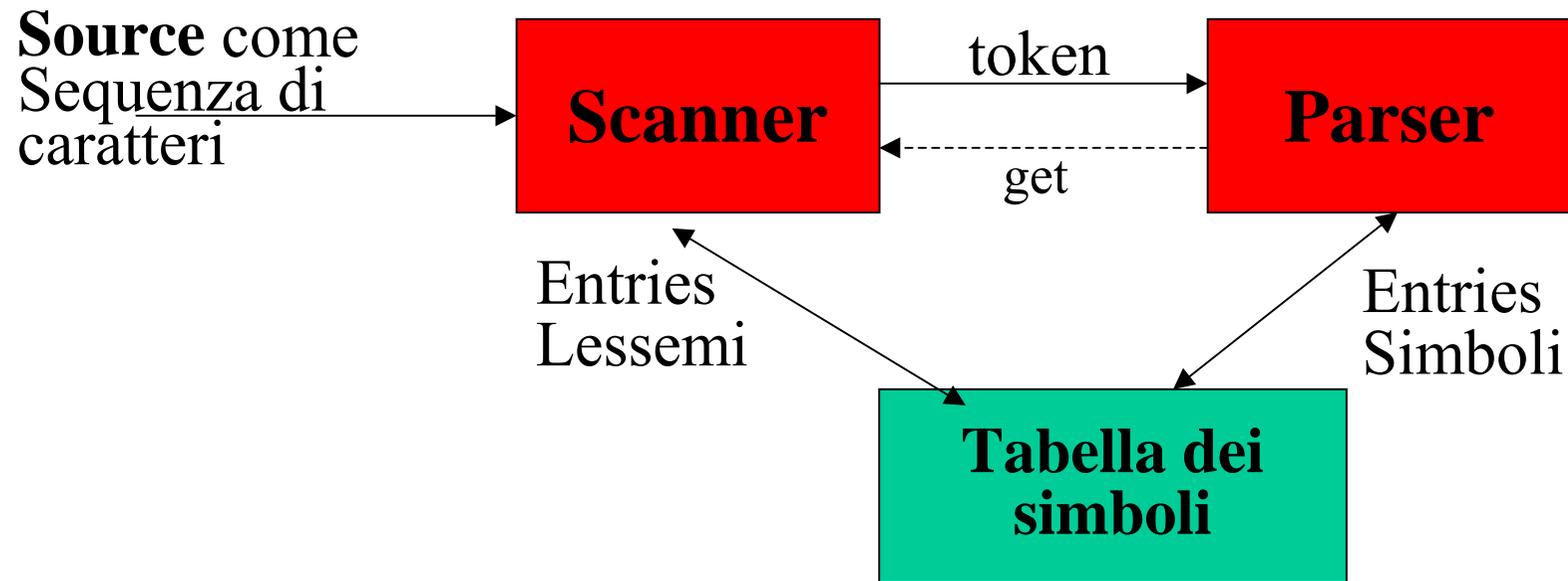
# ANALISI LESSICALE

Controllare e riconoscere il linguaggio utilizzato per scrivere i simboli (identificatori, numerali, keywords, separatori, delimitatori,...)

**linguaggio lessicale** in generale semplice

**linguaggio regolare,  
grammatiche regolari o lineari  
diagrammi di transizione**

# Una struttura una passata



Analisi lessicale (Scanner)  
guidata dall'analisi sintattica (Parser)

**token** = categoria lessicale

**lessema** = valore di una categoria

**pattern** = regole lessicali con cui definire ogni categoria lessicale

**esempio**

const pigreco = 3.1416

**const id rel num**

## Dove finiscono i valori?

In effetti generiamo coppie:  $\langle \text{token}, \text{attributo} \rangle$

$\langle \text{const}, \rangle$   $\langle \text{id}, k \rangle$   $\langle \text{rel}, = \rangle$   $\langle \text{num}, k+1 \rangle$

