

PARSER LR(1)

# PARSER LR(0), SLR, LR(1), LALR(1): COSA HANNO IN COMUNE?

- Usano azioni di **shift** e **reduce**;
- Sono **macchine guidate da una tabella**:
  - sono **raffinamenti di LR(0)**
    - Calcolano un FSA usando la costruzione basata sugli item
    - **SLR**: usa gli **stessi item di LR(0)** e usa anche le informazioni dell'insieme **follow**
    - **LR(1)/LALR(1)**: un item contiene anche informazioni date dai simboli **lookahead**.
      - LALR(1) è una semplificazione di LR(1) per ridurre il numero degli stati
- Consentono di definire classi di grammatiche

Se il parser LR(0) (o SLR, LR(1), LALR(1)) calcolato dalla grammatica non ha conflitti shift/reduce o reduce/reduce, allora  $G$  è per definizione una grammatica LR(0) (o SLR, LR(1), LALR(1)).

# PANORAMICA SU LR PARSING

Le grammatiche **LR** sono più potenti delle **LL**.

**LR(0)** ha esclusivo interesse didattico.

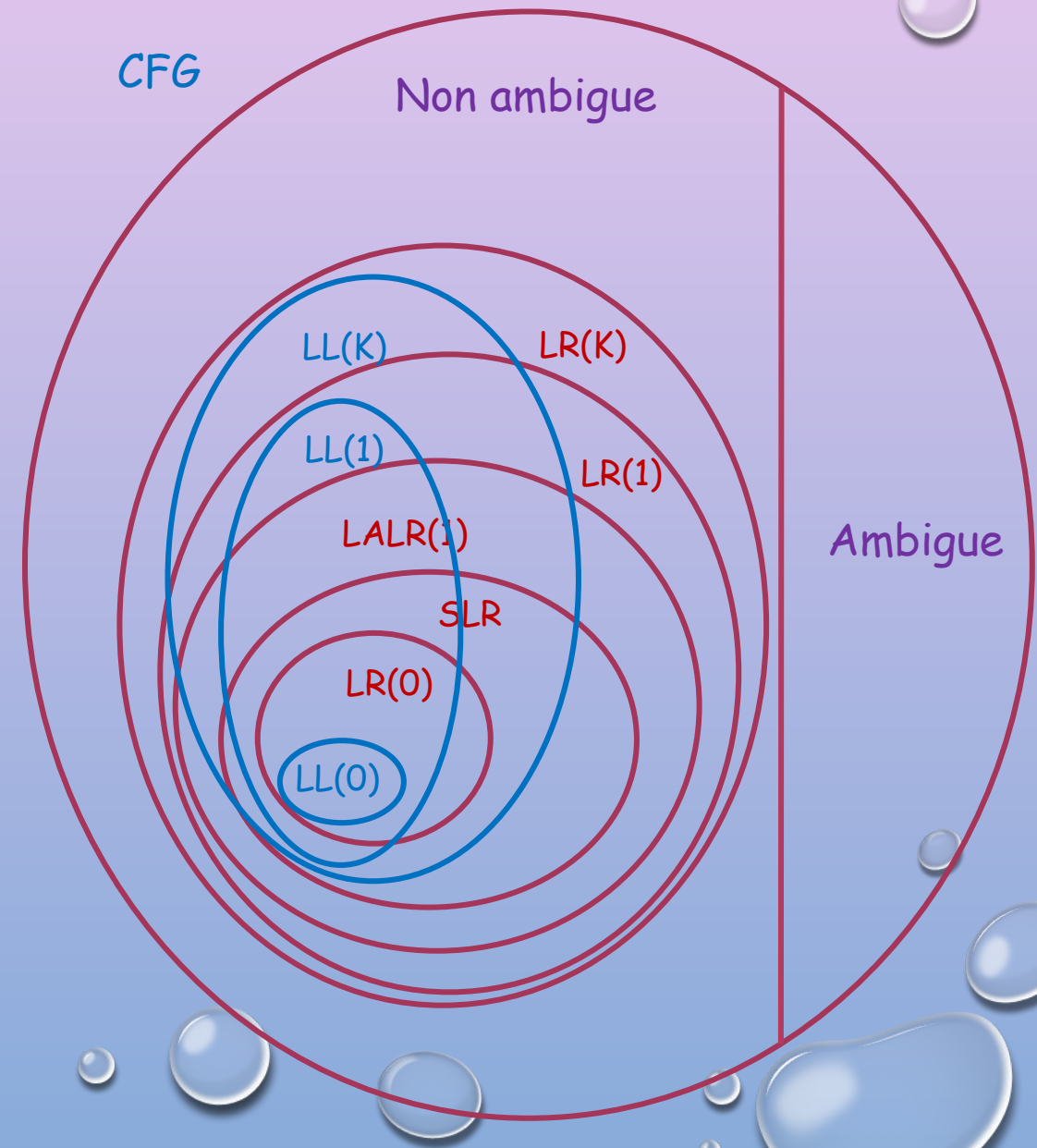
Esiste una classe di parser chiamato **Simple LR** (o **SLR**) che consentono il parsing di una famiglia un po' più vasta di Linguaggi.

La maggior parte dei linguaggi di programmazione ammettono una grammatica **LALR(1)**.

Molti generatori di parser usano questa classe.

**LR(1)** fornisce un parsing molto potente, ma l'implementazione è poco controllabile.

Si cercano grammatiche **LALR(1)** equivalenti.



# IL METODO LR

Nel metodo SLR, se siamo in uno stato che contiene l'item  $A \rightarrow \alpha.$  e il prossimo simbolo di input  $a$  è in  $\text{FOLLOW}(A)$ , dobbiamo fare una riduzione tramite  $A \rightarrow \alpha$

Se la grammatica non è SLR in generale questo può essere sbagliato: se sulla pila dei simboli c'è  $\beta\alpha$ ,  $\beta A$  potrebbe non essere seguito da  $a$  in nessuna forma sentenziale destra.



