

# A New Look at LR(k)

Bill McKeeman,

MathWorks Fellow

for

Worcester Polytechnic

Computer Science Colloquium

February 20, 2004

## **Abstract**

Knuth offered LR(k) as an algorithm for parsing computer languages (*Information and Control*, Nov. 1965). Most work since has concentrated on parser usability and space/time optimization.

This talk presents LR(k) from first principles, teasing apart some complex formalisms applied by Knuth into simpler components.

## **Assumptions**

The audience is comfortable with mathematical notation and computer programming.

## Points to Ponder

- A context-free grammar can describe meaning in addition to syntax.
- The set of reducible parse stacks is a regular language.
- Any finite automata can be described by a context-free grammar.
- The LR(0) machine can be constructed in two stages, first an NFA, then a DFA.
- In most cases, the LALR(1) lookahead can be teased out of the LR(0) machine.

A context-free grammar  $\langle \mathbf{N}, \mathbf{T}, \mathbf{G}, \mathbf{R} \rangle$  is a four-tuple describing a language  $\mathbf{L}$ .

$\mathbf{N}$  is a set of non-terminal symbols (naming syntactic structure like *expression*).

$\mathbf{T}$  is a set of terminal symbols (naming the written symbols like *+*, *for*, *int...*)

$$\mathbf{V} = \mathbf{N} \cup \mathbf{T}, \quad \mathbf{N} \cap \mathbf{T} = \{ \}.$$

$\mathbf{G} \subseteq \mathbf{N}$ , is a set of symbols describing the language (usually just *sentence* or *translation-unit*).

$\mathbf{R}$  is a set of rewriting rules,  $\mathbf{R} \subseteq \mathbf{N} \times \mathbf{V}^*$ .

Given an input text from language  $\mathbf{L}$ , we repeatedly substitute the LHS of a rewriting rule for its RHS, leading eventually to something in  $\mathbf{G}$ . The sequence of substitutions is called the *parse*.

# Why Use Context-free Grammars?

$\mathbf{N} = \{P \ D \ C \ B\}$ ,  $\mathbf{T} = \{\perp \mid \ \& \ t \ f\}$ ,  $\mathbf{G} = \{P\}$ .

<i>7 rewriting rules:</i>	<i>rewritings</i>	<i>rule</i>	<i>consumed terminals</i>
1. $P \rightarrow D \perp$	$\underline{t} \& t \mid f \& f \perp$	6	t
2. $D \rightarrow D \mid C$	$\underline{B} \& t \mid f \& f \perp$	5	
3. $D \rightarrow C$	$C \& \underline{t} \mid f \& f \perp$	6	t
4. $C \rightarrow C \& B$	$\underline{C \& B} \mid f \& f \perp$	4	&
5. $C \rightarrow B$	$\underline{C} \mid f \& f \perp$	3	
6. $B \rightarrow t$	$D \mid \underline{f} \& f \perp$	7	f
7. $B \rightarrow f$	$D \mid \underline{B} \& f \perp$	5	
	$D \mid C \& \underline{f} \perp$	7	f
	$D \mid \underline{C \& B} \perp$	4	&
	$\underline{D} \mid \underline{C} \perp$	2	
	$\underline{D} \perp$	1	$\perp$
	P		

*Grammar for boolean expressions*

*pcode*  
tt&ff&| $\perp$

## How to Get a Parse?

- Recursive Descent (top down)
  - Write a C function for each nonterminal
  - Report rule applications
  - Pro: no tools needed, Con: can be buggy
- YACC (bottom up)
  - Write a CFG
  - Implement rule applications
  - Pro: reliable, Con: deal with YACC diagnostics

<i>Parse Stack</i>	<i>Input</i>	<i>S/R sequence</i>
•	t&t   f&f⊥	<i>start</i>
• t	&t   f&f⊥	<i>shift t</i>
• B	&t   f&f⊥	<i>rule 6</i>
• C	&t   f&f⊥	<i>rule 5</i>
• C&	t   f&f⊥	<i>shift &amp;</i>
• C&t	f&f⊥	<i>shift t</i>
• C&B	f&f⊥	<i>rule 6</i>
• C	f&f⊥	<i>rule 4</i>
• D	f&f⊥	<i>rule 3</i>
• D	f&f⊥	<i>shift  </i>
• D   f	&f⊥	<i>shift f</i>
• D   B	&f⊥	<i>rule 7</i>
• D   C	&f⊥	<i>rule 5</i>
• D   C&	f⊥	<i>shift &amp;</i>
• D   C&f	⊥	<i>shift f</i>
• D   C&B	⊥	<i>rule 5</i>
• D   C	⊥	<i>rule 4</i>
• D	⊥	<i>rule 2</i>
• D⊥		<i>shift ⊥</i>
• P		<i>rule 1 (quit)</i>

The **red** parse  
stacks are  
*reducible*

S/R is useful for  
languages  
where a lexer is  
needed.

The set  $L$  of all reducible parse stacks for  $L$  is a regular language (recognizable with a finite automaton).  $L$  has a grammar that can be derived from the CFG of the original language  $L$ . The terminal symbols of  $L$  are the things that can appear on parse stacks (that is, all of  $T$  and  $N$ ). The nonterminal symbols of  $L$  are all the partial rules of  $L$  (that is, leave off zero or more symbols on the end of a rule of  $L$  and get a nonterminal of  $L$ ). The resulting rules  $R$  are all of the forms

$A \rightarrow Bb$

$A \rightarrow B$

$A \rightarrow$

Using finite automata terminology, rules  $A \rightarrow Bb$  are read as “in state  $B$  if you see  $b$ , eat it and go to state  $A$ ”. Rules  $A \rightarrow B$  are the same but no  $b$ . Rule  $A \rightarrow$  says “create  $A$  from nothing”, which is how we get started.

This definition leaves undefined the situation if an unacceptable symbol appears in the input. One can make the NFA 1-1 by adding an error state  $E$  and rules of the form  $E \rightarrow Bb$  to fill out the tables. In practice the error state triggers a syntax diagnostic. It is too verbose for use here.



# LR(0) **NFA** Construction

*notation:  $\alpha, \beta \in V^*$*

**L** defined by  $\langle \mathbf{N}, \mathbf{T}, \mathbf{G}, \mathbf{R} \rangle$

**L** defined by  $\langle \mathbf{N}, \mathbf{T}, \mathbf{G}, \mathbf{R} \rangle$

**T** = **V** *// anything on the parse stack*

**N** =  $\{A \rightarrow \alpha \mid A \rightarrow \alpha \beta \in \mathbf{R}\}$  *// any partial rule*

**G** = **R** *// any completed rule*

**R** =  $\{A \rightarrow \alpha a \rightarrow A \rightarrow \alpha a \mid A \rightarrow \alpha a \beta \in \mathbf{R}\}$

$\cup \{B \rightarrow \rightarrow A \rightarrow \alpha \mid A \rightarrow \alpha B \beta \in \mathbf{R} \wedge A \in \mathbf{N}\}$

$\cup \{A \rightarrow \rightarrow \mid A \in \mathbf{G}\}$

**R** for **L**

$P \rightarrow D \perp$   
 $D \rightarrow D \mid C$   
 $D \rightarrow C$   
 $C \rightarrow C \ \& \ B$   
 $C \rightarrow B$   
 $B \rightarrow t$   
 $B \rightarrow f$

**N** for **L**

0  $P \rightarrow$   
1  $P \rightarrow D$   
2  $P \rightarrow D \perp$   
3  $D \rightarrow$   
4  $D \rightarrow D$   
5  $D \rightarrow D \mid$   
6  $D \rightarrow D \mid C$   
7  $D \rightarrow C$   
8  $C \rightarrow$   
9  $C \rightarrow C$   
10  $C \rightarrow C \&$   
11  $C \rightarrow C \& B$   
12  $C \rightarrow B$   
13  $B \rightarrow$   
14  $B \rightarrow t$   
15  $B \rightarrow f$

**R** for **L**

0  $\rightarrow$   
1  $\rightarrow 0D$   
2  $\rightarrow 1 \perp$   
4  $\rightarrow 3D$   
5  $\rightarrow 4 \mid$   
6  $\rightarrow 5C$   
7  $\rightarrow 3C$   
9  $\rightarrow 8C$   
10  $\rightarrow 9 \&$   
11  $\rightarrow 10B$   
12  $\rightarrow 8B$   
14  $\rightarrow 13t$   
15  $\rightarrow 13f$   
3  $\rightarrow 0$   
8  $\rightarrow 3$   
8  $\rightarrow 5$   
13  $\rightarrow 8$   
13  $\rightarrow 10$

**G** for **L**

2 6 7 11  
12 14 15

Apply LR(0) NFA

$D \mid C \& B$   
0  $D \mid C \& B$   
3  $D \mid C \& B$   
4  $\mid C \& B$   
5  $C \& B$   
8  $C \& B$   
9  $\& B$   
10  $B$   
11

For rule 11 we rewrite  
 $C \rightarrow C \& B$  and start over:

$D \mid C$

*elements of **N**  
renamed for  
readability*

## NFA to DFA

The NFA is not efficient. So we apply the (textbook) NFA-to-DFA transformation, getting yet another context-free grammar.

In this case the nonterminals are *sets* of partial rules, and the terminals are the same as the for the NFA (since they describe the same language).

## LR(0) DFA Construction

$E(M) = M \cup \{X \mid Y \in E(M) \wedge X \rightarrow Y \in R\}$  // tx on empty

$S(M, b) = E(\{X \mid Y \in M \wedge X \rightarrow Yb \in R\})$  // tx on  $b$  from  $M$

$0 = E(\{0\})$  // start state

$N = \{0\} \cup \{S(M, b) \mid M \in N \wedge b \in T\}$

$G = \{M \mid M \in N \wedge M \cap G \neq \{\}\}$  // contains a rule

$P = \{0 \rightarrow\} \cup \{N \rightarrow Mb \mid M \in N \wedge b \in T \wedge N = S(M, b)\}$

Note: this construction gives just the states reachable from 0.

The textbook formalism is simpler but computes useless states

# The LR(0) DFA

**T** for **L**

P D C B  $\perp$   
| & t f

**N** for **L**

0 1 2 3 4 5 6  
7 8 9 10 11  
12 13 14 15

**N** renamed for  
readability

**T** for **L**

P D C B  $\perp$   
| & t f

**N** for **L**

0	0 3 8 13
1	1 4
2	2
3	5 8 13
4	6 9
5	10 13
6	11
7	14
8	15
9	7 9
10	12

**R** for **L**

0  $\rightarrow$   
1  $\rightarrow$  0D  
2  $\rightarrow$  1 $\perp$   
3  $\rightarrow$  1|  
4  $\rightarrow$  3C  
5  $\rightarrow$  4&  
6  $\rightarrow$  5B  
7  $\rightarrow$  5t  
8  $\rightarrow$  5f  
10  $\rightarrow$  3B  
7  $\rightarrow$  3t  
8  $\rightarrow$  3f  
9  $\rightarrow$  0C  
10  $\rightarrow$  0B  
7  $\rightarrow$  0t  
8  $\rightarrow$  0f

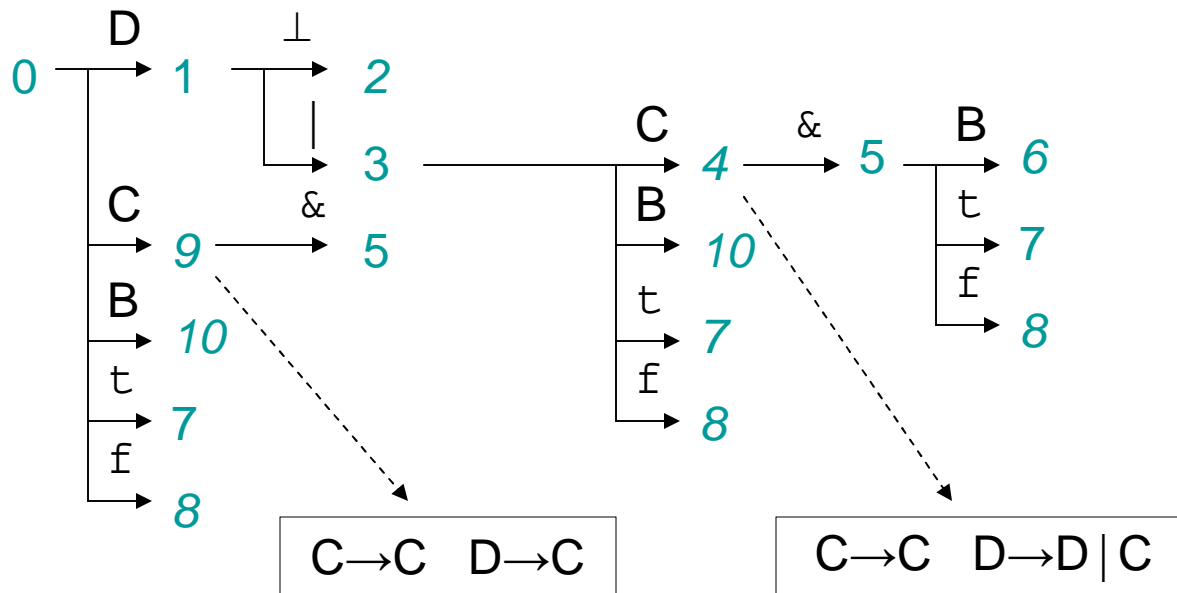
**G** for **L**

2 4 6 7  
8 9 10

Apply LR(0) **DFA**

D | C&B  
0D | C&B  
1 | C&B  
3C&B  
4&B  
5B  
6

# The LR(0) DFA



Apply LR(0) DFA

D | C&B  
 0 D | C&B  
 1 | C&B  
 3 C&B  
 4 &B  
 5 B  
 6

States 2, 6, 7, 8, 10 are obviously in **G** (no place to go).

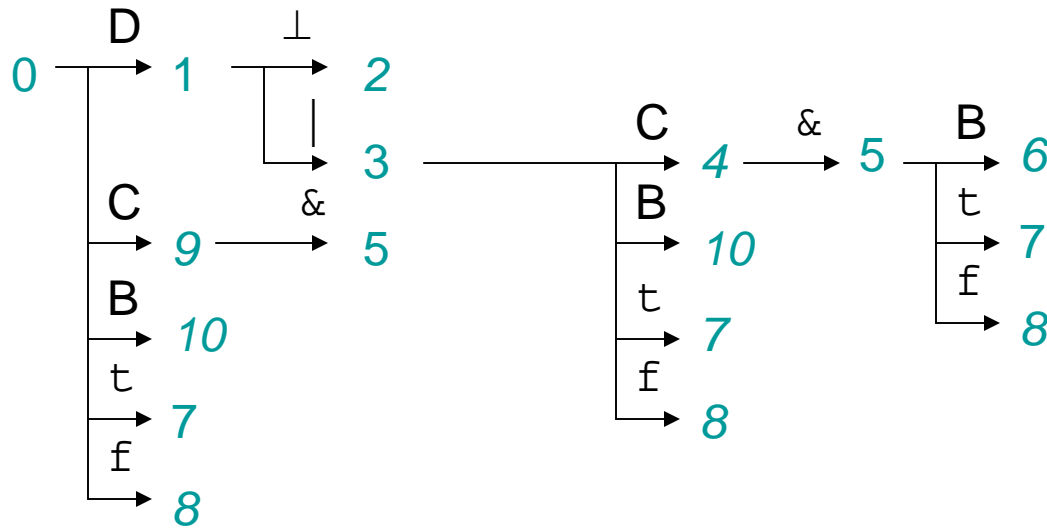
States 4 and 9 are in **G** but also can shift.