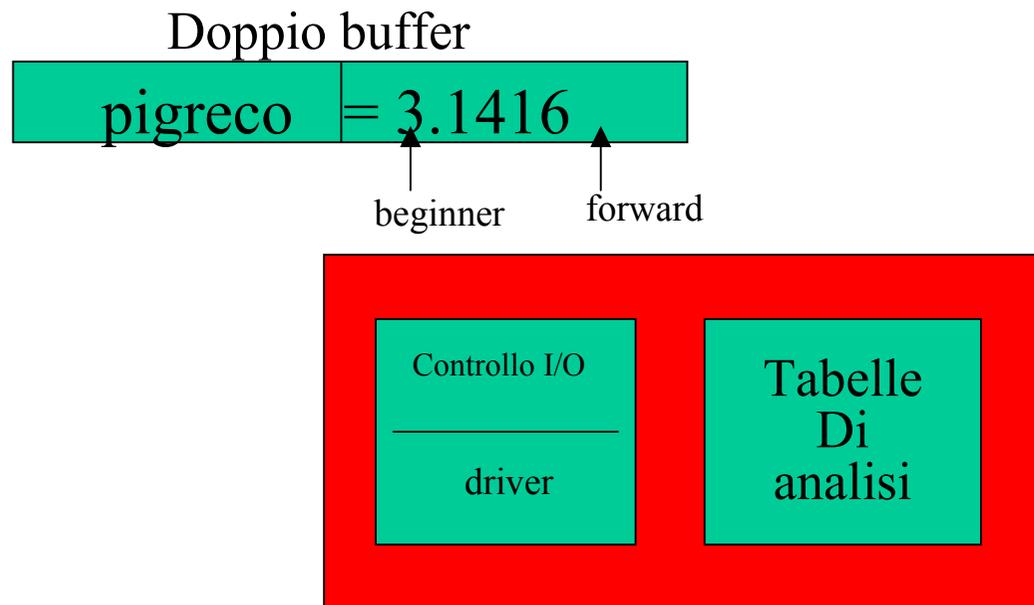
k	pigreco
K+1	3.1416

Symboltable

I **pattern** sono definiti con vari formalismi:

- diagrammi di transizione
- grammatiche regolari

Possiamo costruire riconoscitori basati sui patterns



# Espressioni e grammatiche regolari

Espressioni regolari  $E_\Sigma$  su  $\Sigma$  sono il

*minimo insieme* definito ricorsivamente da:

1.  $\epsilon \in E_\Sigma$
2.  $a \in E_\Sigma, \forall a \in \Sigma$
3.  $e_1 \cdot e_2 \in E_\Sigma, \forall e_1, e_2 \in E_\Sigma$
4.  $e_1 | e_2 \in E_\Sigma, \forall e_1, e_2 \in E_\Sigma$
5.  $e^* \in E_\Sigma, \forall e \in E_\Sigma$
6.  $e^i \in E_\Sigma, \forall e \in E_\Sigma$

# Hanno significato L su $2^{\Sigma^*}$

1.  $[\epsilon] = \{\lambda\}$
2.  $[a] = \{a\}$
3.  $[e_1.e_2] = \{uv \mid u \in [e_1], v \in [e_2]\} = [e_1] \times [e_2]$
4.  $[e_1|e_2] = \{u \mid u \in [e_1] \cup [e_2]\}$
5.  $[e^*] = \{u \mid u \in \bigcup_{i \in \mathbb{N}} [e]^i\}$
6.  $[e^i] = \{u \mid u \in [e]^i\}$       *abbreviazione*

$$\Sigma = \{a,b\}$$

$$\Sigma^* = \bigcup_{i \in \mathbb{N}} \Sigma \times \Sigma \dots \times \Sigma = \{\lambda, a, b, aa, ab, bb, ba, aaa, \dots\}$$

$$2^{\Sigma^*} = \{u \mid u \subseteq \Sigma^*\} = \{\{\lambda\}, \{a\}, \{b\}, \dots, \{\lambda, a\}, \dots, \{\lambda, a, b, ab\}, \dots\}$$

0 | 1 | ... | 9

$$[0 | 1 | \dots | 9] = \{0, 1, \dots, 9\}$$

(0 | 1 | ... | 9)\*

$$[(0 | 1 | \dots | 9)^*] = \{0, 1, \dots, 9, 00, 01, \dots, 3808, \dots\}$$

(1 | ... | 9).(0 | 1 | ... | 9)\*

$$[(1 | \dots | 9).(0 | 1 | \dots | 9)^*] = \{1, \dots, 9, 10, 11, \dots, 3808, \dots\}$$

# Grammatica (regolare) su $\Sigma$

$$G = \langle V, \Sigma, s \in V, P \rangle$$

$V$  *non terminali* (token ed ausiliari)

$\Sigma$  *terminali*

$s$  *simbolo distinto*

$P \subseteq V \times E_{\Sigma \cup V}$  *produzioni*

G e' una grammatica, affinche' sia  
**regolare** dobbiamo aggiungere:

- V sono totalmente ordinati, <
- per ogni  $v ::= e \in P$ ,  
 $e \in E_{\Sigma \cup \{v' < v\}}$