Per evitare questi errori, aggiungiamo agli item un simbolo terminale: un **item** nel metodo LR sarà del tipo  $[A \rightarrow \alpha.\beta, a]$  (un solo simbolo terminale se la grammatica è LR(1),  $k$  simboli terminali se è LR( $k$ ))

In realtà il simbolo terminale è usato solo negli item del tipo  $A \rightarrow \alpha.$ , ossia quelli che indurrebbero un reduce.

**Viene fatta una riduzione tramite  $A \rightarrow \alpha$  solo se il prossimo simbolo di input è  $a$**

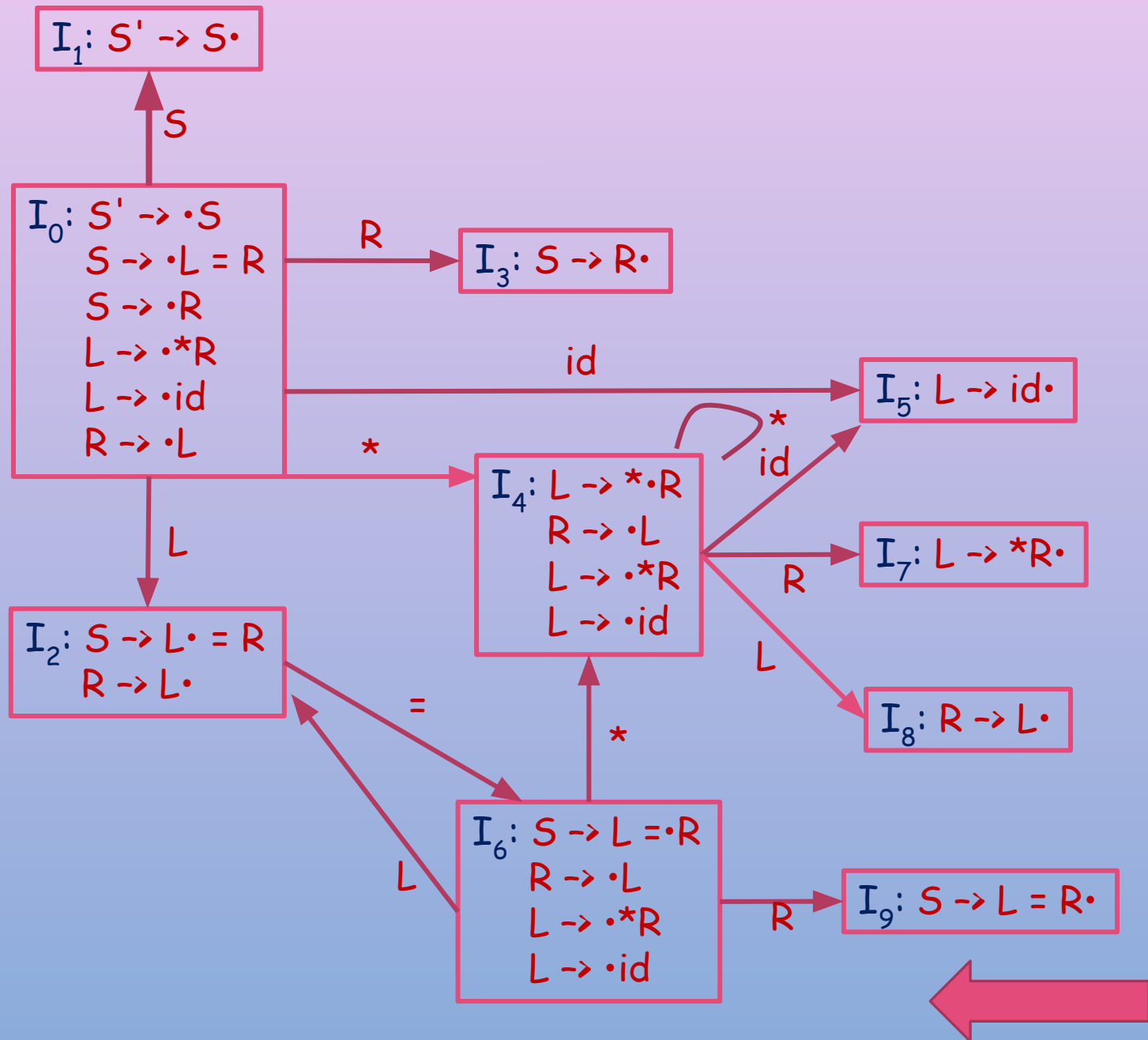


$S' \rightarrow S$   
 $S \rightarrow L = R \mid R$   
 $R \rightarrow L$   
 $L \rightarrow *R \mid id$

$FOLLOW(S') = \{\$, \}$   
 $FOLLOW(S) = \{\$, \}$   
 $FOLLOW(R) = \{\$, =\}$   
 $FOLLOW(L) = \{\$, =, ,\}$

Lo stato  $I_2$  permette un reduce  $R \rightarrow L$  oltre che un'azione di shift. Ma non esiste nessuna forma sentenziale destra della grammatica che inizia con  $=R..$

Quindi in questo caso è corretto fare lo shift e non la riduzione



# PARSER LR(1)

Consente di ovviare a molte ambiguità dei parser LR(0) al prezzo di una crescita sostanziale della complessità dell'algoritmo. Il metodo LR(1) è poco usato in pratica proprio perché poco efficiente:

- Si preferisce il più semplice **LALR(1)**
- L'automa riconoscitore per i parser LR(1) è simile agli automi **LR(0)**: cambiano gli **item** e le operazioni **closure**, **goto** e **reduce**
- **Fu il primo ad essere introdotto [Knuth 1965]**

Idea:

Nei parser **LR(1)** il lookahead è utilizzato **durante** la costruzione dell'automa (quindi si prendono in considerazione i simboli che veramente possano seguire un certo handle)

# COSTRUZIONE DELLA TABELLA LR(1)

La tabella di Parsing LR(1) si costruisce in maniera analoga alla tabella di parsing LR(0), con delle piccole differenze:

- Gli item LR(1) hanno la forma  $(A \rightarrow \alpha \cdot \beta, t)$  dove  $A \rightarrow \alpha \cdot \beta$  è un item LR(0) e  $t$  è un **token** (il simbolo terminale di lookahead) oppure  $t = \$$  (ciò indica il fatto che la sequenza  $\alpha$  si trova in cima alla pila e che alla testa dell'input c'è la stringa derivabile da  $\beta t$ ).
- Si compiono azioni di **Reduce**  $X \rightarrow \gamma$  solo in corrispondenza di quei simboli  $t$  per cui esiste lo stato di accettazione contiene l'item  $[X \rightarrow \gamma \cdot, t]$

# LE OPERAZIONI CLOSURE E GOTO

```
Function Closure(I);  
begin  
  J:=I;  
  repeat  
    for each item  $[A \rightarrow \alpha.X\beta, z]$  in J  
    for each rule  $X \rightarrow \gamma$   
    for each  $x \in \text{FIRST}(\beta z)$   
      add  $[X \rightarrow \gamma, x]$  to J;  
  until no more items can be added to J;  
  Return J;  
end
```

```
Function Goto(I, X);  
begin  
  J:= insieme degli item  $[A \rightarrow \alpha X \beta, a]$   
  tali che  $[A \rightarrow \alpha.X\beta, a]$  è in I;  
  return CLOSURE(J);  
end
```

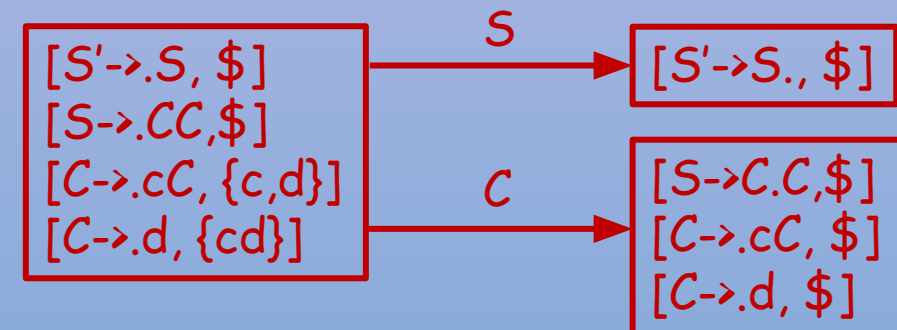
Esempio:

```
S' -> S  
S -> CC  
C -> cC | d
```

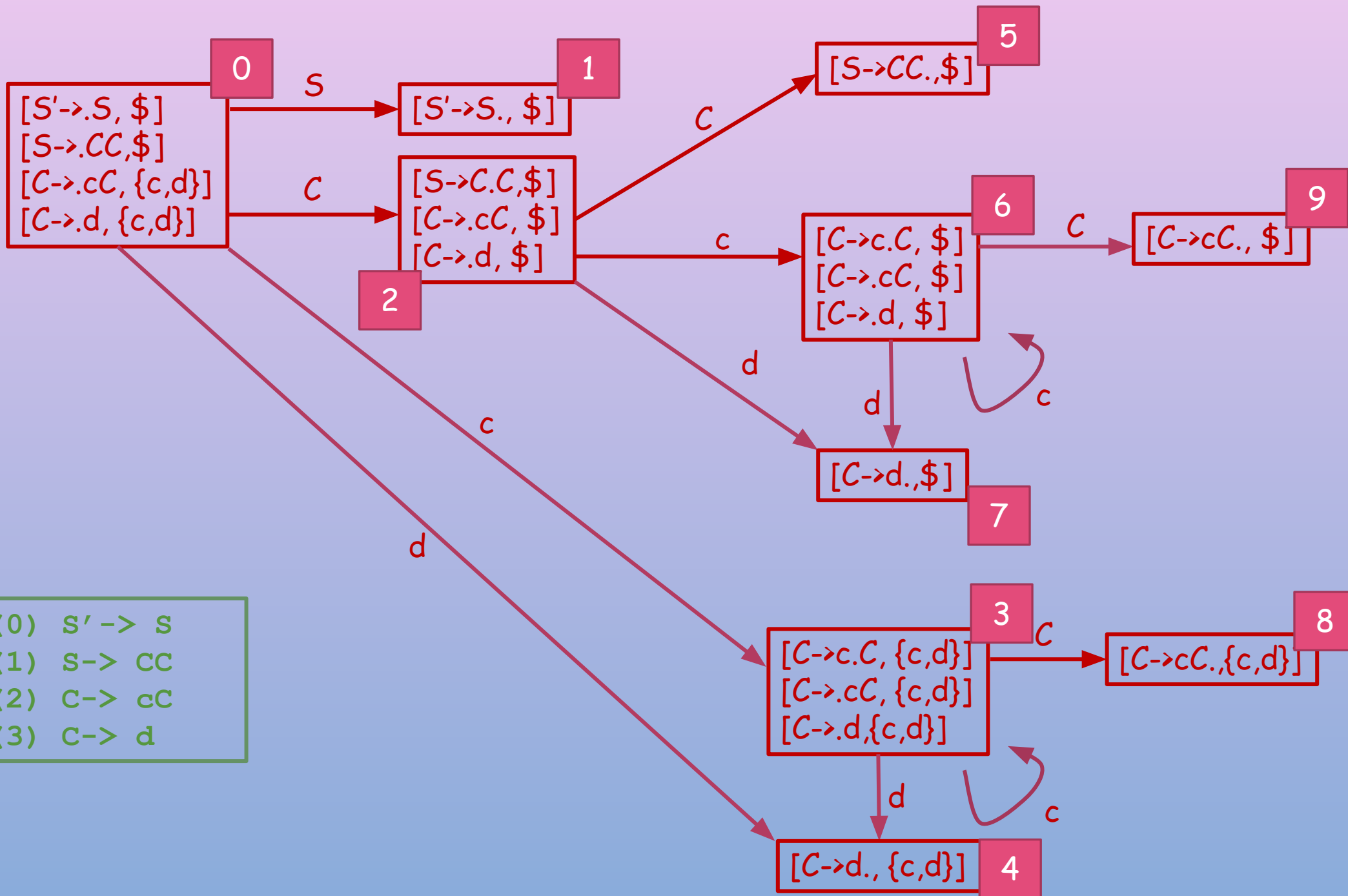
Closure ( $[S' \rightarrow .S, \$]$ )