States 4 and 9 show LR(0) shift/reduce conflicts.

For this CFG, shift if you can resolves the conflicts.

Note that is what we did for state 4 above.

# The LALR(1) Lookahead



The LALR(1) lookahead for  $D \rightarrow D \mid C$  in state **4** is found by stepping back from **4**, over the RHS of the rule ( $D \mid C$ ), then stepping forward (from **0** in this case) over the LHS of the rule ( $D$ ) and collecting the transition symbols ( $\perp \mid$ ). Since the transition symbols out of **1** are different from the transition symbols out of **4** ( $\&$ ), there is no LALR(1) shift/reduce conflict.

If we erroneously did the reduction  $D \rightarrow D \mid C$  while looking at something other than  $\perp$  or  $\mid$ , we would immediately fail on the next step.

Conventional  
LALR(1)  
parsing table

see:

*in state:*

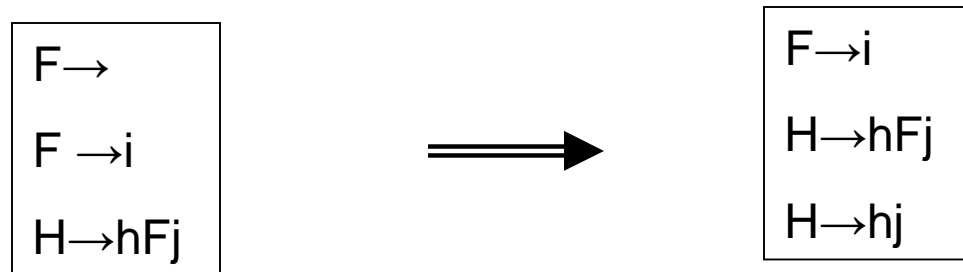
	P	D	C	B	⊥		&	t	f
0		1	9	10				7	8
1					2	3			
2	R1	R1	R1	R1	R1	R1	R1	R1	R1
3			4	10				7	8
4					R2	R2	5		
5				6				7	8
6					R4	R4	R4		
7					R6	R6	R6		
8					R7	R7	R7		
9					R3	R3	5		
10					R5	R5	R5		



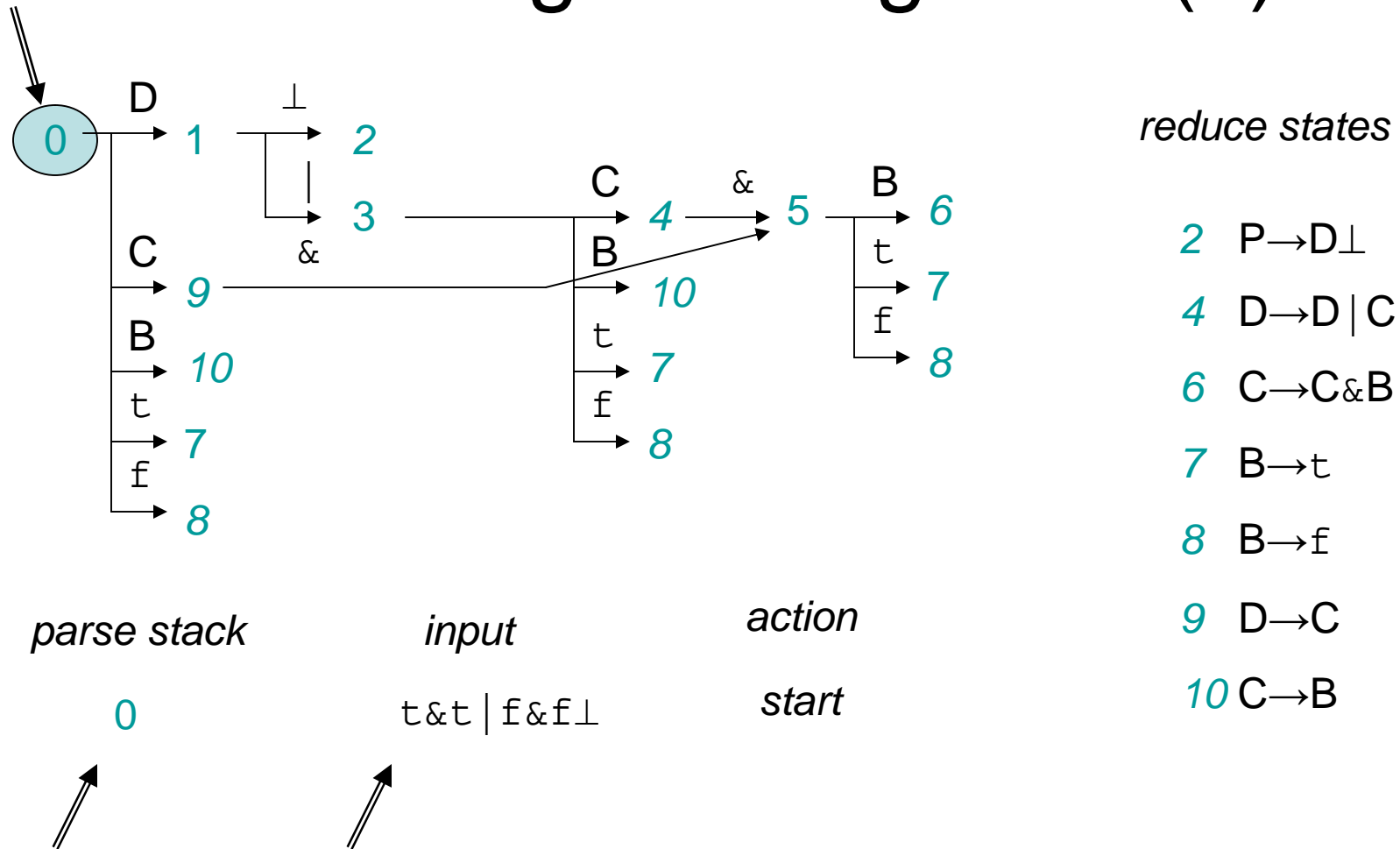
# NQLALR vs LALR

The simple trick for LALR(1) lookahead worked because the grammar for **L** had no erasing rules (rules with empty RHS). One way to insure the no-erasing property is to transform the user's grammar, removing the erasing rules.

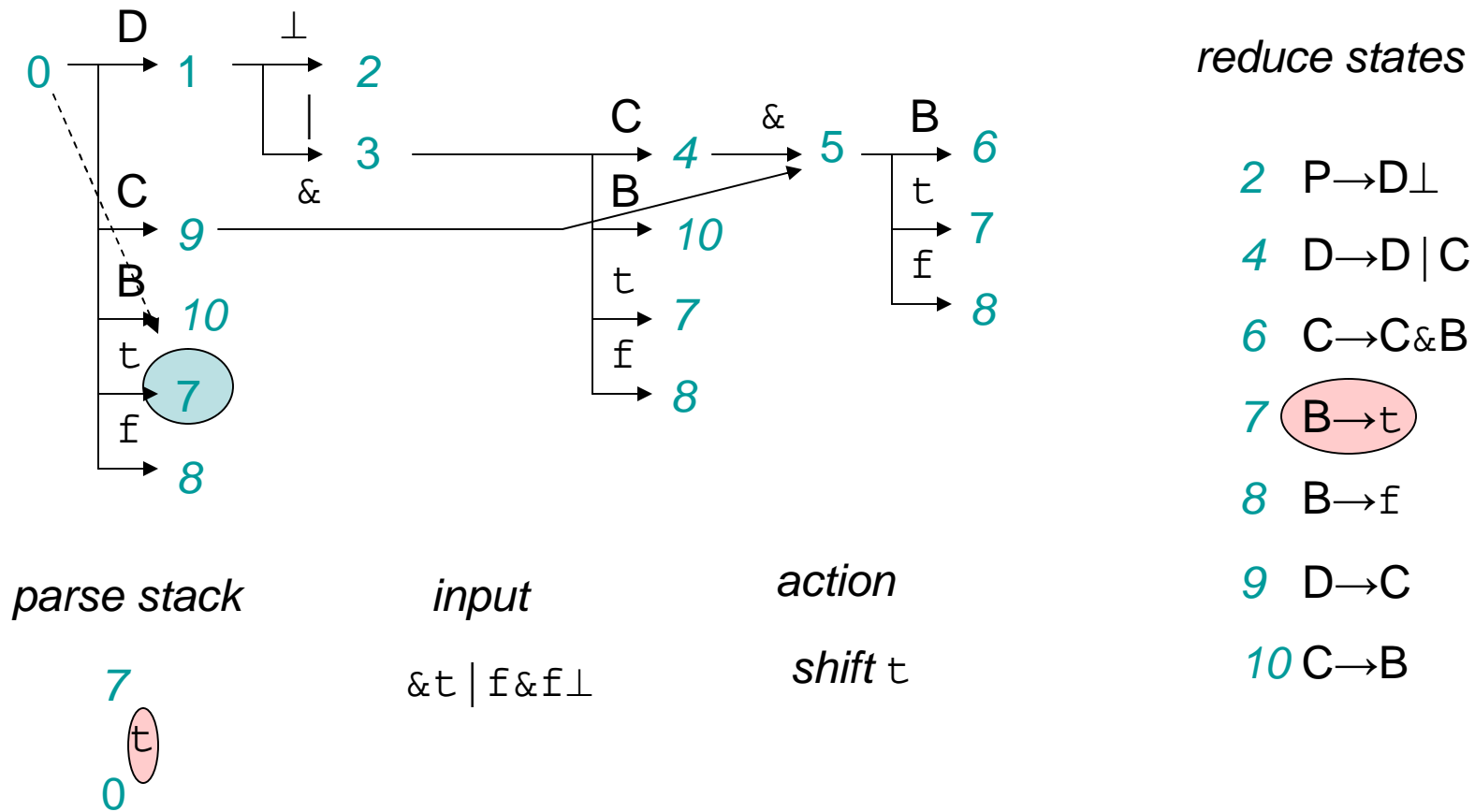
The trick is to remove each erasing rule, say  $F \rightarrow$ , and everywhere in the grammar that  $F$  is used, create a new rule with  $F$  erased.



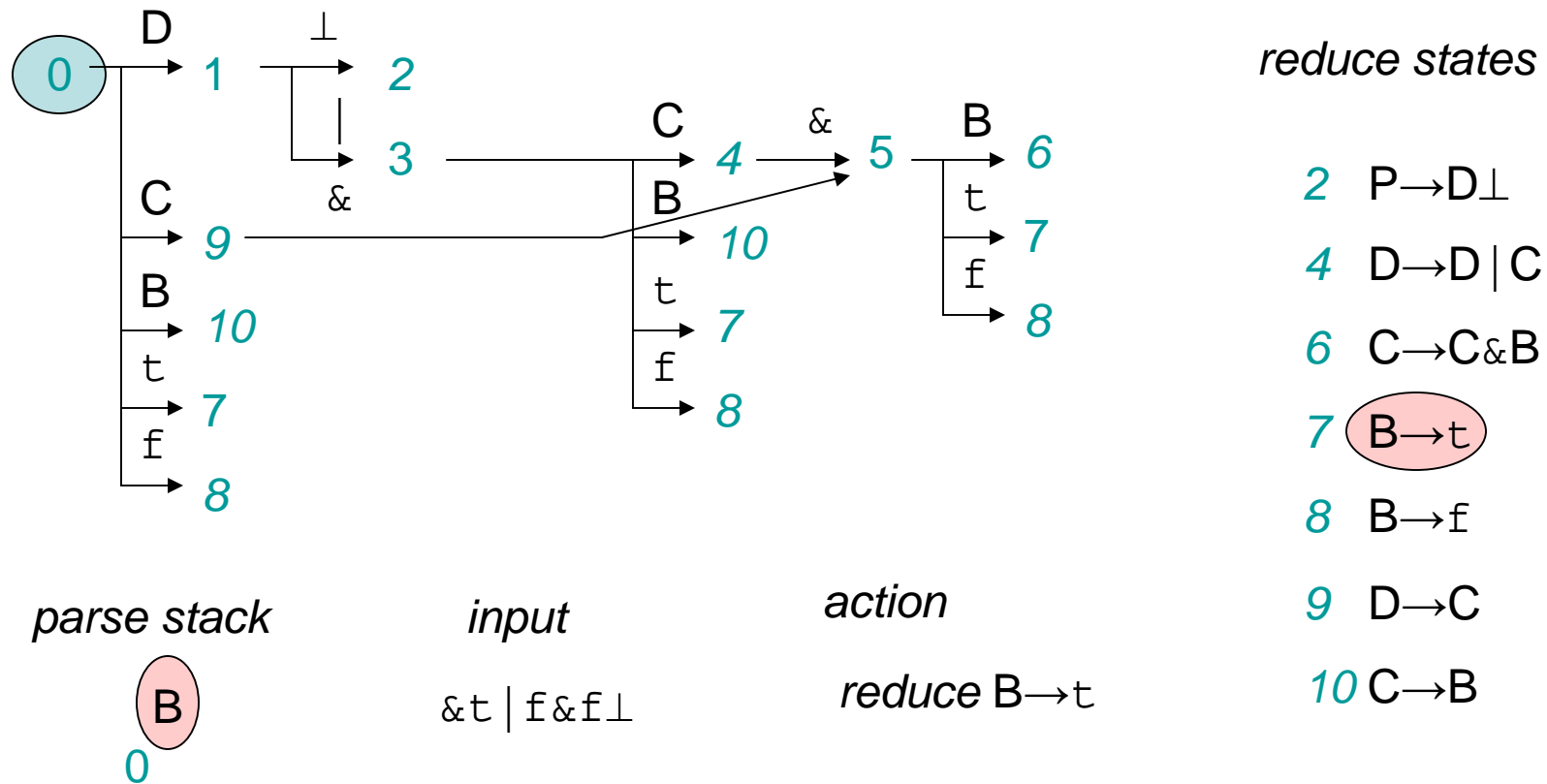
# Putting it all together (1)