# Grammatiche come equazioni tra linguaggi

$$e \in \mathbf{E}_\Sigma \implies \mathbf{L}(e) \in 2^{\Sigma^*}$$

$$v ::= e \implies v = \mathbf{L}(e) \in (V \times 2^{\Sigma^*})$$

$$\{v_1 ::= e_1, \dots, v_k ::= e_k\} \implies \{ \underline{\mathbf{L}}(v_1), \dots, \underline{\mathbf{L}}(v_k) \mid \forall i, \underline{\mathbf{L}}(v_i) \in 2^{\Sigma^*} \text{ and } (v_i \equiv \mathbf{L}(e_i))_{\{\underline{\mathbf{L}}(v_1)/v_1, \dots, \underline{\mathbf{L}}(v_k)/v_k\}} \}$$

$$G = \langle V, \Sigma, s \in V, P \rangle \implies \underline{\mathbf{L}}(G) = \underline{\mathbf{L}}(s)$$

# Esempio

## Relazione tra grammatiche, espressioni e diagrammi

Un lessico per i numerali con grammatiche regolari

```
simple ::= digit digit*  
fract ::= simple.simple  
exp ::= fract E simple |  
        fract E (+|-) simple  
num ::= simple | fract | exp  
digit ::= 0 | 1 | ... | 9
```

```
s ::= d d*  
f ::= s.s  
e ::= f E s |  
        f E (+|-) s  
n ::= s | f | e  
d ::= 0 | 1 | ... | 9
```

chi e'  $\Sigma$  ?

quali strutture sono in G ?

Omettiamo il simbolo "." dell'operatore di giustapposizione

$$\mathbf{G} = \langle \mathbf{V}, \Sigma, s \in \mathbf{V}, \mathbf{P} \rangle$$

$$V = \{\text{simple, fract, exp, num}\}$$

$$\Sigma = \{0, 1, \dots, 9, ., +, -\}$$

$$s = \text{num}$$

**Notare** l'uso di metasimboli:

(,) per comporre espressioni regolari

# Riconoscere il lessico espresso da una regolare

## Automa a stati finiti

$\langle S, \Sigma, \text{move: } S \times \Sigma' \rightarrow S', S_0 \subseteq S, F \subseteq S \rangle$   
per  $S$  finito

### NFA (nondeterministico)

$\Sigma' = \Sigma \cup \{\epsilon\}$  mentre  $S' = 2^S$   
 $S_0$  ha cardinalita'  $> 1$

### DFA (deterministico)

$\Sigma' = \Sigma$  mentre  $S' = S$   
 $S_0$  e' singoletto

# La funzione *move*

Possiamo usare **grafi** simili ai diagrammi [G] o **tabelle** [T] per descrivere la funzione finita *move*

**Ad ogni stato S corrisponde:**

**G** un vertice del grafo

**T** un'entry della tabella con tante colonne  
quanta la cardinalita' di  $\Sigma$  (+1 per nondeter.)  
tante righe quanti gli stati

**Ad ogni  $\langle \langle s, a \rangle, S \rangle \in \text{move}$  corrisponde:**

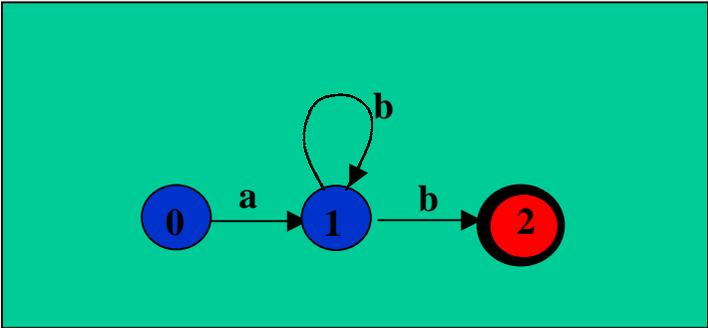
**G** un'arco  $\langle s, t \rangle$  etichettato  $a$  per ogni  $t \in S$

**T** il valore  $S$  per l'entry  $\langle s, a \rangle$

Esempio: consideriamo l'automa **L**

$S = \{0,1,2\}$ ,  
 $\Sigma = \{a,b\}$ ,  
 $move = \{ \langle \langle 0,a \rangle \{1\} \rangle, \langle \langle 1,b \rangle \{1,2\} \rangle \}$   
 $s_0 = 0$ ,  
 $F = \{2\}$  >

E' questo un automa deterministico  
oppure automa nondeterministico ?