$[S' \rightarrow .S, \$]$ ,  
 $[S \rightarrow .CC, \$]$   
 $[C \rightarrow .cC, \{c, d\}]$ ,  
 $[C \rightarrow .d, \{c, d\}]$

First( $\$$ )= $\$$   
First( $C\$$ )= $\{c, d\}$







# COSTRUZIONE DELLA TABELLA LR(1)

La tabella  $M$  è strutturalmente simile a quella LR(0)

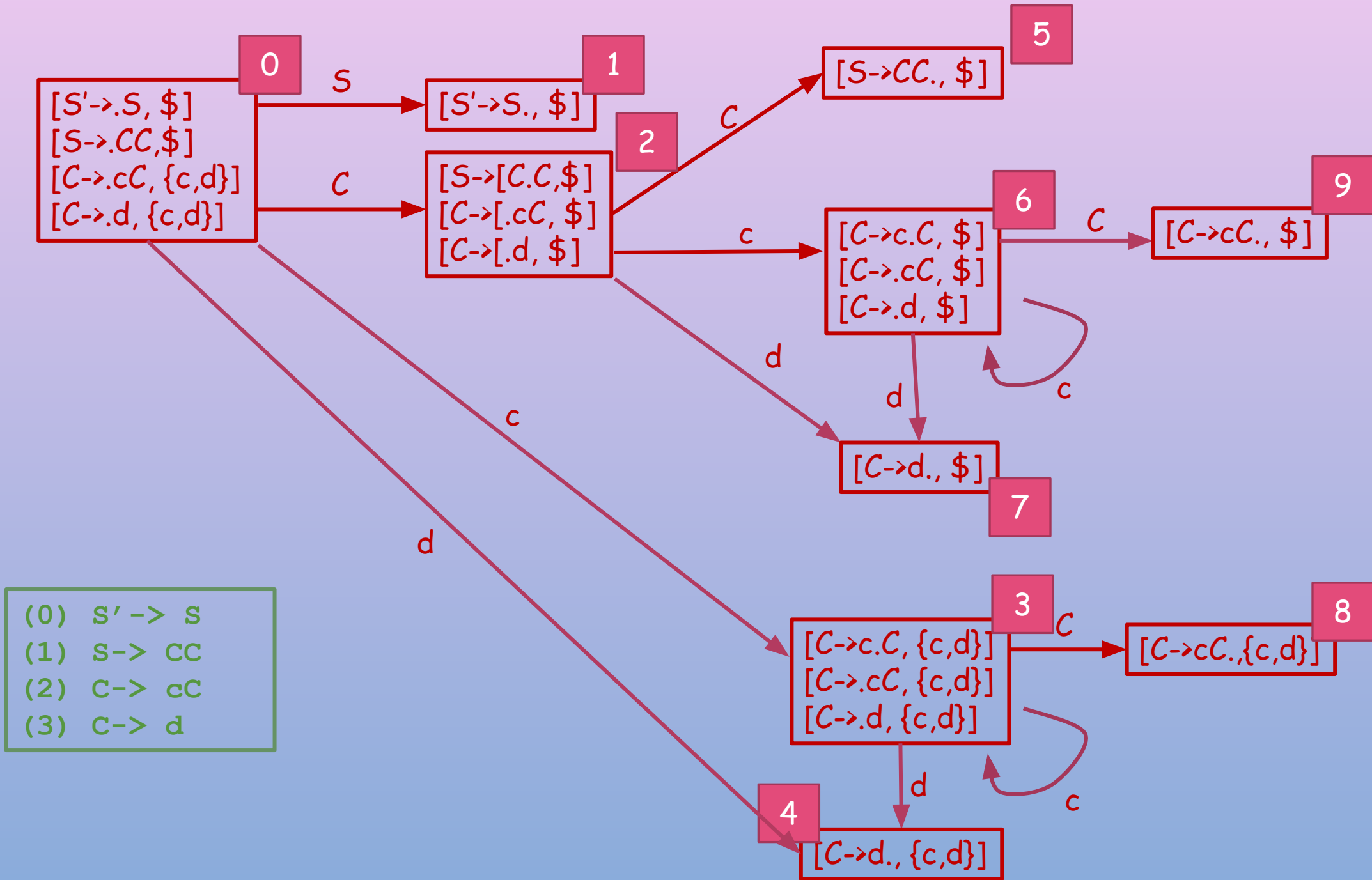
**Azioni shift:** se dallo stato  $q$  esiste una transizione  $t$  (simbolo terminale) nello stato  $q'$ , inserire **shift**  $q'$  in  $M(q, t)$

**Azioni goto:** se dallo stato  $q$  esiste una transizione  $S$  (simbolo non-terminale) in  $q'$ , inserire **goto**  $q'$  in  $M(q, S)$

**Azioni reduce:** se lo stato  $q$  contiene un item LR(1) del tipo  $[X \rightarrow \gamma ., t]$ , con  $t$  simbolo terminale e  $X$  diverso da  $S'$  inserire **reduce**  $X \rightarrow \gamma$  in  $M(q, t)$

**Azione accept:** se lo stato  $q$  contiene l'item  $[S' \rightarrow S ., \$]$ , allora inserire **accept** in  $M(q, \$)$

**Osservazione:** il parser LR(1) ridurrà soltanto quando il simbolo in testa all'input sarà  $\underline{t}$ . Lo stato iniziale è quello costruito da  $[S' \rightarrow .S, \$]$ .