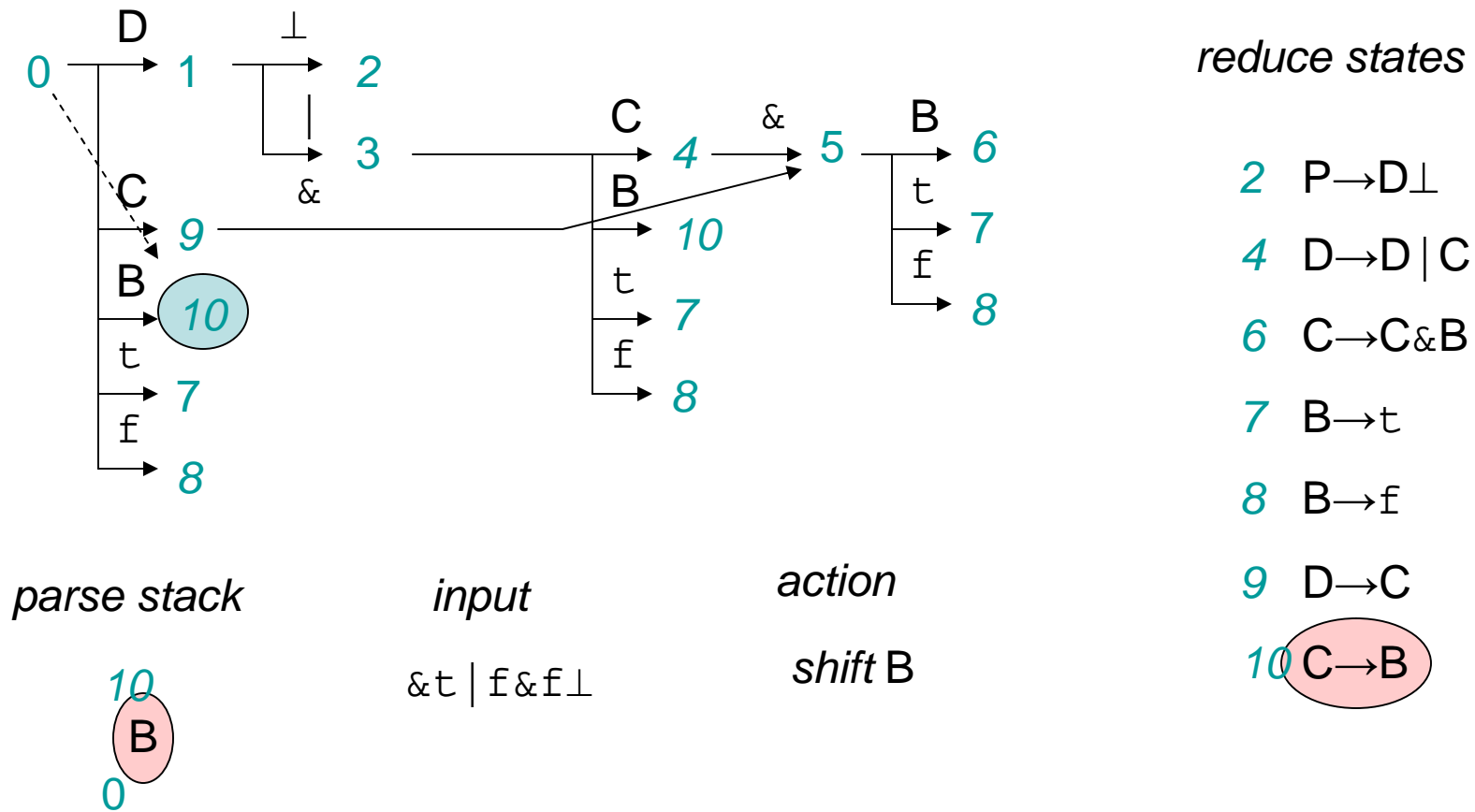
# Putting it all together (2)



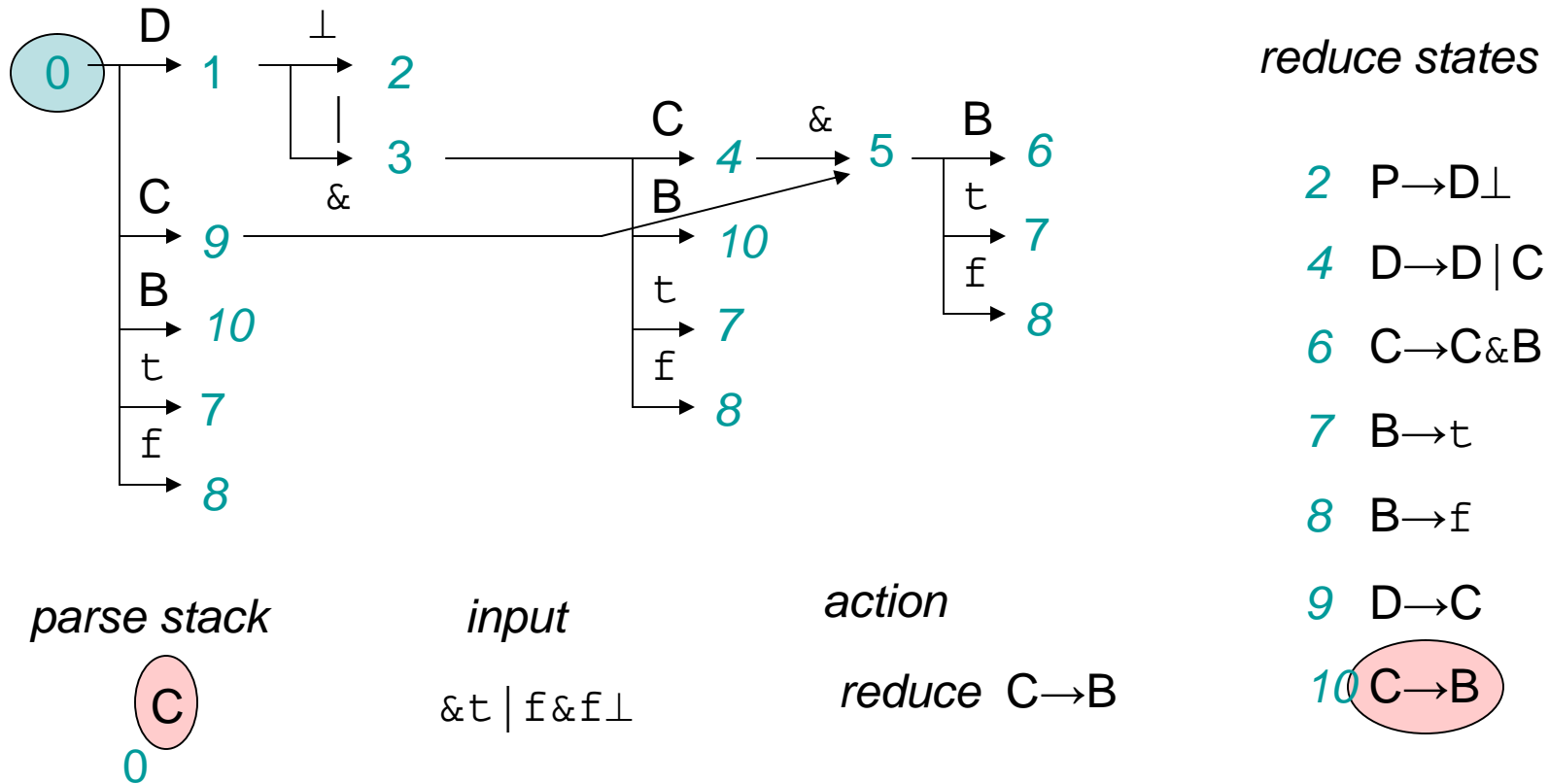
# Putting it all together (3)



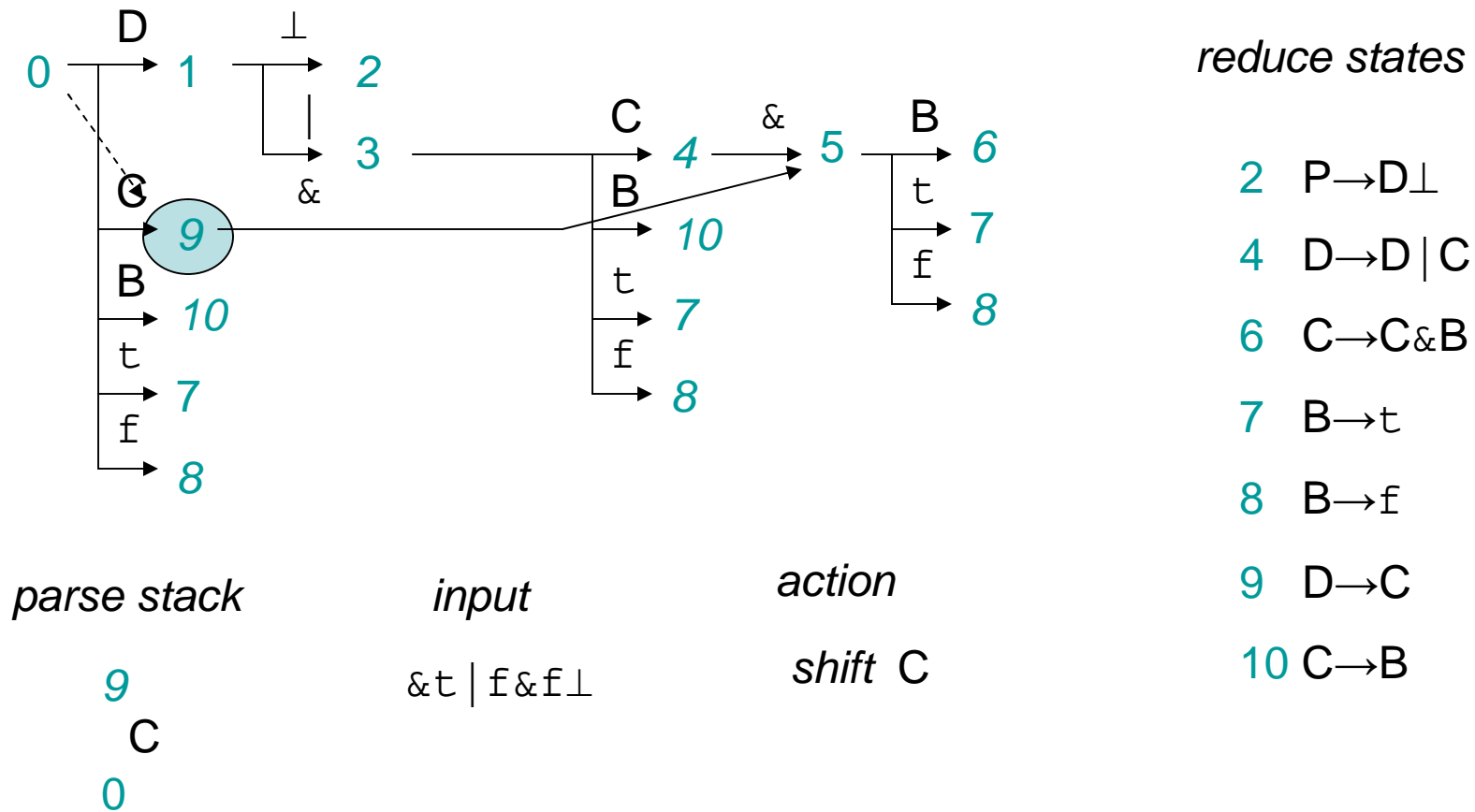
# Putting it all together (4)



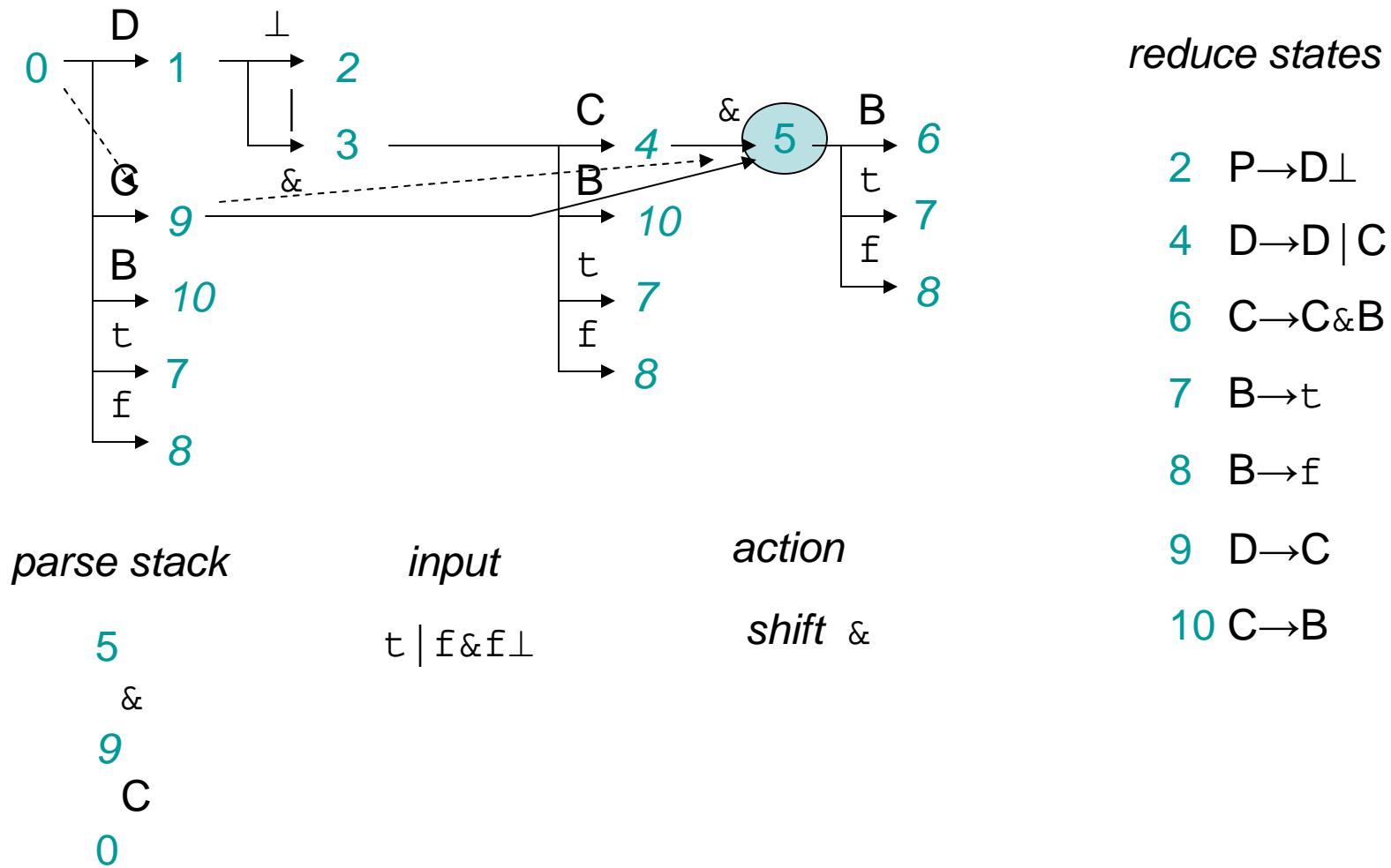
# Putting it all together (5)