	a	b	$\epsilon$
0	{1}		
1		{1,2}	
2			

# DOVE POSSIAMO DARE SIGNIFICATO AD UN AUTOMA A STATI FINITI

$$L = \langle S, \Sigma, \text{move: } S \times \Sigma' \rightarrow S', s_0 \in S, F \subseteq S \rangle$$

sui linguaggi regolari  $\mathfrak{R}$  su  $2^{\Sigma^*}$

## Usiamo funzioni di decisione

$\text{sem}(L) = f \in \Sigma^* \rightarrow \{accept, noaccept\}$

tale che: per ogni  $c_1c_2\dots c_n \in \Sigma^*$

$f(c_1c_2\dots c_n) = accept$

**se e solo se**

$c_1c_2\dots c_n$  appartiene al linguaggio riconosciuto

**definiamo la relazione  $\implies$**

$$\gamma \in \Sigma^*, s \in S \implies \gamma', s'$$

**se e solo se:**

$$\gamma = c\gamma' \text{ and } s' \in \textit{move}(c,s)$$

oppure

$$\gamma = \gamma' \text{ and } s' \in \textit{move}(\epsilon,s)$$

Sia  $\implies^*$  la chiusura transitiva di  $\implies$

Allora:

$f(\gamma) =$   
*accept se*  $\gamma, s_0 \implies^* \lambda, s \in F$   
*noaccept se* per nessun  $s \in F,$   
 $\gamma, s_0 \implies^* \lambda, s$

# Da un'espressione regolare $E$ a un NFA che riconosce il linguaggio $L(E)$

1.  $\varepsilon$

$\langle \{s_0\}, \{\}, \{\}, s_0 \in S, \{s_0\} \rangle$

2.  $a$  per ogni  $a \in \Sigma$

$\langle \{s_0, s_1\}, \{a\}, \{ \langle \langle s_0, a \rangle, s_1 \rangle \}, s_0 \in S, \{s_1\} \rangle$

### 3. $e1.e2$ per ogni $e1, e2 \in E$

$e1: \langle S_1, \Sigma_1, M_1, s_1, \{f_1\} \rangle$

$e2: \langle S_2, \Sigma_2, M_2, s_2, \{f_2\} \rangle$

$e1.e2: \langle S_1 \cup S_2, \Sigma_1 \cup \Sigma_2,$   
 $M_1 \cup M_2 \cup \{ \langle \langle f_1, \epsilon \rangle, \{s_2\} \rangle \},$   
 $s_1, \{f_2\} \rangle$

#### 4. $e_1|e_2$ per ogni $e_1, e_2 \in E$

$e_1: \langle S_1, \Sigma_1, M_1, s_1, \{f_1\} \rangle$

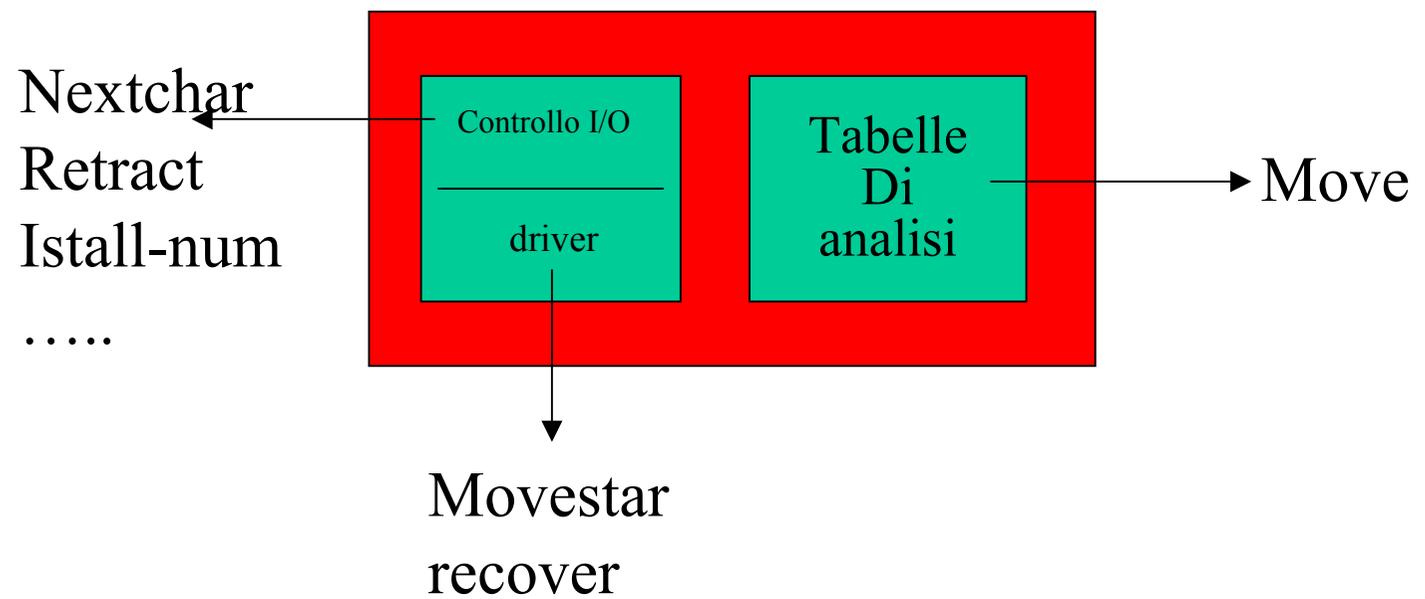
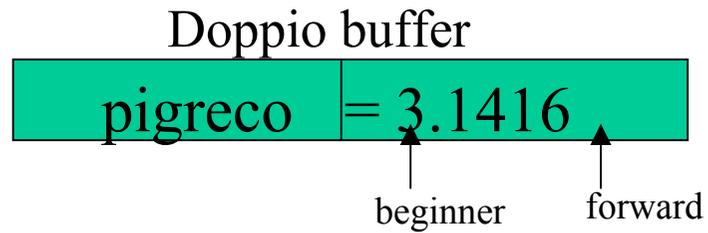
$e_2: \langle S_2, \Sigma_2, M_2, s_2, \{f_2\} \rangle$