(0)  $S' \rightarrow S$

(1)  $S \rightarrow CC$

(2)  $C \rightarrow cC$

(3)  $C \rightarrow d$

Stato	c	d	\$	S	C
0	Shift 3	Shift 4		Goto 1	Goto 2
1			Accept		
2	Shift 6	Shift 7			Goto 5
3	Shift 3	Shift 4			Goto 8
4	R. $C \rightarrow d$	R. $C \rightarrow d$			
5			R. $S \rightarrow CC$		
6	Shift 6	Shift 7			Goto 9
7			R. $C \rightarrow d$		
8	R. $C \rightarrow cC$	R. $C \rightarrow cC$			
9			R. $C \rightarrow cC$		

PARSER LALR(1)

# PARSING LALR(1)

Le tabelle di analisi **LR(1)** sono di solito di **ordini di grandezza maggiori** di quelle **LR(0)**: se una tabella **LR(0)** di un linguaggio di programmazione è intorno ai 10 KB, una tabella **LR(1)** dello stesso linguaggio è intorno ai MB, con il doppio della memoria per costruirla.

**Osservazione:** se consideriamo l'automa **LR(1)** della grammatica

$S' \rightarrow S, S \rightarrow CC, c \rightarrow cC \mid d$

e ignoriamo i lookahead, alcune coppie di stati sono identici: per esempio le coppie di stati **(8, 9)**, **(4, 7)**, **(3, 6)**, a prescindere dal simbolo di lookahead, sono uguali. Diciamo che queste coppie hanno lo stesso **core**.

Il parser **LALR(1)** consiste nell'identificare questi stati, combinando i loro lookahead, con l'obiettivo di ottenere un DFA **LR(1)** simile al DFA **LR(0)**. Infatti:

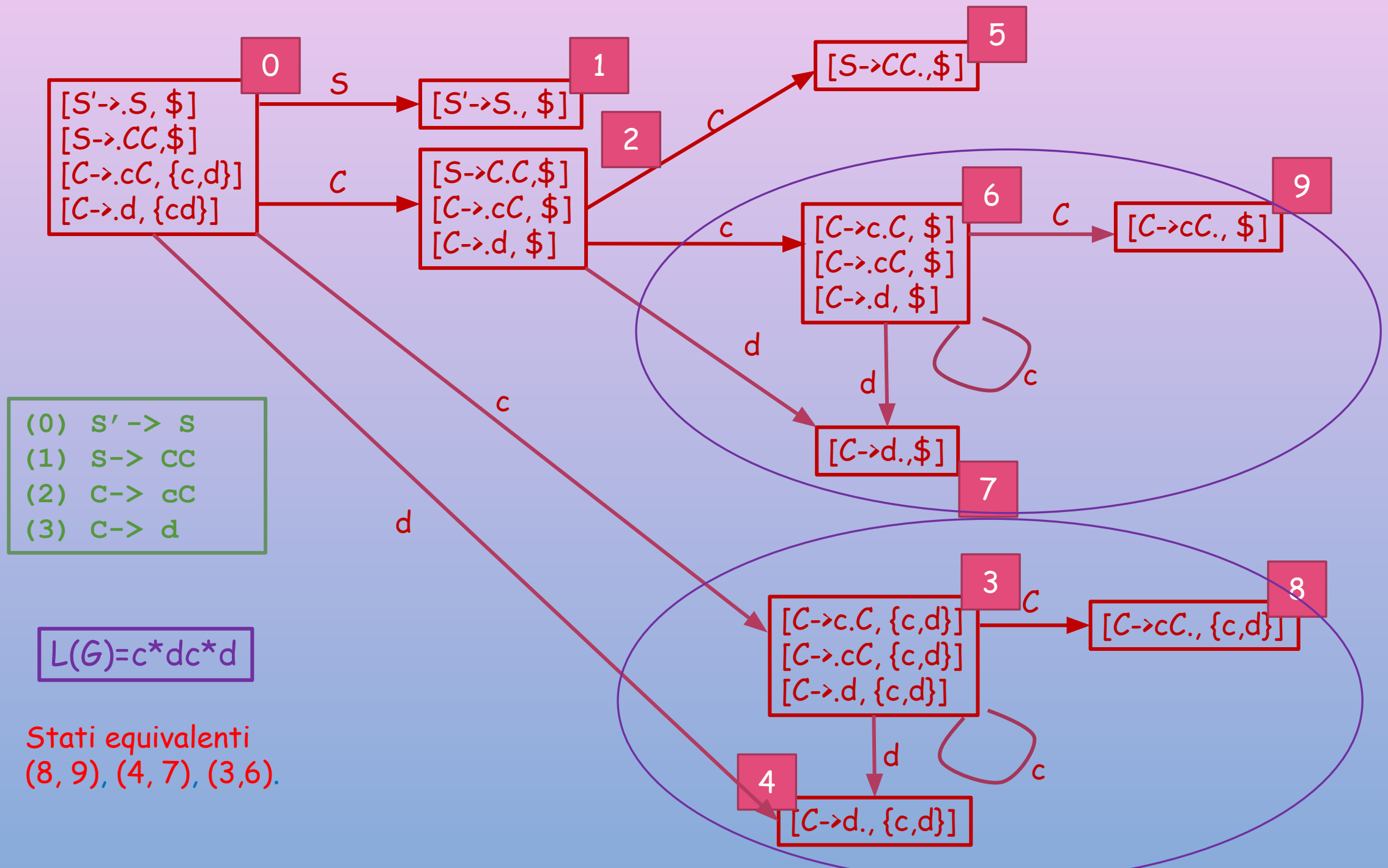
1. Le prime componenti degli item **LR(1)** sono item **LR(0)**
2. Se due stati  $p$  e  $q$  **LR(1)** hanno lo stesso core e se da  $p$  esce una transizione con  $X$  verso lo stato  $p'$ , allora anche da  $q$  uscirà una transizione con  $X$  verso  $q'$  e  $p'$  e  $q'$  avranno lo stesso core.

