

Compilation

0368-3133

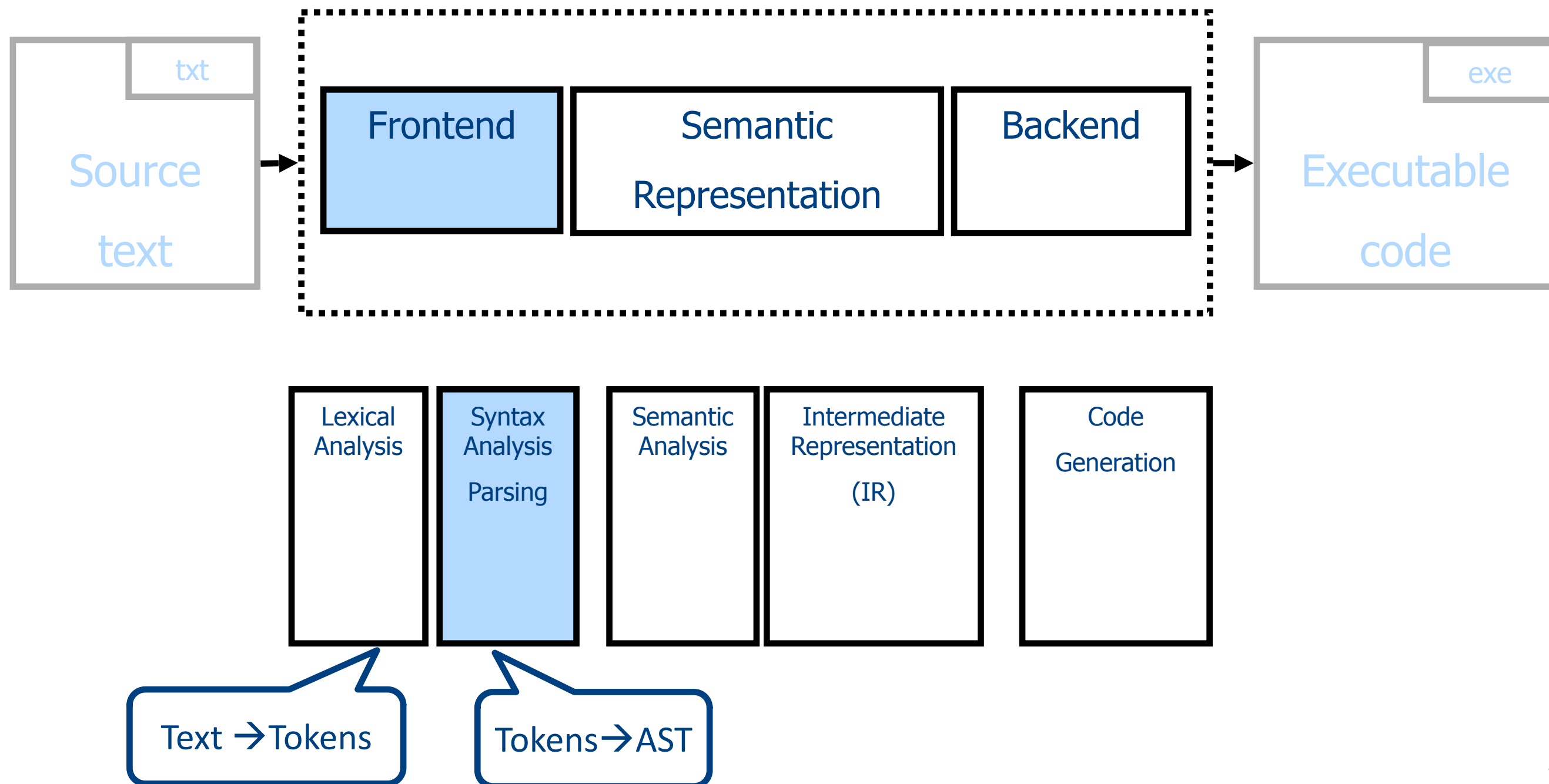
Lecture 3b:

Syntax Analysis:

Bottom-Up Parsing

Conceptual Structure of a Compiler

Compiler

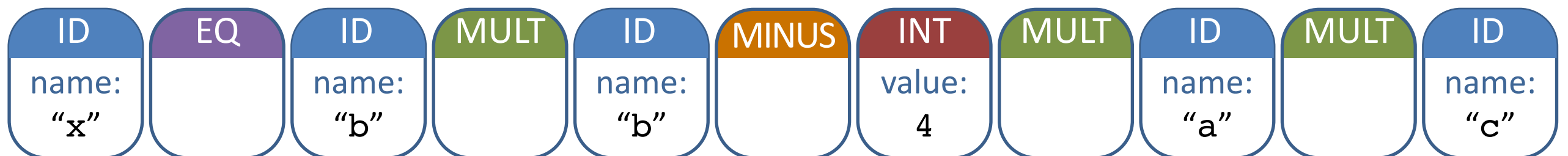


Lexing: from Characters to Tokens

	txt
<code>x = b*b - 4*a*c</code>	

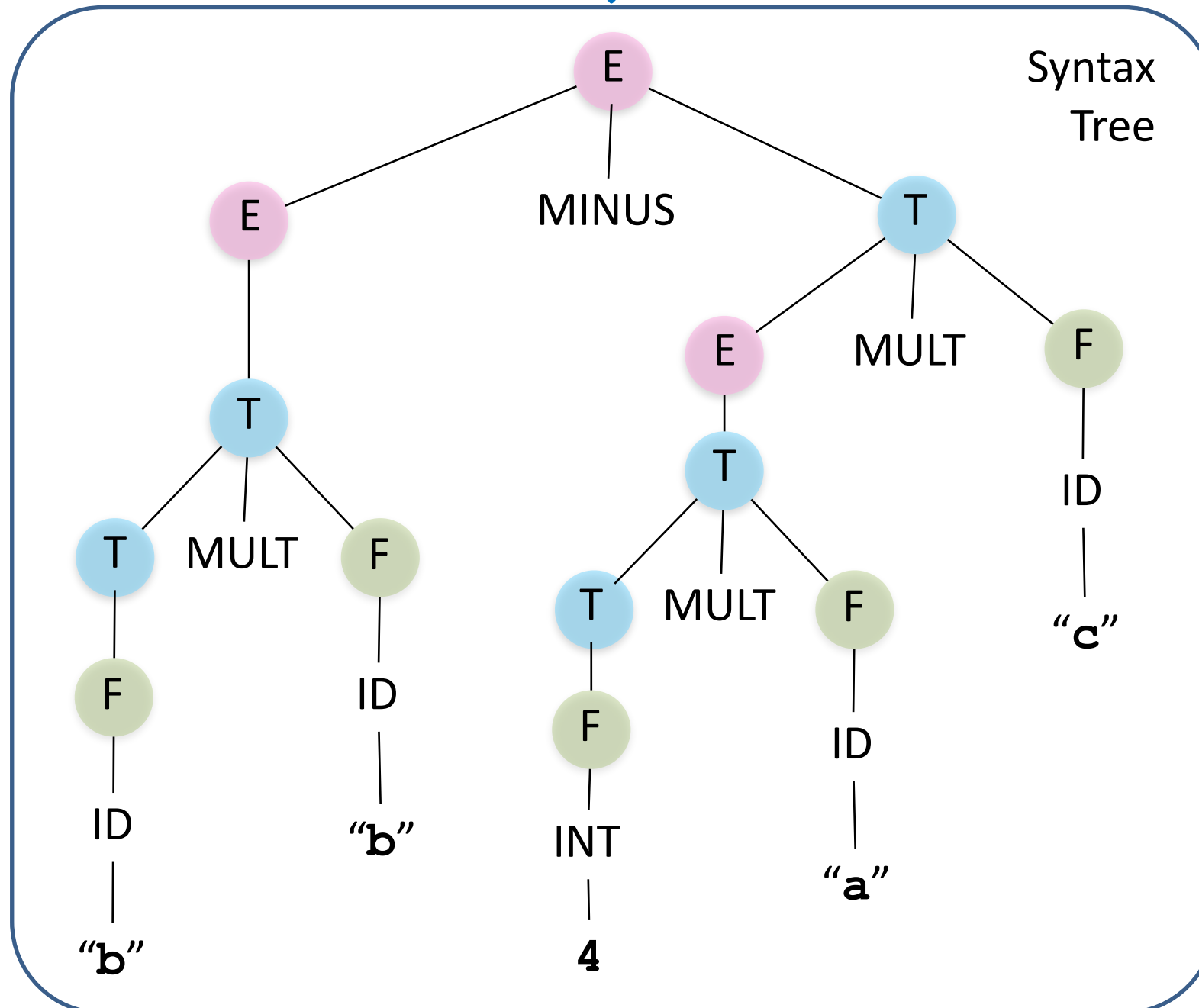


Token Stream



Parsing: from Tokens to Syntax Tree

$\langle \text{ID}, "b" \rangle \langle \text{MULT} \rangle \langle \text{ID}, "b" \rangle \langle \text{MINUS} \rangle \langle \text{INT}, 4 \rangle \langle \text{MULT} \rangle \langle \text{ID}, "a" \rangle \langle \text{MULT} \rangle \langle \text{ID}, "c" \rangle$



Broad kinds of parsers

- **Top-Down** parsers (*LL(k) grammars*)