$e_1|e_2: \langle S_1 \cup S_2 \cup \{s_{\text{new}}, f_{\text{new}}\}, \Sigma_1 \cup \Sigma_2,$   
 $M_1 \cup M_2 \cup \{ \langle \langle s_{\text{new}}, \epsilon \rangle, \{s_1, s_2\} \rangle,$   
 $\langle \langle f_1, \epsilon \rangle, \{f_{\text{new}}\} \rangle$   
 $\langle \langle f_2, \epsilon \rangle, \{f_{\text{new}}\} \rangle \},$   
 $s_{\text{new}}, \{f_{\text{new}}\} \rangle$

## 5. $e^*$ per ogni $e \in E$

$e: \langle S, \Sigma, M, s, \{f\} \rangle$

$e^* : \langle S \cup \{s_{\text{new}}, f_{\text{new}}\}, \Sigma,$   
 $M \cup \{ \langle \langle s_{\text{new}}, \epsilon \rangle, \{s, f_{\text{new}}\} \rangle,$   
 $\langle \langle f, \epsilon \rangle, \{s, f_{\text{new}}\} \rangle \},$   
 $s_{\text{new}}, \{f_{\text{new}}\} \rangle$



# Tre differenti driver:

- 1) **Diretto** Esecutore  
(gestisce il nondeterminismo con apposite strutture per il backtrack)
- 2) **Emulato** Interprete  
(interpreta il nodeterminismo come don't care: usa stati che sono insieme)
- 3) **Tradotto** Compilatore  
(compila il nodeterminismo: traduce un nondeterministico in un deterministico)

# 1) Un driver per NFA (DA)

```
Answer driver()
{states=push([Clos( {s0} ),pointerInput()]);
repeat
  answer='accept';
  [s,input]=pop(states);
  nextchar(c);
  while(c≠eof) and (s≠⊥) {
    S=move[s,c];
    push([Clos(S),pointerInput()]);
    [s,input]=pop(state);
    nextchar(c);}
  if (s∉ F) or (c≠eof) answer='noaccept';
until emptystack() or (answer='accept');
return (answer);
}
```

**States:** stack di controllo  
**Push:** inserisce [S,&  
S: insieme di stati  
&: puntatore input  
**Pop:** rimuove uno stato da S  
e copia & nell'input  
Se S e' singoletto l'intera  
coppia [S,&] e' rimossa  
dallo stack  
⊥: valore calcolato da move[s,c]  
quando la tabella non ha  
entry per move[s,c]

## 2) Un simulatore per NFA

move\* ci fornisce un riconoscitore

```
Answer movestar()  
  {S=Clos( {s0 } );  
  nextchar(c);  
  while(c≠eof) and (S≠∅) {  
    S=move1[S,c];  
    nextchar(c)};  
  if (S∩F) ≠ ∅ return 'accept';  
  else if return 'noaccept';  
}
```

Lineare  
(passo exp)

**Esercizio 1**

# Eliminiamo il nondeterminismo della relazione

$f(\gamma) =$   
*accept* se  $\text{move}^*(\gamma) \cap F \neq \{\}$   
*noaccept* altrimenti