# METODO LALR

1. Costruire l'automa con item LR(1)
2. Si accorpano gli stati con core comune.
3. Tabella ACTION costruita come prima ma dal nuovo insieme di stati
4. Tabella GOTO:  $\text{goto}(I, X) = K$  se
  - $I$  e' l'unione di  $I_1, \dots, I_n$
  - $K$  e' l'unione degli stati che hanno lo stesso core di  $\text{Goto}(I_1, x)$
$$\text{Core}(\text{goto}(I_1, X)) = \text{core}(\text{goto}(I_2, X)) = \dots$$

Una grammatica e' LALR(1) se la tabella ottenuta non ha conflitti

Osservazione: Non si possono introdurre conflitti shift/reduce (altrimenti ci sarebbero stati anche prima dell'accorpamento degli stati). Si potrebbero invece introdurre conflitti reduce/reduce.





(0)  $S' \rightarrow S$

(1)  $S \rightarrow CC$

(2)  $C \rightarrow cC$

(3)  $C \rightarrow d$

Sulla tabella questo equivale ad accorpare righe

Stato	c	d	\$	S	C
0	Shift 3	Shift 4		Goto 1	Goto 2
1			Accept		
2	Shift 6	Shift 7			Goto 5
3	Shift 3	Shift 4			Goto 8
4	Reduce 3	Reduce 3			
5			Reduce 1		
6	Shift 6	Shift 7			Goto 9
7			Reduce 3		
8	Reduce 2	Reduce 2			
9			Reduce 2		

(0)  $S' \rightarrow S$

(1)  $S \rightarrow CC$

(2)  $C \rightarrow cC$

(3)  $C \rightarrow d$

Stato	c	d	\$	S	C
0	Shift 3	Shift 4		Goto 1	Goto 2
1			Accept		
2	Shift 6	Shift 7			Goto 5
3	Shift 3	Shift 4			Goto 8
4	Reduce 3	Reduce 3	Reduce 3		
5			Reduce 1		
8	Reduce 2	Reduce 2	Reduce 2		

# CONDIZIONI LALR

Nota che una grammatica soddisfa la condizione **LALR(1)** se valgono entrambe le condizioni:

1. Ogni candidata di riduzione ha un insieme di prospezione disgiunto dalle etichette terminali uscenti;
2. Se vi sono due candidate di riduzione i loro insiemi di prospezione sono disgiunti;