# Putting it all together (6)

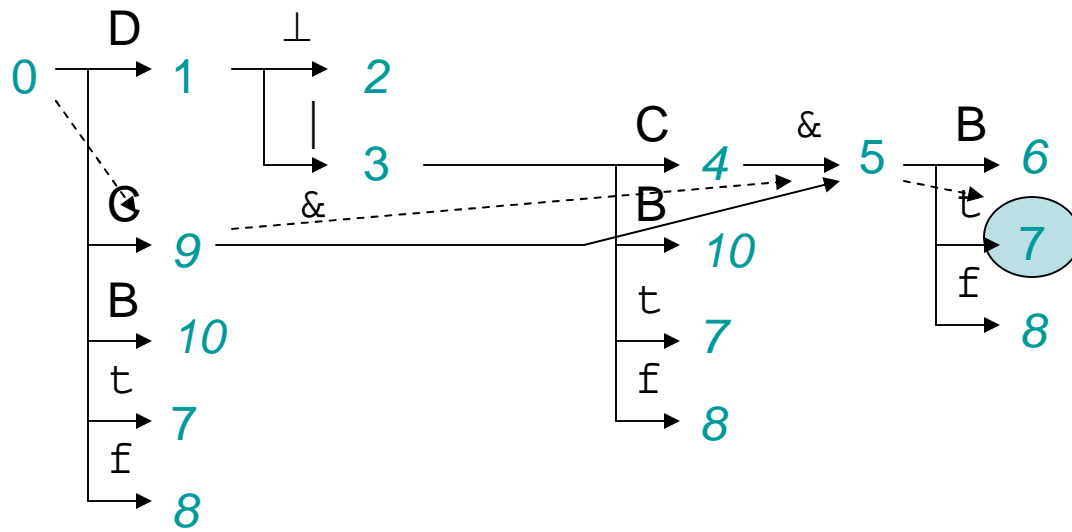


# Putting it all together (7)





# Putting it all together (8)



*reduce states*

- 2  $P \rightarrow D \perp$
- 4  $D \rightarrow D \mid C$
- 6  $C \rightarrow C \& B$
- 7  $B \rightarrow t$
- 8  $B \rightarrow f$
- 9  $D \rightarrow C$
- 10  $C \rightarrow B$

*parse stack*

7  
t  
5  
&  
9  
C  
0

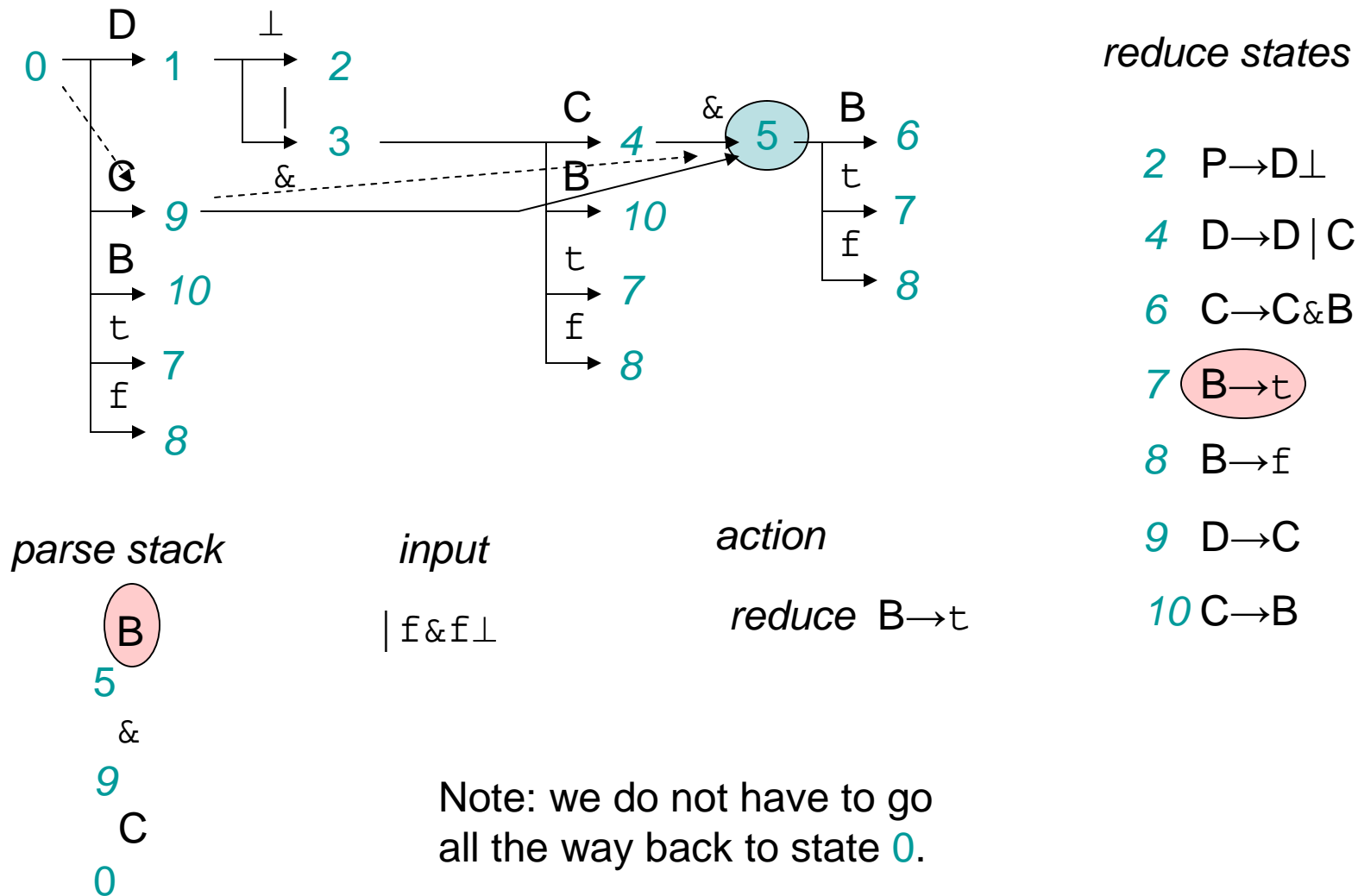
*input*

| f & f ⊥

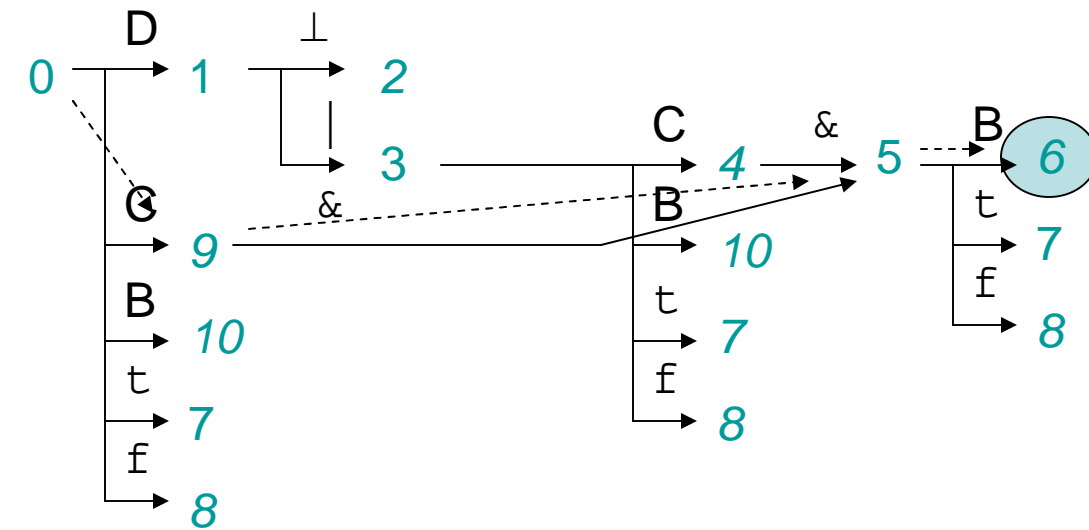
*action*

*shift* t

# Putting it all together (9)



# Putting it all together (10)



*reduce states*

2  $P \rightarrow D \perp$

4  $D \rightarrow D \mid C$

6  $C \rightarrow C \& B$

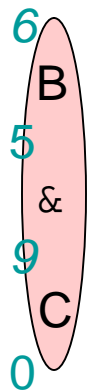
7  $B \rightarrow t$

8  $B \rightarrow f$

9  $D \rightarrow C$

10  $C \rightarrow B$

*parse stack*



*input*

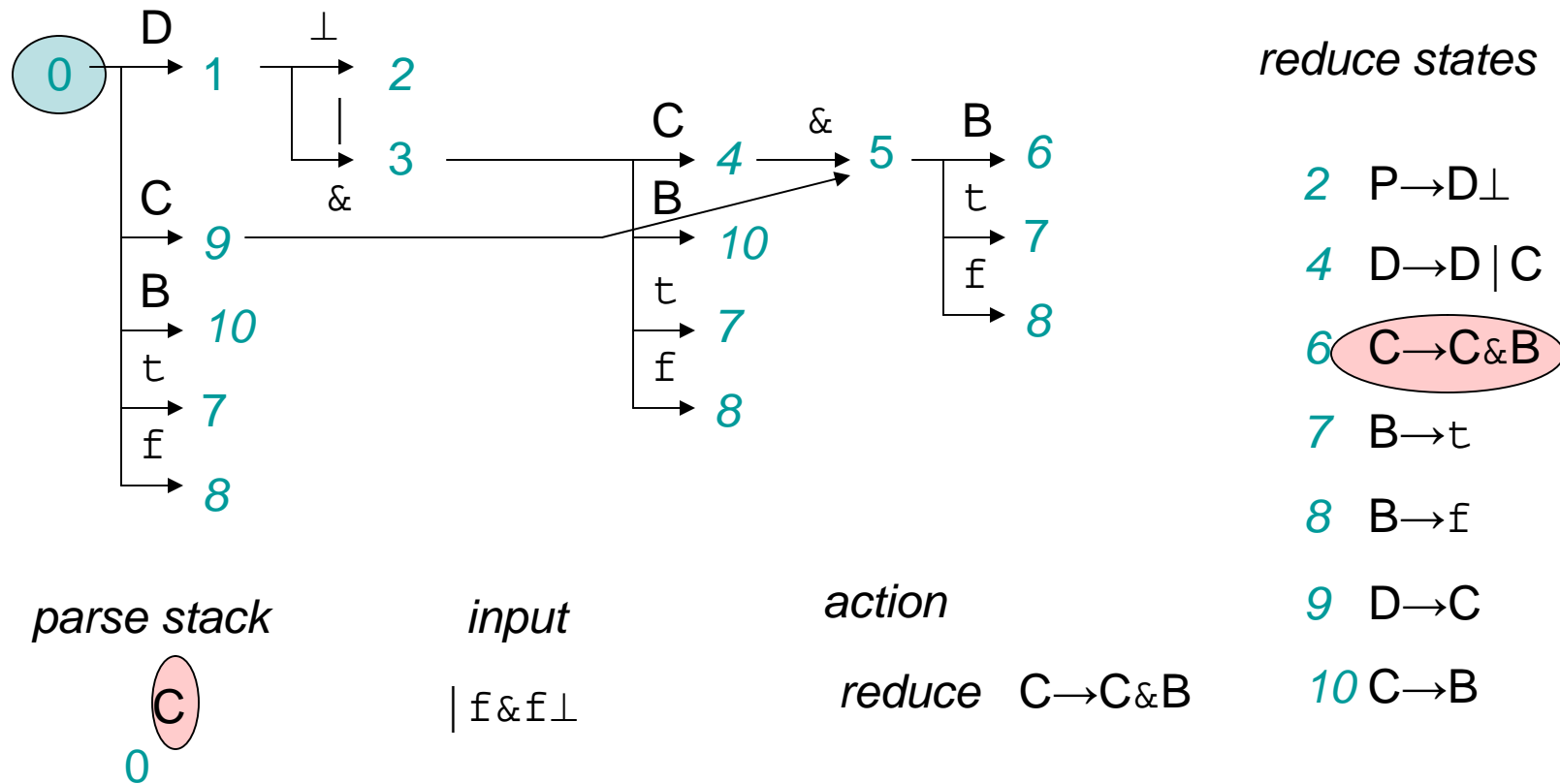
| f & f ⊥

*action*

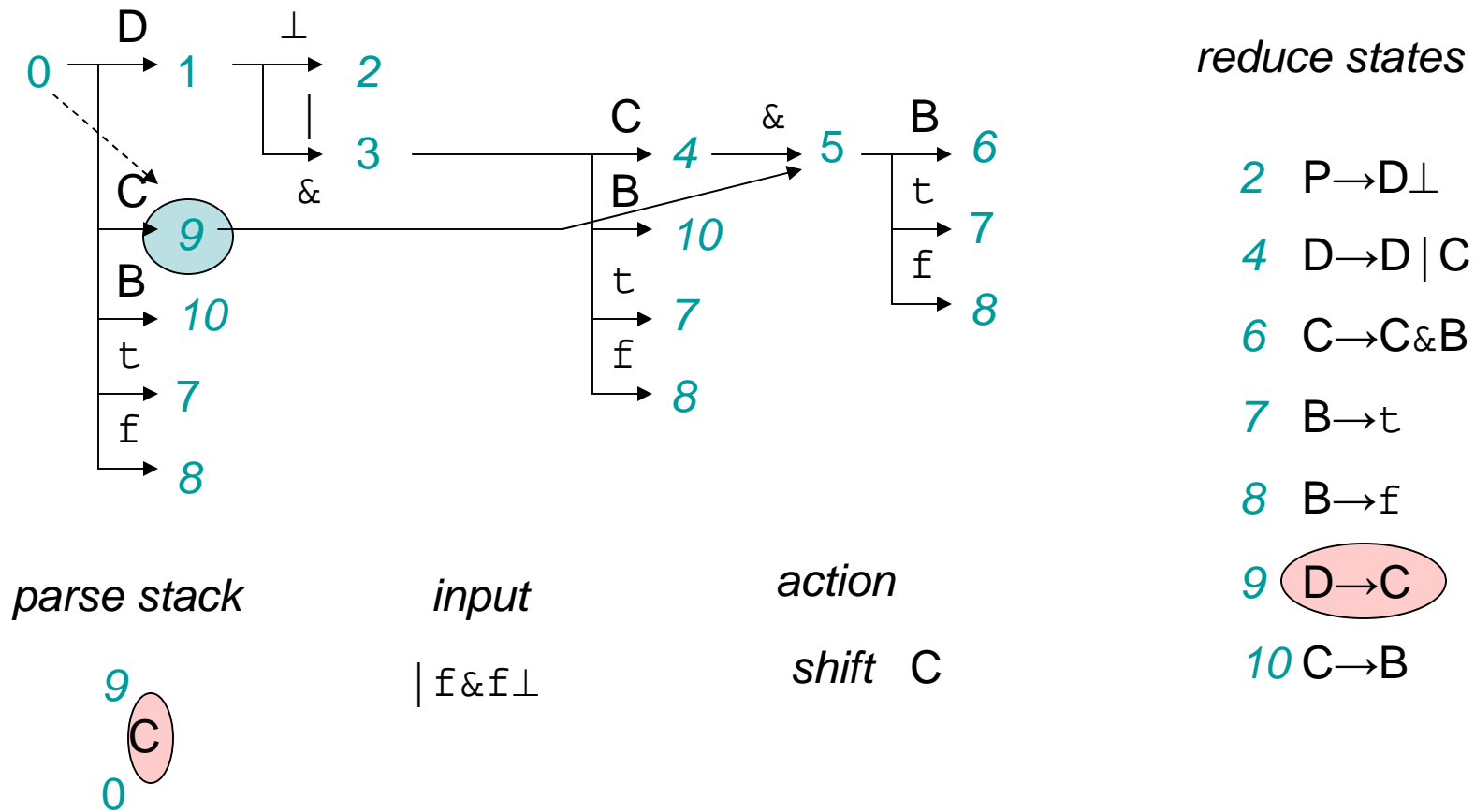
shift B

Note: we do not have to go all the way back to state 0.