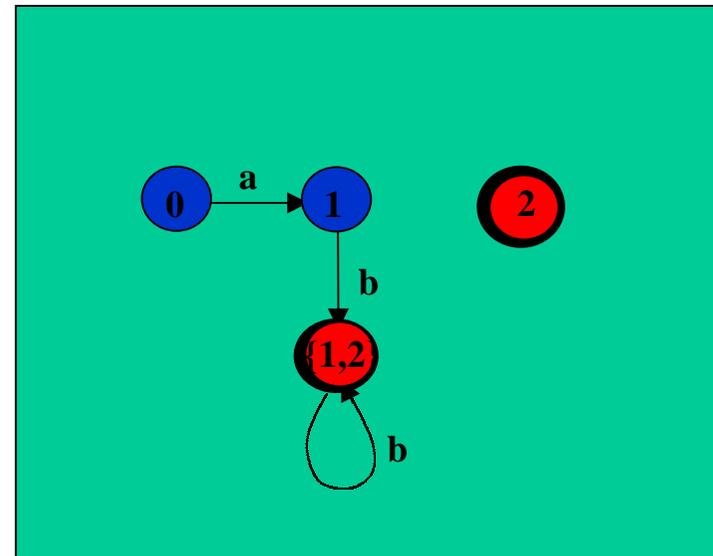
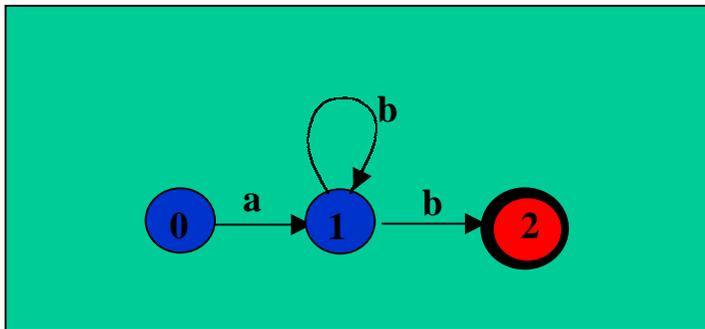
move

	a	b	$\epsilon$
0	{1}		
1		{1,2}	
2			

move\*

	a	b	$\epsilon$
0	{1}		
1		{1,2}	
2			
{1,2}		{1,2}	

# graficamente



Alcuni stati

possono diventare **inutili**

possono risultare **ridondanti**

## La funzione $move^*$ :

$$move1(S,c) = Clos(\cup_{s \in Clos(S)} \{move(s,c)\})$$

$$Clos(S) = S \cup Clos(\cup_{s \in S} move(s,\epsilon))$$

**Allora:**

$$move^*(c) = move1(\{s_0\},c)$$

$$move^*(c_1, \dots, c_{n-1}, c_n) = \\ move1(move^*(c_1, \dots, c_{n-1}), c_n)$$

Star: Set  $\rightarrow (\Sigma \rightarrow \text{Set})$

Osservare: Star non contiene colonna  $\epsilon$

$$\text{Star}[S]c = S'$$

$$\text{Star}[S]^+ = \bigcup_{c \in \Sigma} \text{Star}[S]c$$

$$\text{merge}(\text{move}, S) = \text{merge-row}(\{\text{move}[s] \mid s \in S\})$$

$$\forall_{c \in \Sigma} \text{merge-row}(\{R_1, \dots, R_k\})c = \left( \bigcup_{1 \leq i \leq k} \text{Clos}(R_i c) \right)$$

	digit	.	$\epsilon$
0	1		1
1	3	2	{1,3}
2		2	
3	0	4	
4	4	5	
5	6		
6	6		

move

$$\text{merge-row}(\{0,1,3\}) = \{0,1,3\} \mid \{2,4\}$$

### 3) Un convertitore NFA in DA

```
Table movestar()
  {EntryStar =  $\emptyset$ ;
  List= Clos(S0);
  while (List≠emptylist) {
    S=firstout(List);
    if (S $\notin$  EntryStar) {
      add(S,EntryStar);
      Star[S]= merge(move,S)};
    List=List+Star[S]+;
  }
  return Star;
}
```

Exp  
(*passo lineare*)

Esercizio 2

# **Una seconda tecnica per costruire riconoscitori**

e' basata sui

*diagrammi di transizione*

# Diagrammi di transizione su $\Sigma$

sono grafi diretti etichettati  $G(V,E,\Sigma)$

- vertici, detti stati, bipartiti in:
  - interni
  - finali
- archi  $E \subseteq V \times V$  (source,target)
- per ogni  $e \in E$ ,  $\text{label}(e) \in \Sigma$