# GAMMATICA LR(1)

Sia data la grammatica

$S \rightarrow aXb$

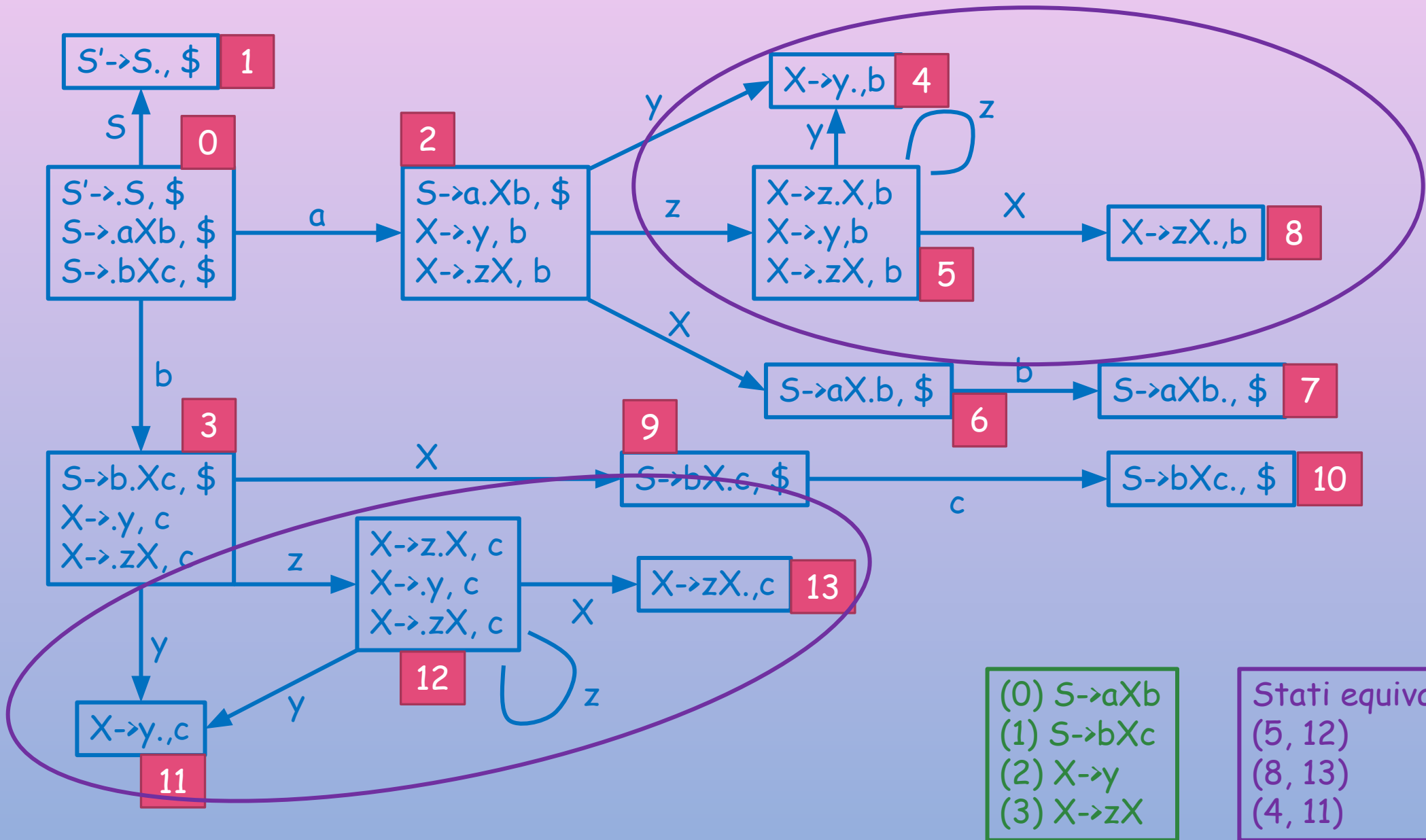
$S \rightarrow bXc$

$X \rightarrow y$

$X \rightarrow zX$

Costruire l'automa LR(1).

Sugg. Ha 14 stati e 3 coppie di stati possono essere fusi, dando luogo all'automa LALR(1) con 11 stati.



# Tabella LR(1)

	a	b	c	y	z	\$	S	X
0	Shift 2	Shift 3					Goto 1	
1						Accept		
2				Shift 4	Shift 5			Goto 6
3				Shift 11	Shift 12			Goto 13
4		Reduce X→y						
5				Shift 4	Shift 5			Goto 8
6		Shift 7						
7						Reduce S→aXb		
8		Reduce X→zX						
9			Shift 10					
10						Reduce S→bXc		
11			Reduce X→y					
12				Shift 11	Shift 12			
13			Reduce X→zX					

# Tabella LALR(1)

	a	b	c	y	z	\$	S	X
0	Shift 2	Shift 3					Goto 1	
1						Accept		
2				Shift 4	Shift 5			Goto 6
3				Shift 11	Shift 12			Goto 13
4		Reduce $X \rightarrow y$	Reduce $X \rightarrow y$					
5				Shift 4	Shift 5			Goto 8
6		Shift 7						
7						Reduce $S \rightarrow aXb$		
8		Reduce $X \rightarrow zX$	Reduce $X \rightarrow zX$					
9			Shift 10					
10						Reduce $S \rightarrow bXc$		

# ESEMPIO LR(1) MA NON LALR(1)

$S \rightarrow A \mid Ba \mid bAa \mid bB$

$A \rightarrow a$

$B \rightarrow a$



Un parsing LALR(1) potrebbe generare conflitti che il LR(1) corrispondente non genererebbe (ciò non accade in pratica).

Si dimostra che se una grammatica è LR(1), la tabella LALR(1) non può avere conflitti shift/reduce ma solo reduce/reduce.

E' possibile computare il DFA del LALR(1) direttamente dal DFA del LR(0) attraverso un processo chiamato **lookahead propaganti**

# REGOLE PER RISOLVERE L'AMBIGUITÀ

Spesso può essere comodo usare **grammatiche ambigue** ed usare delle regole per risolvere l'ambiguità.

Per **esempio** regole di precedenza ed associatività. La grammatica

$E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$

(è ambigua poiché non specifica la precedenza e l'associatività tra gli operatori)

La grammatica non ambigua equivalente è:

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid \text{id}$

La prima è preferibile perché:

- Si possono cambiare precedenza e associatività senza cambiare le produzioni
- Il parser con meno produzioni è più veloce.

Stabilire delle regole di precedenza e di associatività fa risolvere al parser i conflitti.

Un altro esempio riguarda l'ambiguità del "dangling else".

# PROPRIETÀ DEI LINGUAGGI E DELLE GRAMMATICHE LR(K)

1. La famiglia dei linguaggi verificabili da parser deterministici coincide con quella dei linguaggi generati dalle grammatiche LR(1). Ciò non significa che ogni grammatica il cui linguaggio è deterministico, sia necessariamente LR(1): potrebbe richiedere una prospezione di lunghezza  $k > 1$ ; tuttavia esisterà una grammatica equivalente LR(1);
2. La famiglia dei linguaggi generati dalle grammatiche LR(k) coincide con quella dei linguaggi generati da LR(1). Quindi un linguaggio context free ma non-deterministico non può avere una grammatica LR(k) per nessun k.
3. Per ogni  $k \geq 1$ , esistono grammatiche che sono LR(k) ma non LR(k-1).
4. Data una grammatica, è indecidibile se esista un  $k > 0$  per cui tale grammatica risulti LR(k); di conseguenza non è decidibile se il linguaggio generato da una grammatica CF è deterministico. E' decidibile soltanto se k è fissato (si applica la costruzione del parser).

# CONSIDERAZIONI FINALI SU LINGUAGGI E GRAMMATICHE $LL(k)$ E $LR(k)$ :

1. Ogni linguaggio regolare è  $LL(1)$ ;
2. Ogni linguaggio  $LL(k)$  è deterministico, ma vi sono linguaggi deterministici per cui non esiste alcuna grammatica  $LL(k)$ ;
3. Per ogni  $k \geq 0$ , una grammatica  $LL(k)$  è anche  $LR(k)$ ;
4. Le grammatiche  $LL(1)$  e  $LR(0)$  non sono incluse una nell'altra;
5. Quasi tutte le grammatiche  $LL(1)$  sono  $LALR(1)$

# DOMANDE A RISPOSTA APERTA

- Cosa si intende per front end e back end di un compilatore?
- Qual è il ruolo dell'analizzatore lessicale in un compilatore?