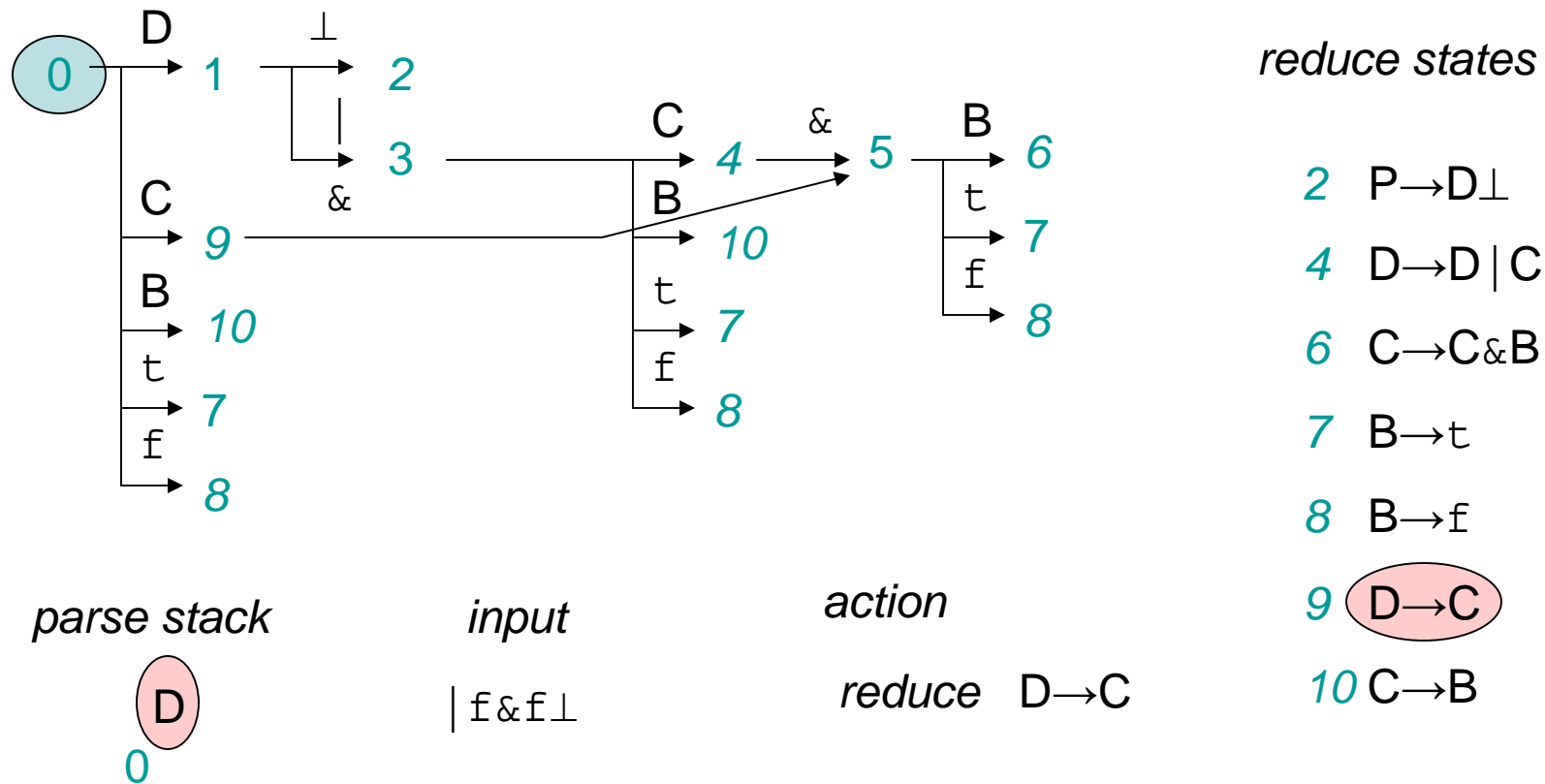
# Putting it all together (11)



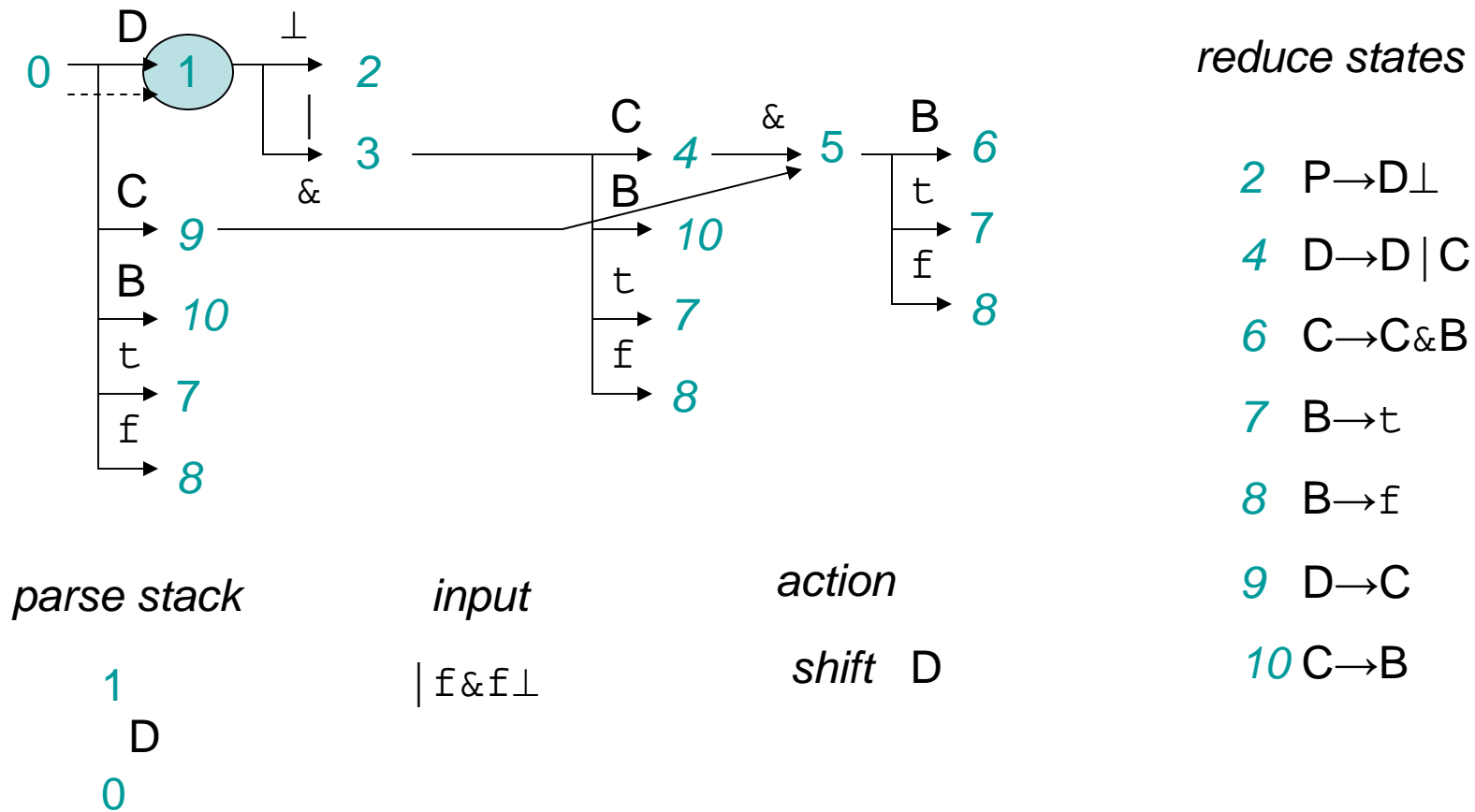
# Putting it all together (12)



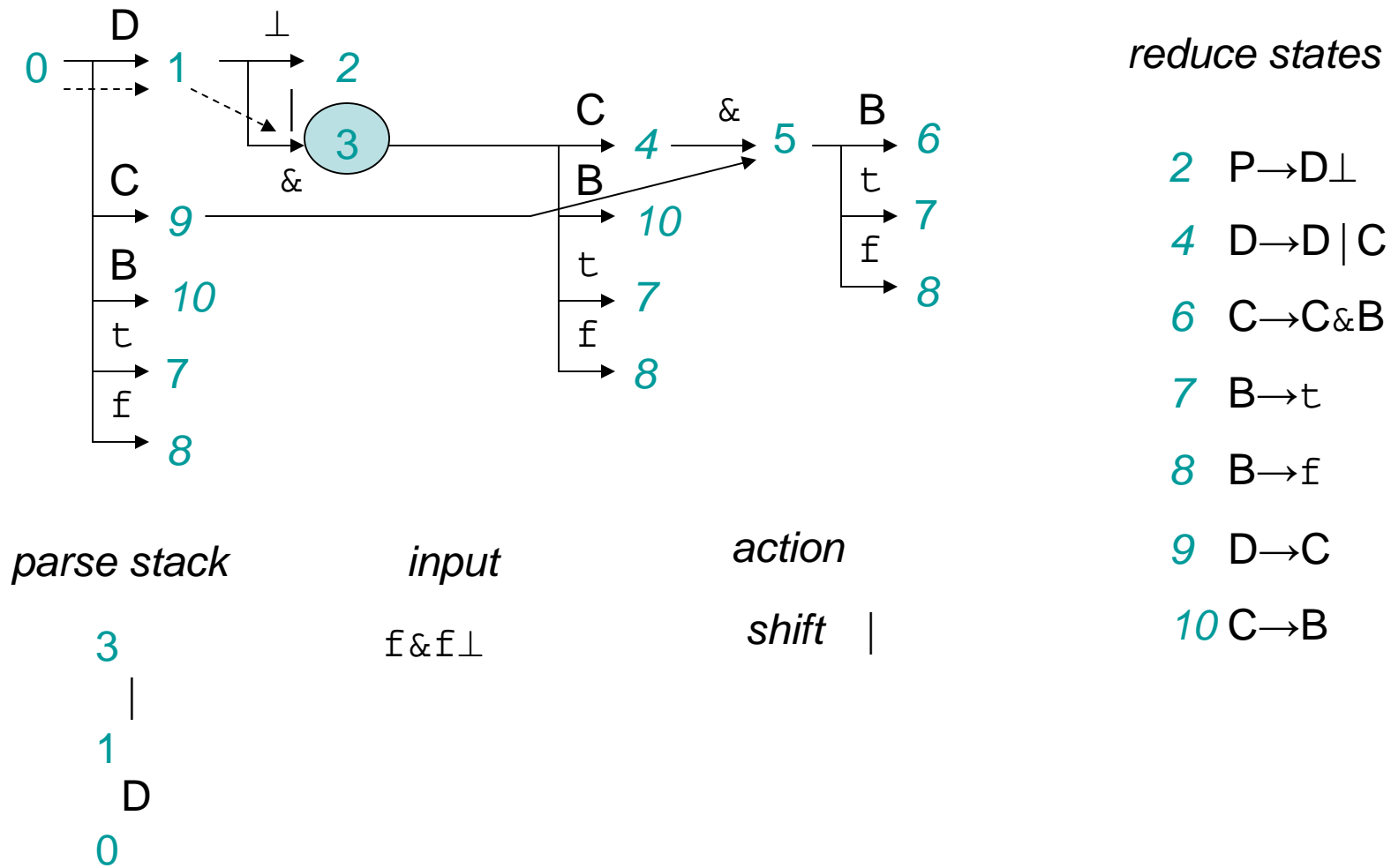
# Putting it all together (13)