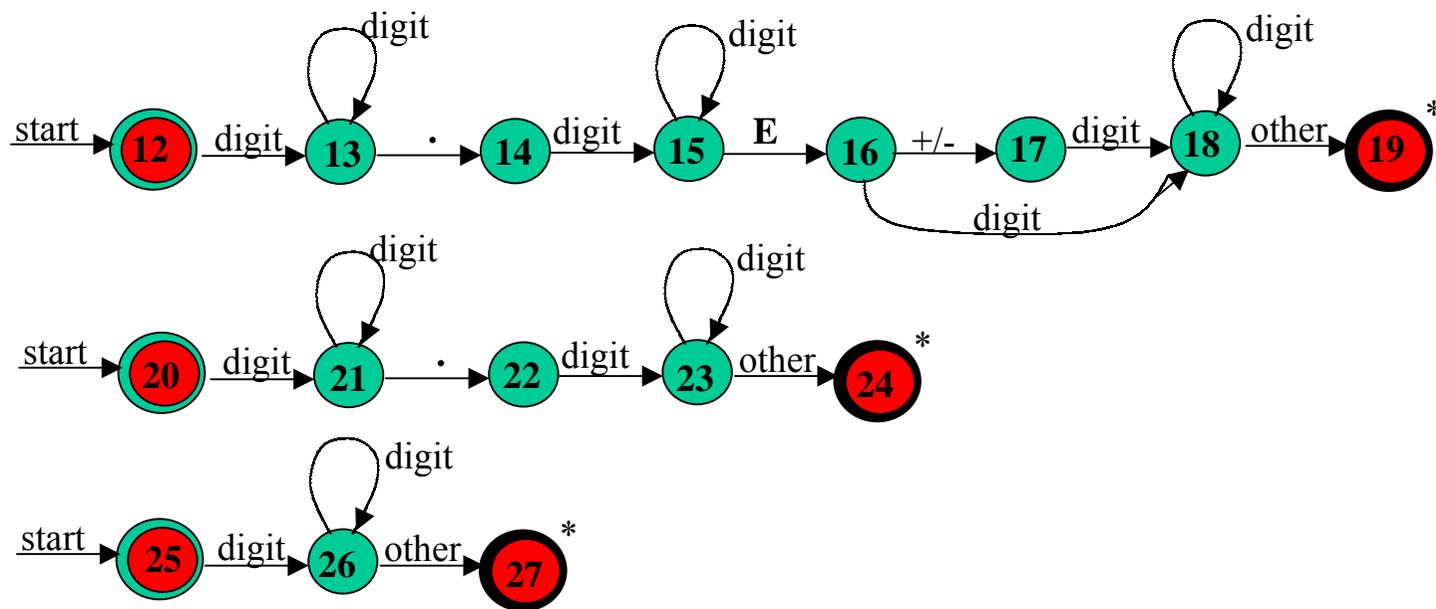
- nondeterministici

- deterministici

label(e1) = label(e2) only if  
source(e1) ≠ source(e2)

# Esempio (continua)

un lessico per i numerali con **diagrammi di transizione**



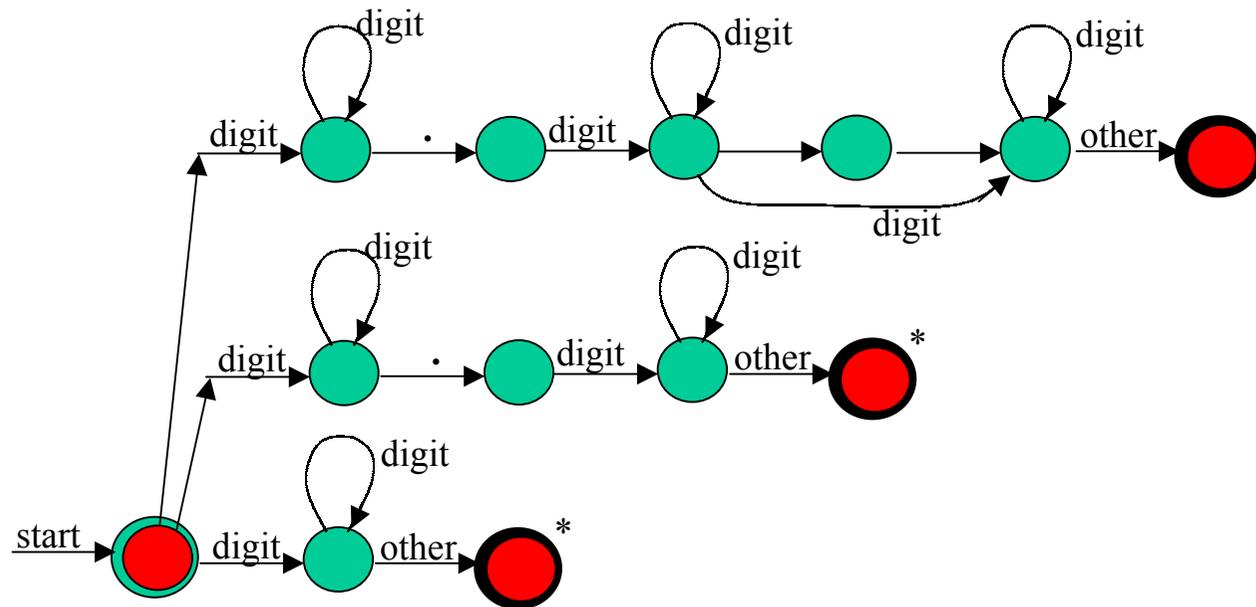
I diagrammi sono deterministici??