ESERCIZI PER CASA

# ESERCIZI

# ESERCIZI

Siano date la grammatica  $G$  e la stringa  $badc$ . Si applichi l'algoritmo di Earley per stabilire se la stringa appartiene al linguaggio generato da  $G$ .

$G$ :

$S \rightarrow Sc | bA$

$A \rightarrow aA | d$

La grammatica  $G$  non è  $LL(1)$ . Si può trasformare in una grammatica  $LL(1)$ ?

La grammatica è  $LR(0)$ ?

Effettuare il parsing in tutti i modi possibili.



# ESERCIZI

Sia  $G$  una grammatica avente  $X$  come assioma,  $\{a,b\}$  simboli terminali e le cui produzioni sono di seguito descritte.

$G$ :

$$X \rightarrow AA$$

$$A \rightarrow aA \mid b$$

Rispondere alle seguenti domande, motivando le risposte.

- . La grammatica  $g$  è LR(0)?
- . La grammatica è LR(1)?

# ESERCIZI

Sia  $G$  una grammatica avente  $S$  come assioma,  $\{a,b,c,d\}$  simboli terminali e le cui produzioni sono di seguito descritte.

$G$ :

$S \rightarrow TV$

$T \rightarrow a \mid b$

$V \rightarrow dcV \mid d$

La grammatica  $G$  è LR(0)?

# ESERCIZIO

Sia  $G$  una grammatica avente  $S$  come assioma,  $\{a,b,c,d\}$  simboli terminali e le cui produzioni sono di seguito descritte.

$G$ :

$S \rightarrow Bd \mid Dc \mid bBc \mid bDd$

$B \rightarrow a$

$D \rightarrow a$

Rispondere alle domande seguenti

1. Dopo aver costruito la tabella LR(1), dedurre se la grammatica è LR(1).
2. La grammatica è LALR(1)? Perché?

# ESERCIZIO

Sia data la grammatica  $G$  descritta di seguito. Si esegua il riconoscimento della stringa  $bcc$  utilizzando l'algoritmo di Earley. Qual è la complessità di tempo dell'algoritmo per questa grammatica?

$G$ :

$S \rightarrow T|X$

$T \rightarrow bTcc \mid bcc$

$X \rightarrow bXc \mid bc$

Qual è il linguaggio  $L(G)$  generato da  $G$ ?

. Date le caratteristiche della grammatica  $G$  e del linguaggio, è possibile rispondere alle seguenti domande:

1.  $G$  è LR(0)?
2.  $G$  è Lr(k) per qualche  $k$ ?

Nota: l'esercizio non richiede la costruzione di tabelle di parsing o di automi LR(0) o LR(1),

# ESERCIZI

Si consideri la grammatica  $G$  che ha il simbolo  $S$  come assioma,  $\{a,b,c,d,f\}$  simboli terminali e le seguenti regole di produzione:

$G: S \rightarrow CB$

$B \rightarrow b \mid f$

$C \rightarrow AB \mid DAB \mid \varepsilon$

$A \rightarrow Aa \mid a$

$D \rightarrow c \mid d$

Si risponda alle seguenti domande:

1. Per quale motivo per tale grammatica non è possibile effettuare un parsing di tipo LL(1)?  
 $G$  è LR(0)?
2. Si trasformi la grammatica  $G$  in una grammatica  $G'$  adatta all'analisi sintattica discendente. Si costruisca la tabella di parsing LL(1) per tale grammatica. La grammatica  $G'$  è LL(1)?
3. Si descriva il funzionamento di un generico parser LL(1).

# ESERCIZI

Sia data la seguente grammatica  $G$ :

$S \rightarrow iAi \mid aAa \mid bAb \mid aBb \mid bBa$

$B \rightarrow i$

$A \rightarrow i$

Verificare se essa genera un linguaggio di tipo:

. LL(1)

. LR(0)

. LR(1)

. LALR(1)

# ESERCIZIO

Data la grammatica  $G$ :

$S \rightarrow a|ab$

$E'$  LR(0)?  $E'$  LR(1)?  $E'$  LALR(1)?

ESERCIZI CON FLEX



# ESERCIZIO CON FLEX

Con l'ausilio di flex scrivere un analizzatore lessicale che effettui la trasformazione di un testo scritto in un Linguaggio, chiamato linguaggio AA. Tale linguaggio è case-sensitive e le parole del linguaggio sono costituite solo da caratteri alfanumerici. Due parole possono essere separate l'una dall'altra da spazi, tabulazioni o newline. Ogni parola che contiene almeno un carattere diverso da quelli alfanumerici è considerata una parola non appartenente al lessico.

Le parole del lessico del linguaggio vengono trasformate come di seguito descritto:

. Ogni parola  $w$  costituita da lettere maiuscole che inizia per vocale viene trasformata in una parola ottenuta da  $w$  rimuovendone le occorrenze di sequenze di consonanti consecutive di lunghezza dispari. Le sequenze di consonanti consecutive di lunghezza pari vengono trasformate in minuscolo. Per esempio la parola ABSFBERTPLDAWRRTDDORR viene trasformata in absfbeawrrtddorr.

. Ogni parola costituita da lettere minuscole e cifre e che iniziano per lettera minuscola viene trasformata in una parola in cui, se le occorrenze di cifre consecutive costituiscono un numero naturale pari, questo viene trasformato nella sequenza di cifre che rappresentano la metà del suddetto numero, altrimenti viene trasformato nel numero successivo.

Per esempio, la parola cd34rftsy567errer20r viene trasformata in cd17rftsy568errer10r.

. Tutte le altre parole del lessico vengono trasformate nella parola ottenuta ripetendo il primo simbolo della parola tante volte quanti sono i simboli della parola stessa. Per esempio, la parola oeriiiiodivjk34hhfjj4k5 viene trasformata in 0000000000000000000000000000.