# Putting it all together (14)

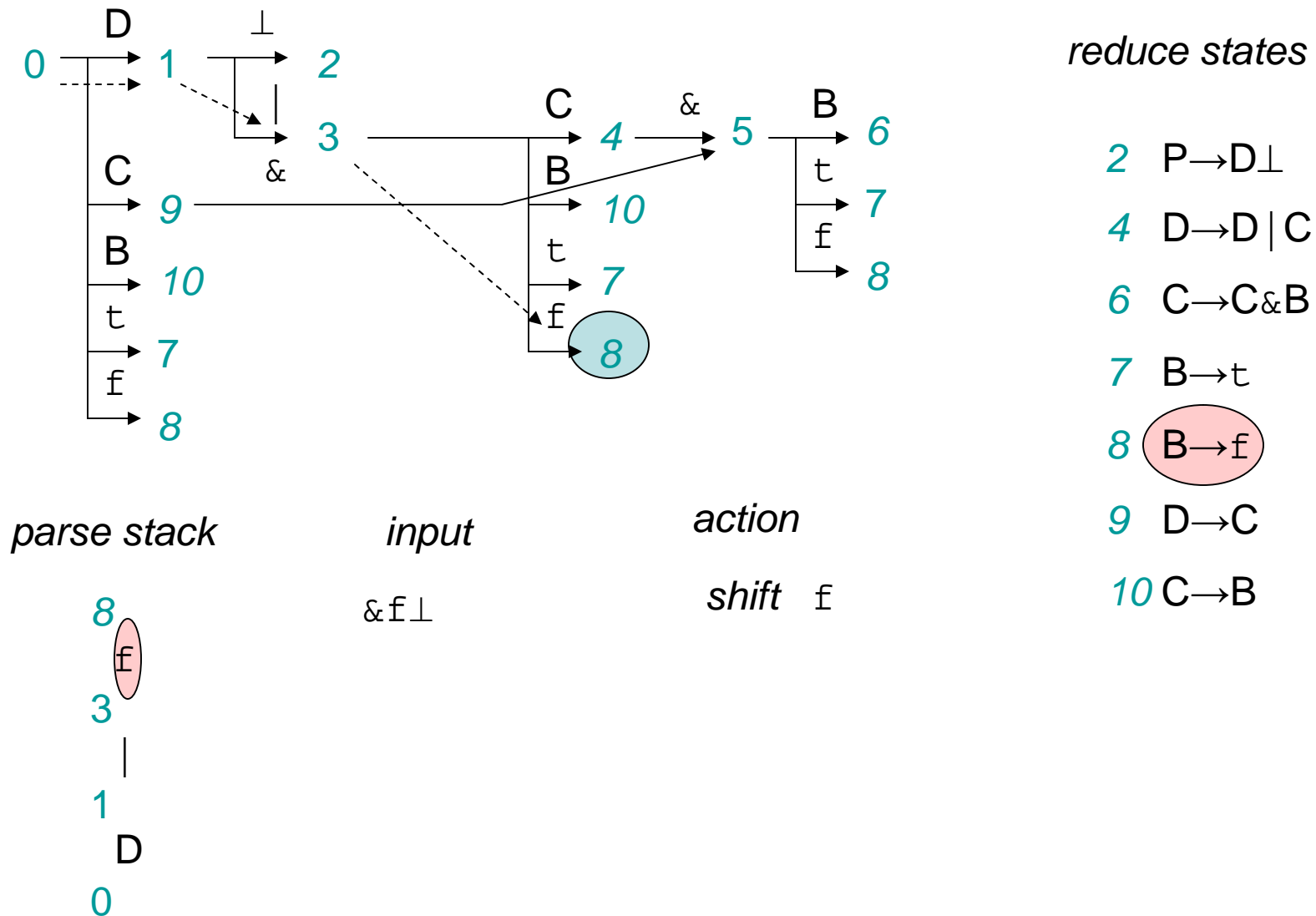


# Putting it all together (15)

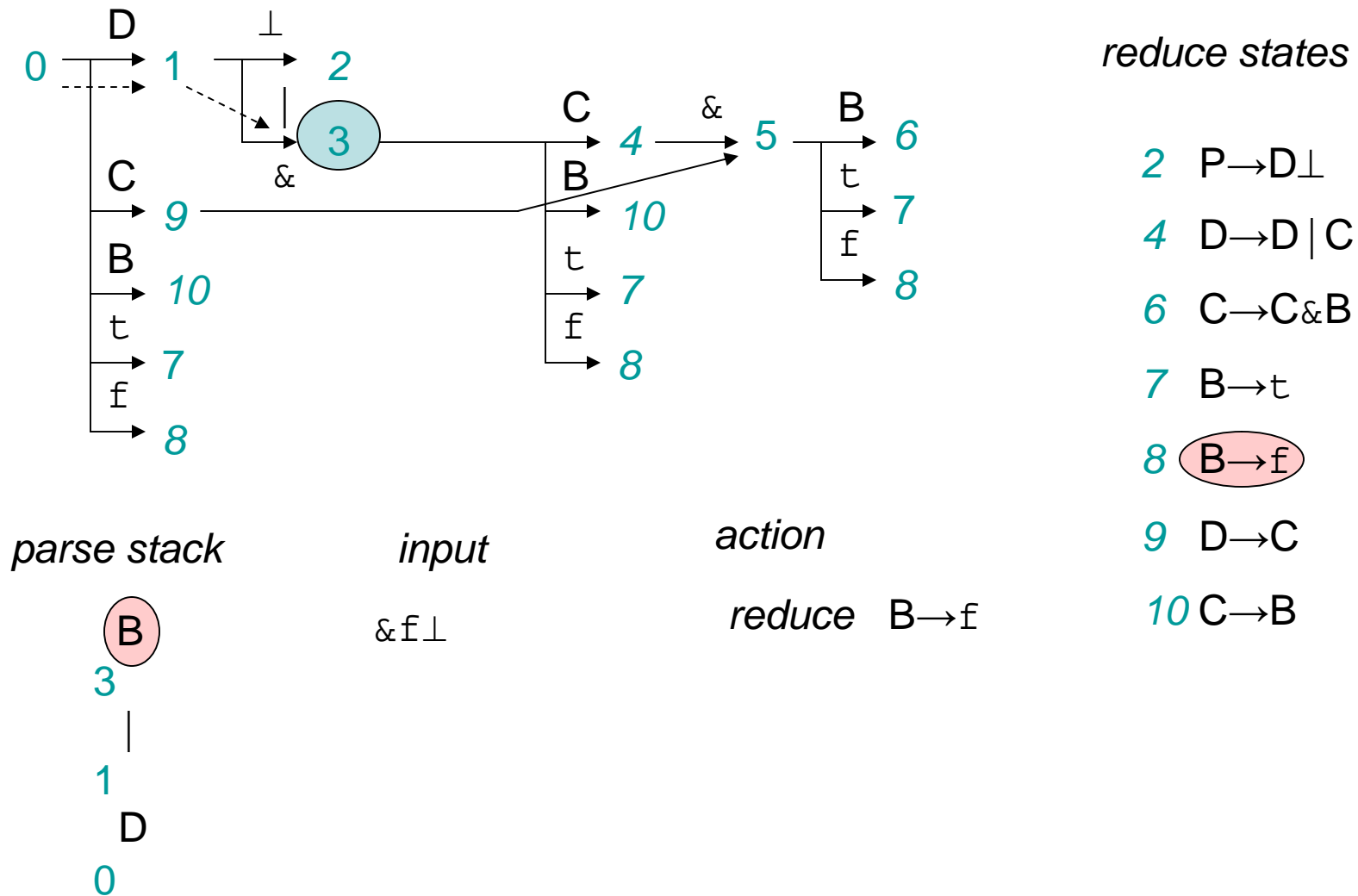




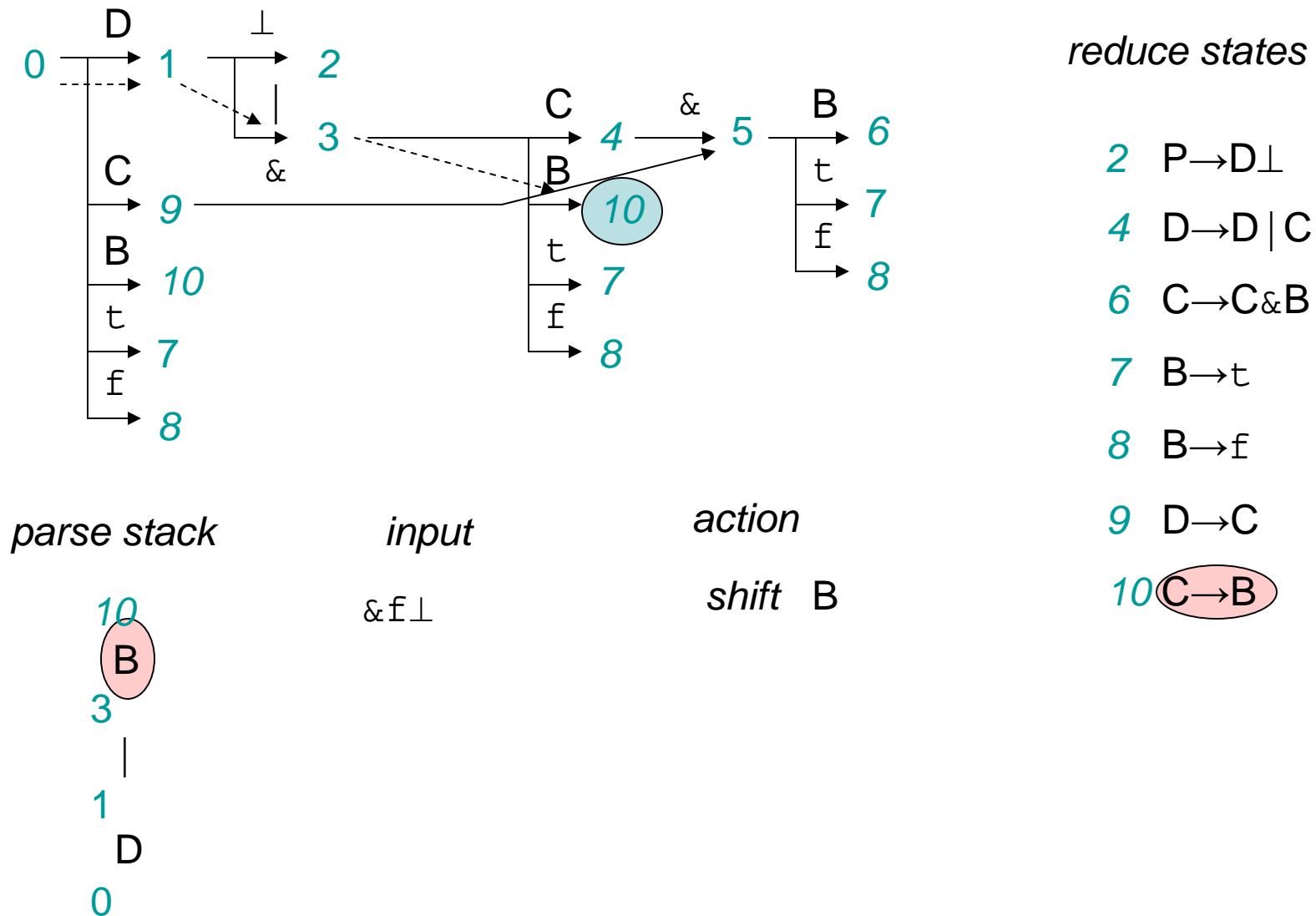
# Putting it all together (16)



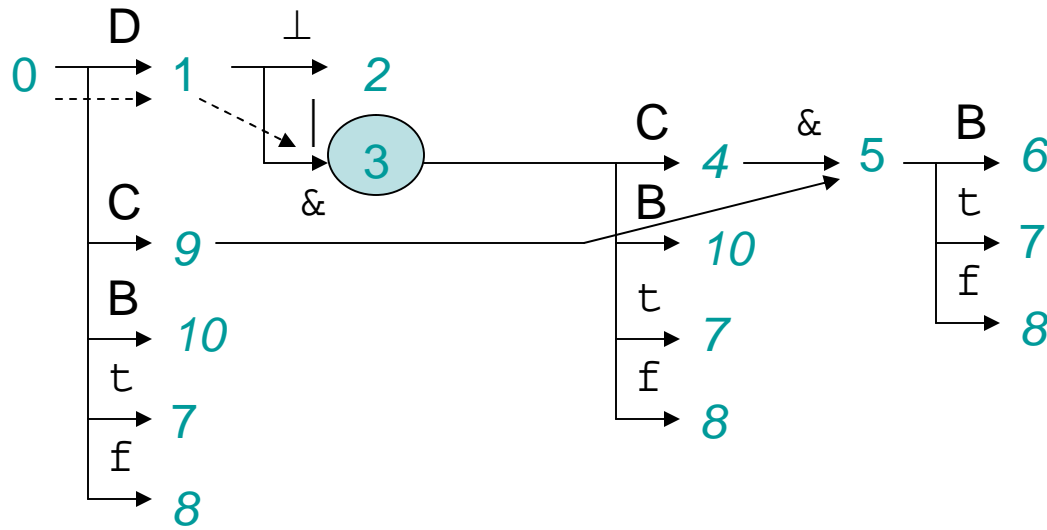
# Putting it all together (17)