Come mai tanti diagrammi??

Le annotazioni \* e other fanno parte della definizione??

# I diagrammi non sono deterministici

Infatti essi corrispondono al diagramma:



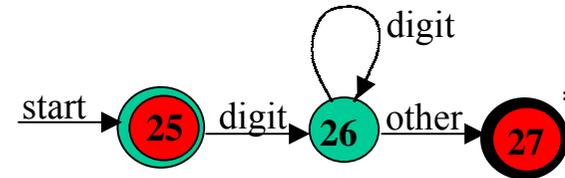
# COSTRUIAMO I RICONOSCITORI

diagrammi sintattici sono già adatti

ogni transizione può essere descritta  
da una struttura del tipo

**if** *predicate* **then** *action*

1. scriviamo **un frammento di codice** (nel metalinguaggio)  
per ogni stato di ogni diagramma



**stati interni**

```
25 : if (isdigit(c)) state=26;  
    else if state=fail()
```

```
26 : if (isdigit(c)) state=26;  
    else if state=27
```

C contiene l'ultimo carattere letto

## stati finali



```
27 :retract(1);  
    install-num();  
    return(num)
```

Abbiamo usato operazioni  
per il controllo dell'I/O  
per l'update della Symboltable

**La definizione** della corretta forma per

predicate e action

dipende dal metalinguaggio:

strutture dati e di controllo permesse.

abbiamo usato C:

i condizionali hanno quella forma

mentre variabili e funzioni sono:

**c**: variabile globale contenente carattere del buffer su cui e' posizionato forward (ultimo carattere letto)

**nextchar()**: avanza forward di 1

**retract(k)**: indietreggia forward di k posizioni

**state**: variabile contenente stato correntemente attraversato

**Start:** variabile contenente stato  
iniziale corrente riconoscimento

**install-num():** aggiorna la symboltable

**fail():** gestisce il fallimento

**isdigit(a):** riconosce un simbolo *digit*

## 2. componiamo tutti questi frammenti

useremo un *case\_statement* come discriminante

useremo un iteratore per  
leggere un carattere per volta  
riesaminare caratteri dopo *fail()*

otteniamo la seguente struttura di programma

```

int state = 0;
int lexical-value;

token nexttoken()
{ while(1) {
    switch (state) {
    case 0:....
    .....
    case 25: if (isdigit(c)) state=26;
            else if state=fail()
    case 26: if (isdigit(c)) state=26;
            else if state=27
    case 27: retract(1);
            install-num();
            return(num)
    .....
    }
    }
}

```

# Osservazioni sull'ambiente di lavoro del programma

Il parser invoca `nexttoken()` attendendosi una coppia `<token, valore>`.  
Il programma restituisce una coppia?

**NO:** restituisce solo il token (num,id...).

però una variabile *lexical-value* può essere resa globale e condivisa dal parser

`install-...:` oltre a modificare la *sym-table* pone in *lexical-value* l'entry relativa al nuovo token

Il **nonderminismo** sulla scelta dei vari diagrammi e' gestito:

ordinando i diagrammi in accordo allo  
stato iniziale

saltando da un diagramma all'altro in  
caso di fallimento [*fail()*]

```
int fail()
{forward = beginner;
  switch (start) {
    case 0:.....
    .....
    case 12: start=20; break;
    case 20: start=25; break;
    case 25: recover(); break;
  }
  return start;
}
```

ri-iniziamo la scansione della sequenza di caratteri

Partiamo da un altro diagramma

Stack  
beginner

Quando abbiamo esaminato tutte le alternative

Esponenziale

# Qual'e' la vera routine di gestione degli errori lessicali???

recover()

Differenze tra *fail()* e *recover()*

# Proprieta' interessanti (1)

L1 e L2 sono regolari.

Unione:  $L1 \cup L2$  e' regolare?

SI

Prodotto:  $L1 \times L2$  e' regolare?

SI

Intersezione:  $L1 \cap L2$  e' regolare?

SI

Complemento:  $C(L1)$  e' regolare?

SI

L e' un linguaggio

Decidibilita': L e' regolare?

NO

# Proprieta' interessanti (2)

L e' regolare

$A = \langle S, \Sigma, m, s_0, F \rangle$  e' automa per L

Th.(di iterazione)

Se  $x \in L$ ,  $|x| \geq \#S$ ,

Allora:  $\exists u, w, v$  tali che  $x = u.w.v$ ,  $u.w^k.v \in L$

Esercizi

fine