# Putting it all together (18)



# Putting it all together (19)



*reduce states*

2  $P \rightarrow D \perp$

4  $D \rightarrow D \mid C$

6  $C \rightarrow C \& B$

7  $B \rightarrow t$

8  $B \rightarrow f$

9  $D \rightarrow C$

10  $C \rightarrow B$

*parse stack*



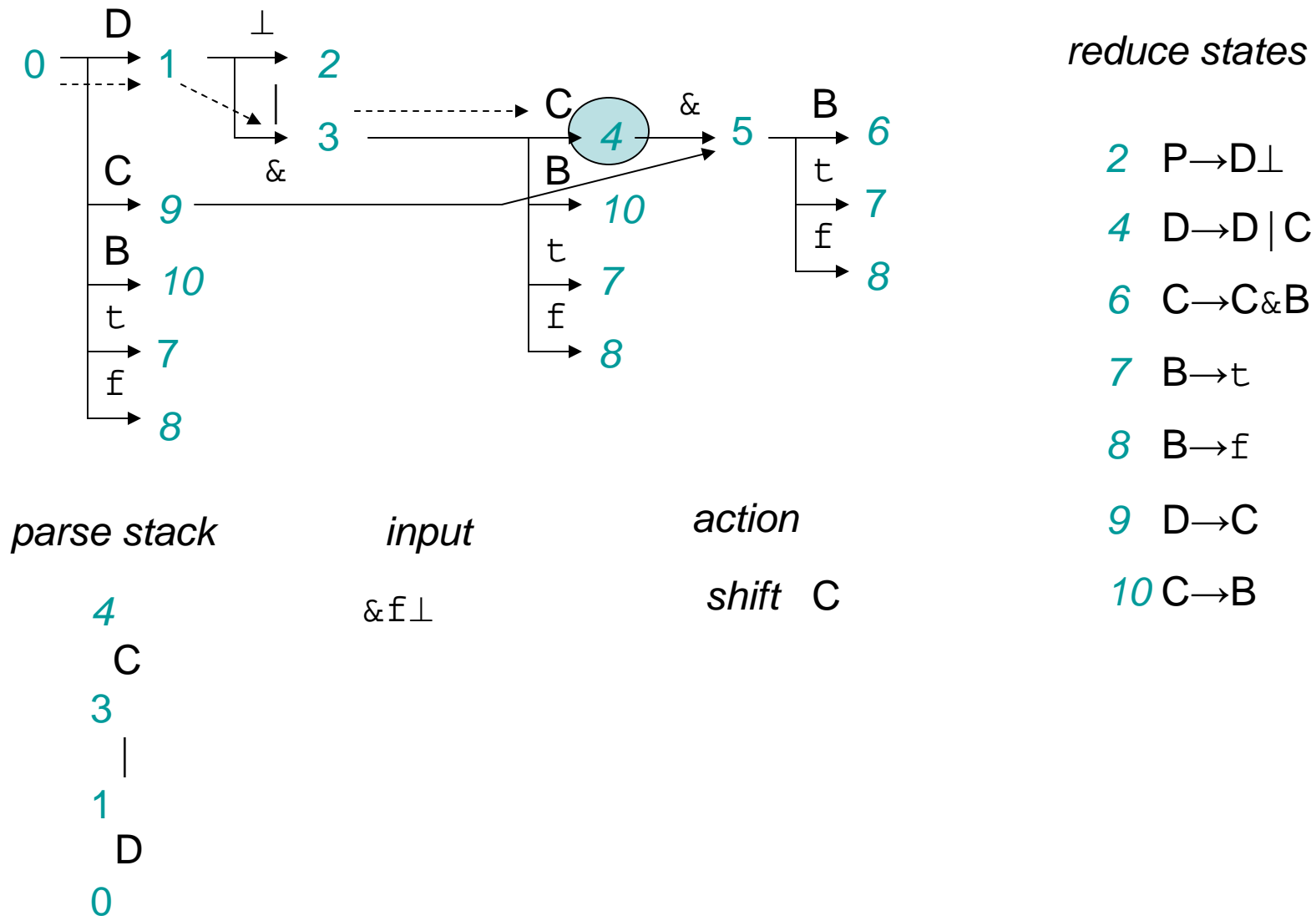
*input*

$\& f \perp$

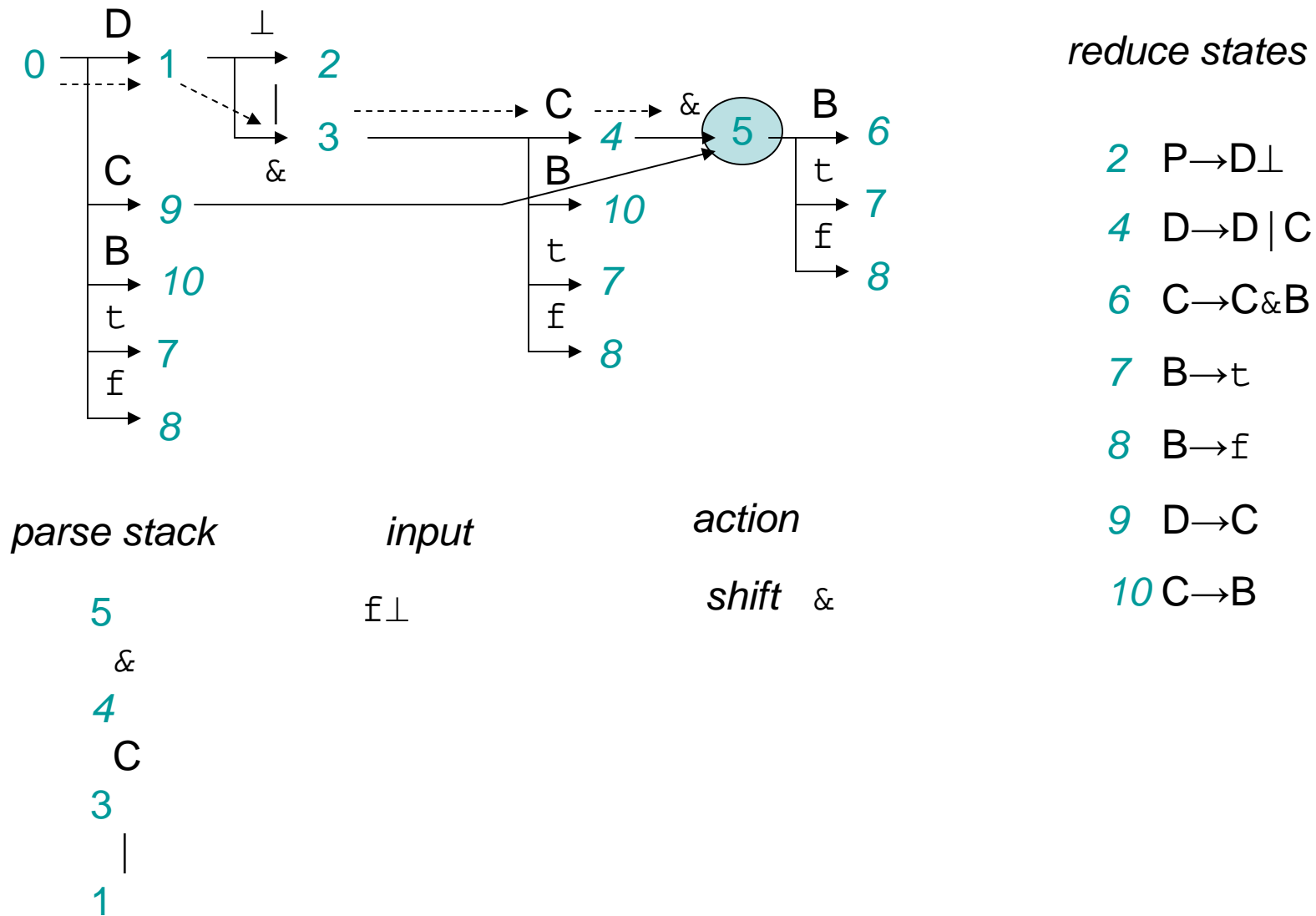
*action*

*reduce*  $C \rightarrow B$

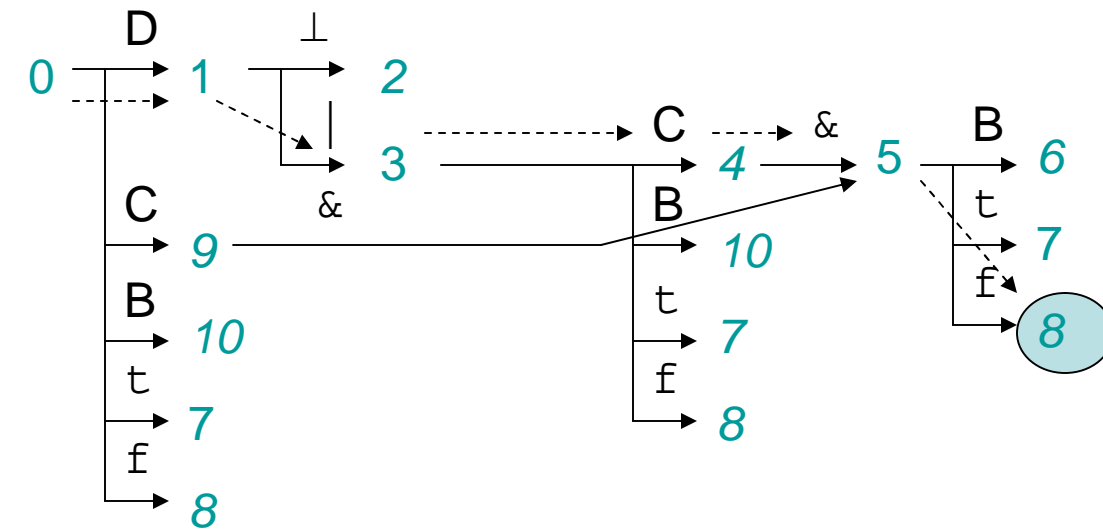
# Putting it all together (20)



# Putting it all together (21)



# Putting it all together (22)



*reduce states*

- 2  $P \rightarrow D \perp$
- 4  $D \rightarrow D \mid C$
- 6  $C \rightarrow C \& B$
- 7  $B \rightarrow t$
- 8  $B \rightarrow f$
- 9  $D \rightarrow C$
- 10  $C \rightarrow B$

*parse stack*

8  
f  
5  
&  
4  
C  
3

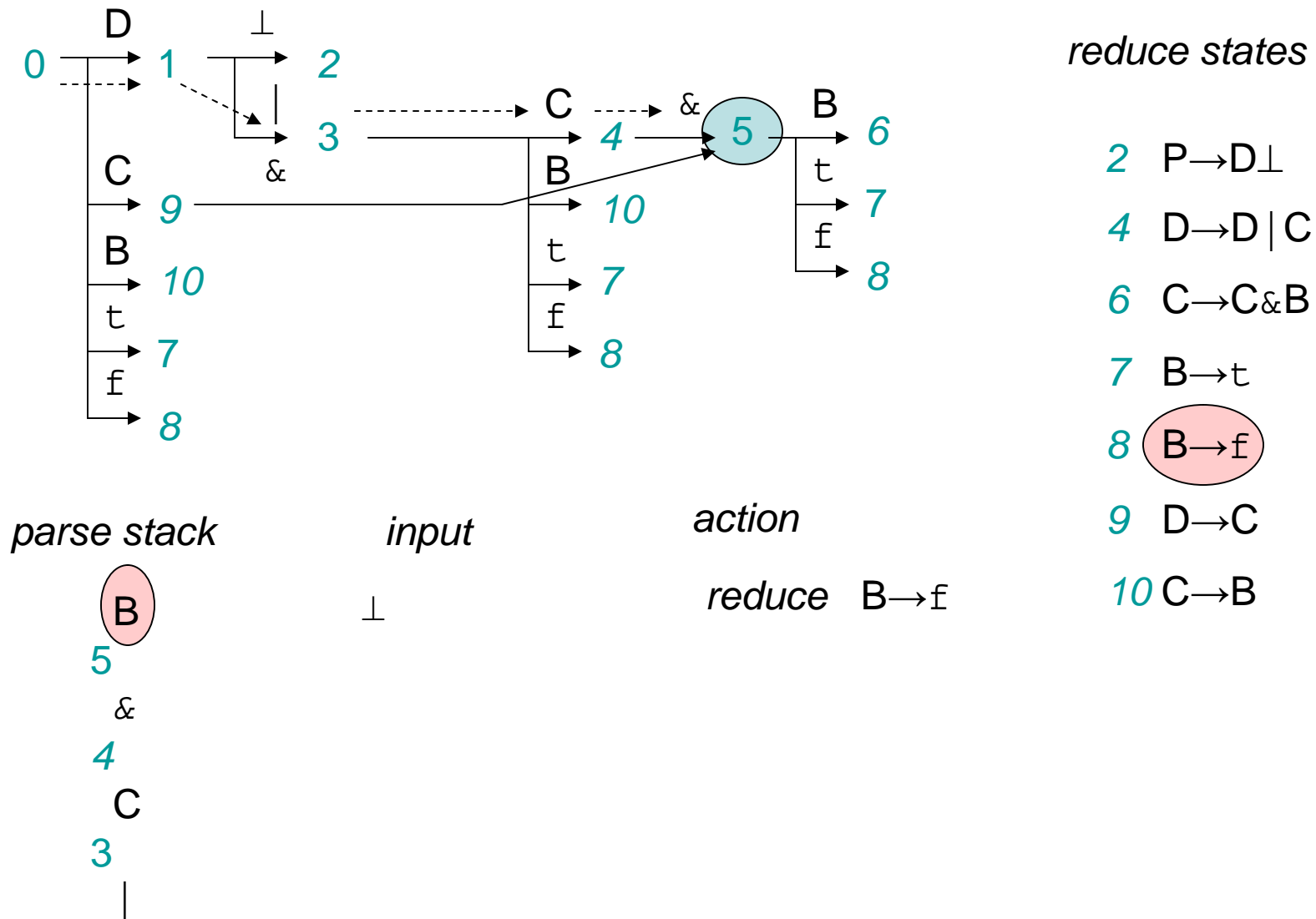
*input*

$\perp$

*action*

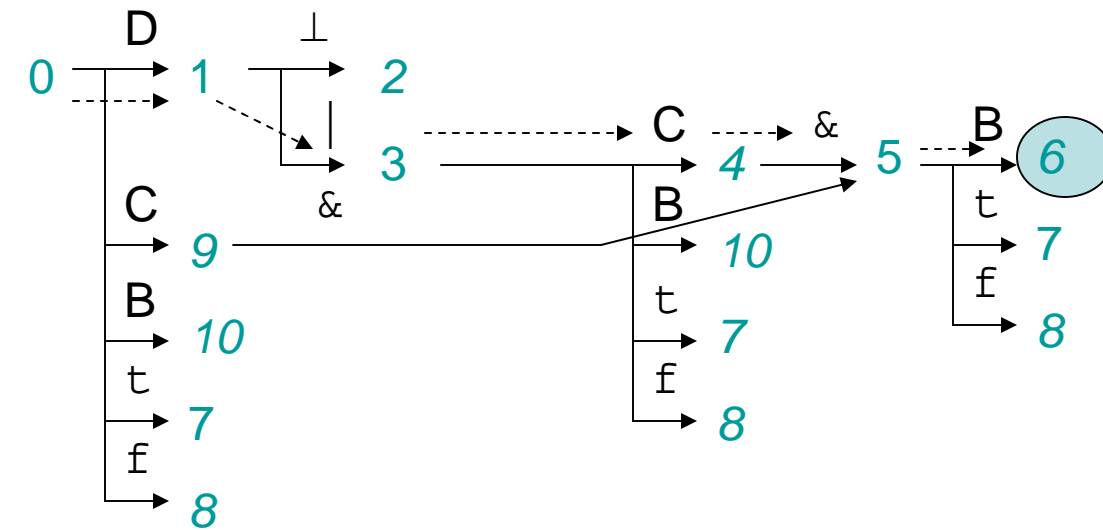
shift f

# Putting it all together (23)





# Putting it all together (24)



*reduce states*

2  $P \rightarrow D \perp$

4  $D \rightarrow D \mid C$

6  $C \rightarrow C \& B$

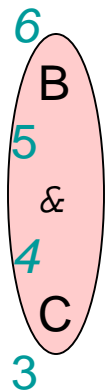
7  $B \rightarrow t$

8  $B \rightarrow f$

9  $D \rightarrow C$

10  $C \rightarrow B$

*parse stack*



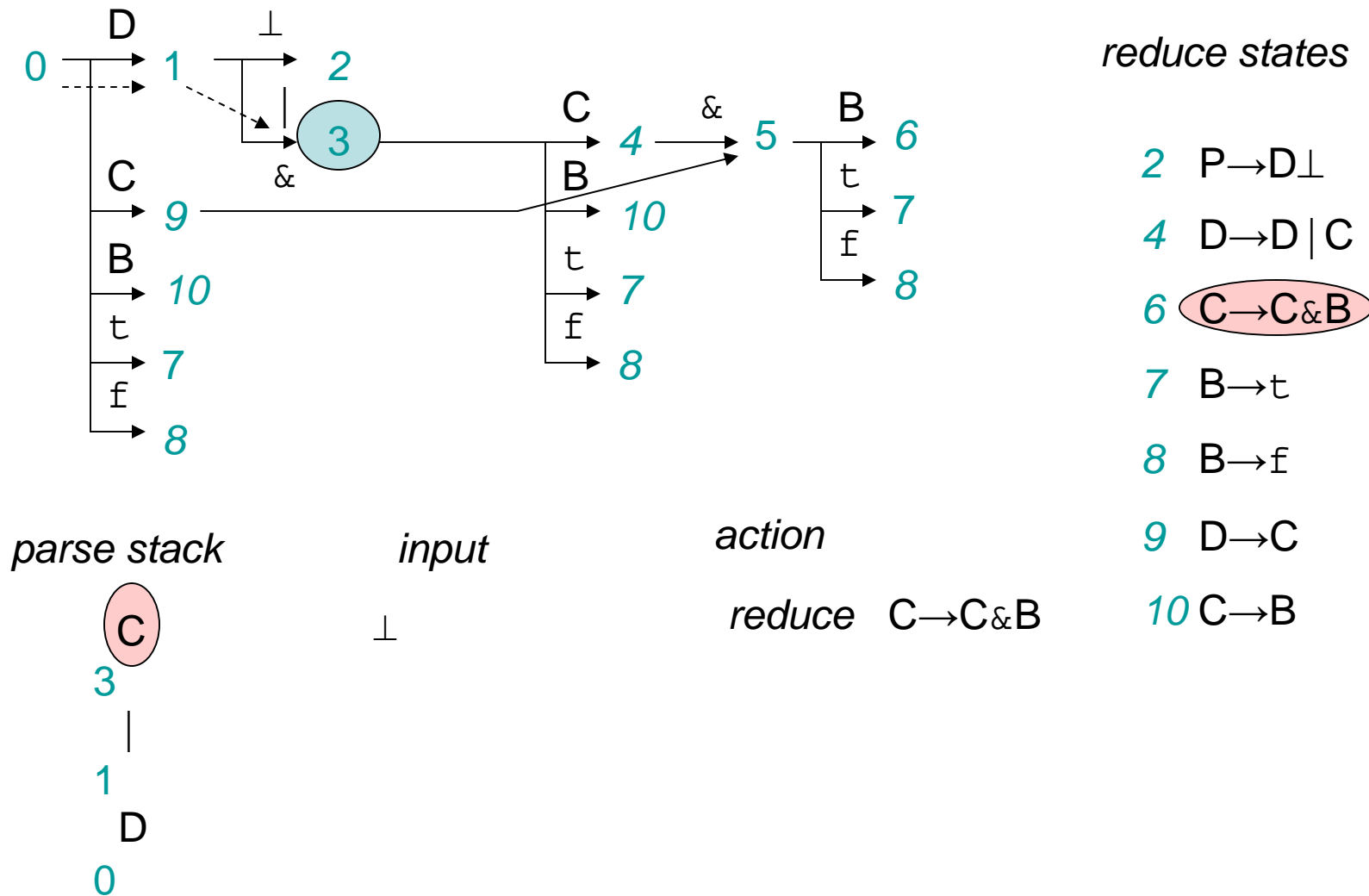
*input*

$\perp$

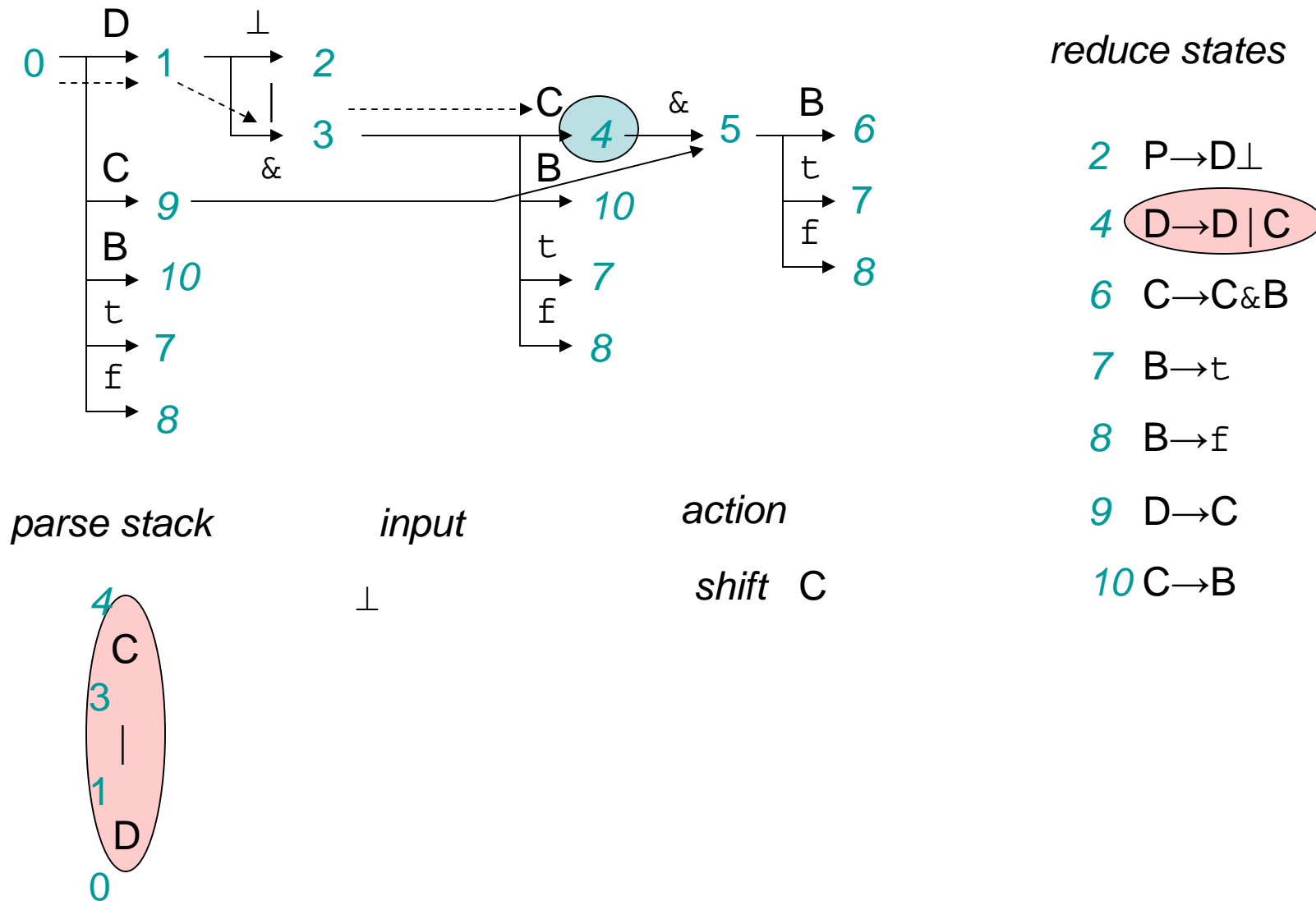
*action*

*shift* B

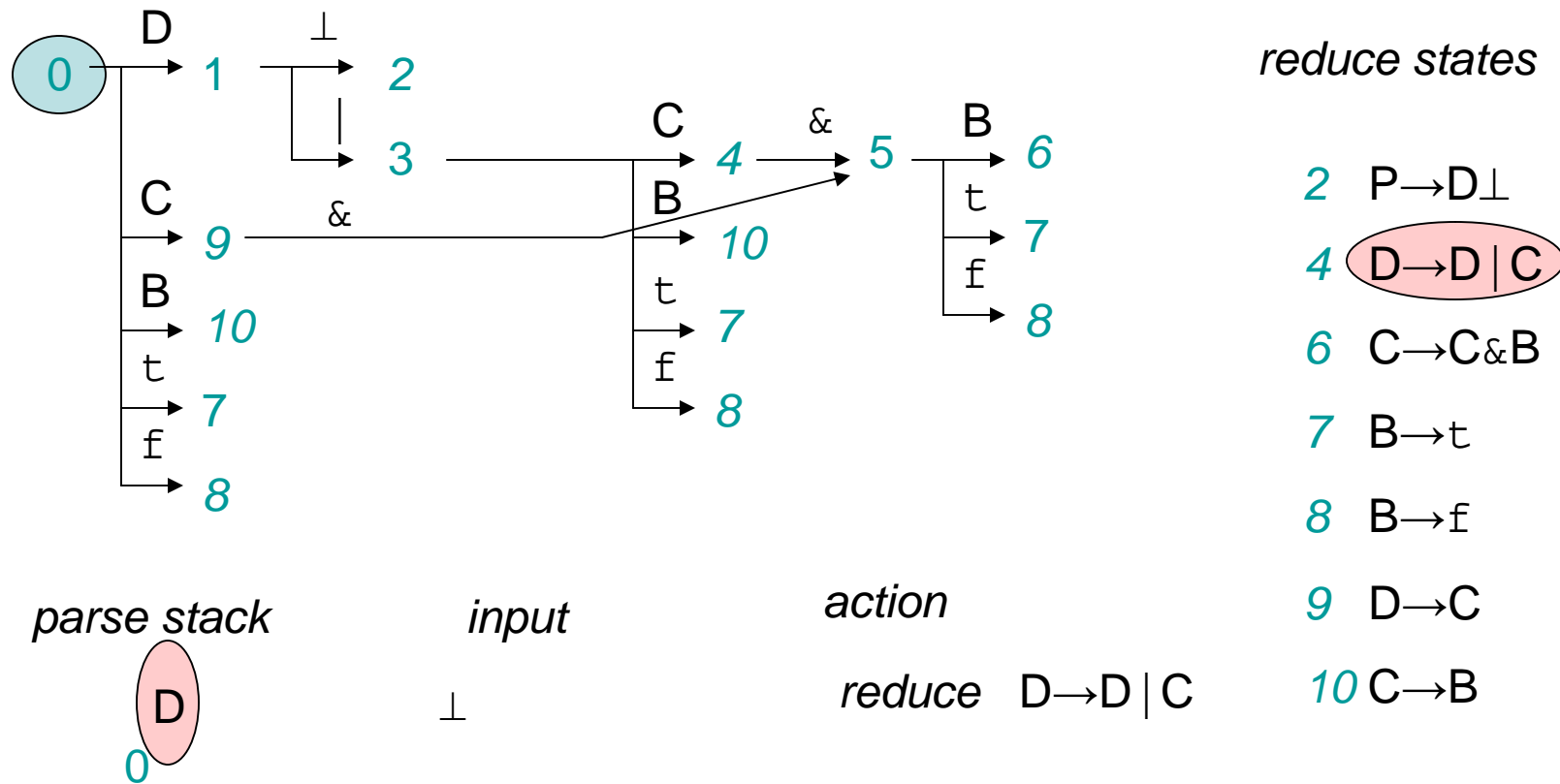
# Putting it all together (25)



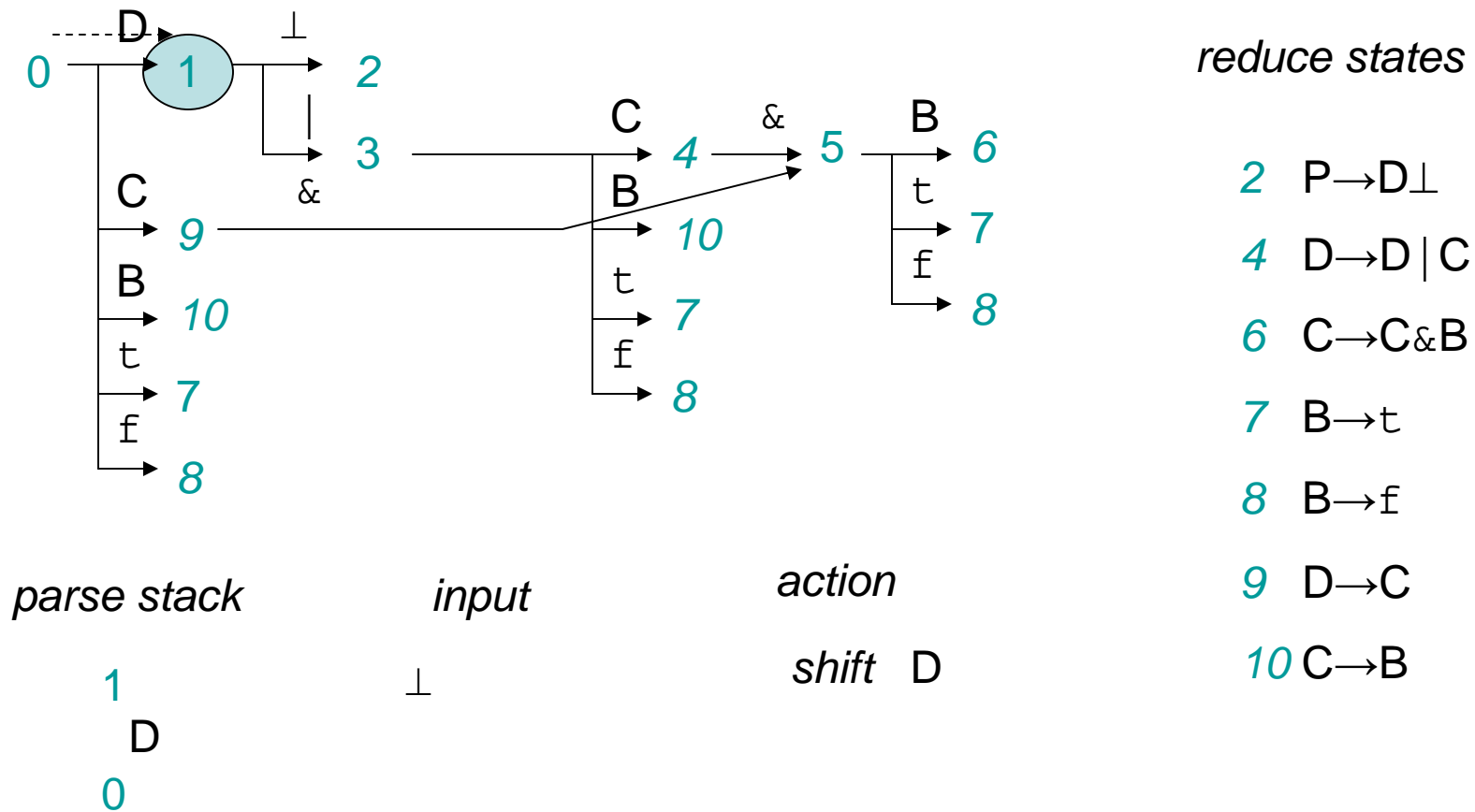
# Putting it all together (26)



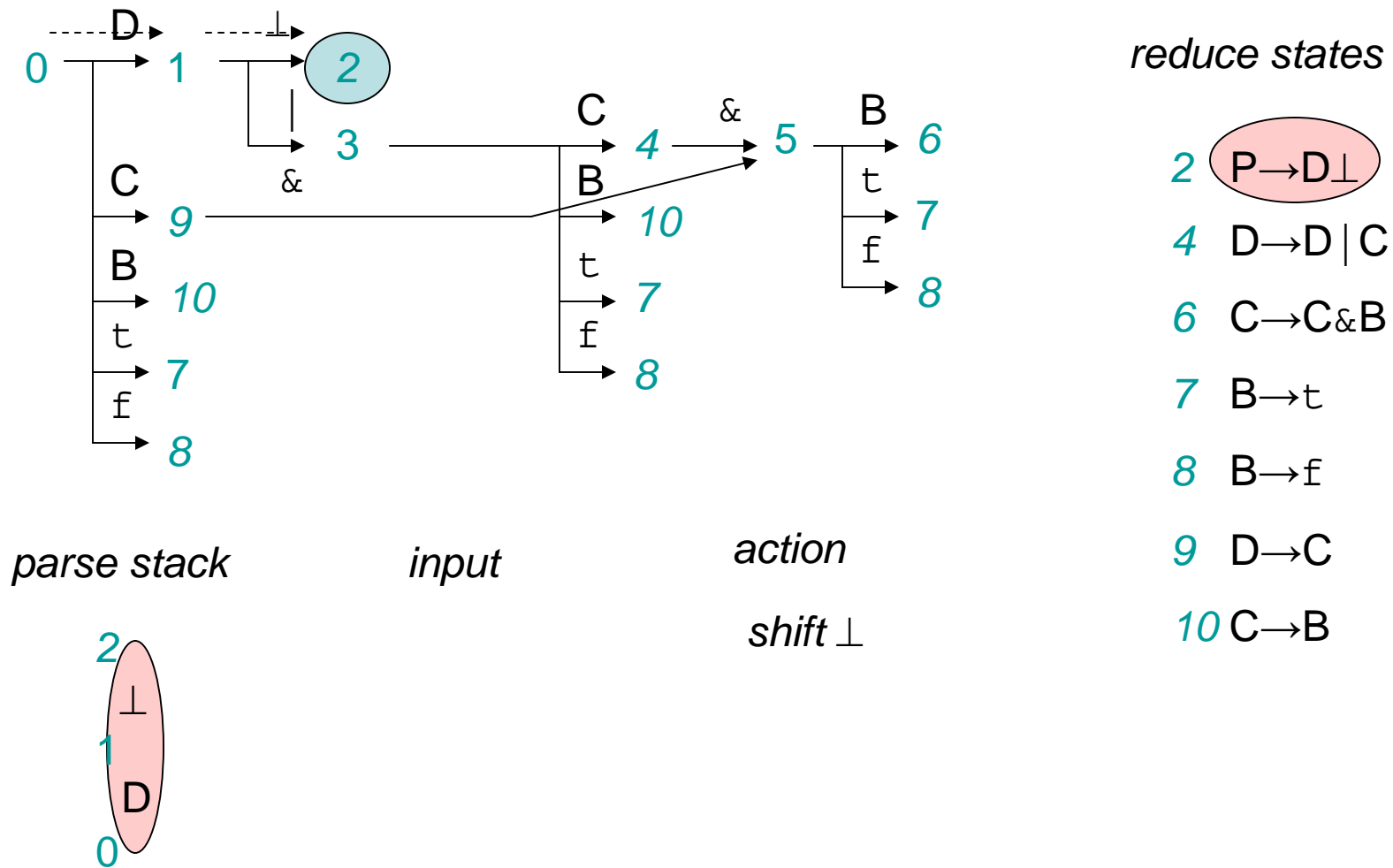
# Putting it all together (27)



# Putting it all together (28)



# Putting it all together (29)



# Putting it all together (30)