 - Construct parse tree in a top-down matter
 - Find the **leftmost** derivation
- **Bottom-Up** parsers (*LR(k) grammars*)
 - Construct parse tree in a bottom-up manner
 - Find the **rightmost** derivation (in a reverse order)
- Parsers for **arbitrary** grammars
 - **Earley**'s method ($O(n^3)$ for ambiguous grammars. $O(n^2)$ for unambiguous grammars), **CYK** method ($O(n^3)$)
 - Usually, not used in practice (though might change)

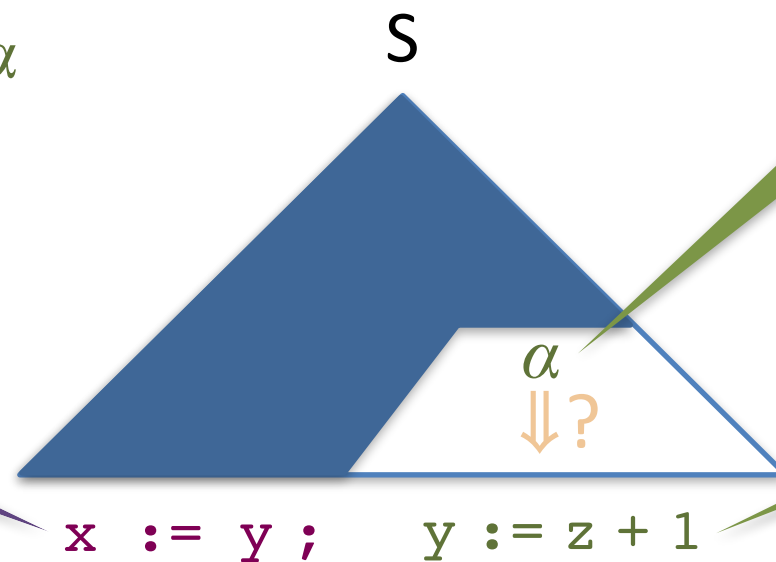
Efficient
($O(n=|input|)$)

Efficient Parsers

- Top-down (predictive)

Sentential form: $x := y ; \alpha$

already read...



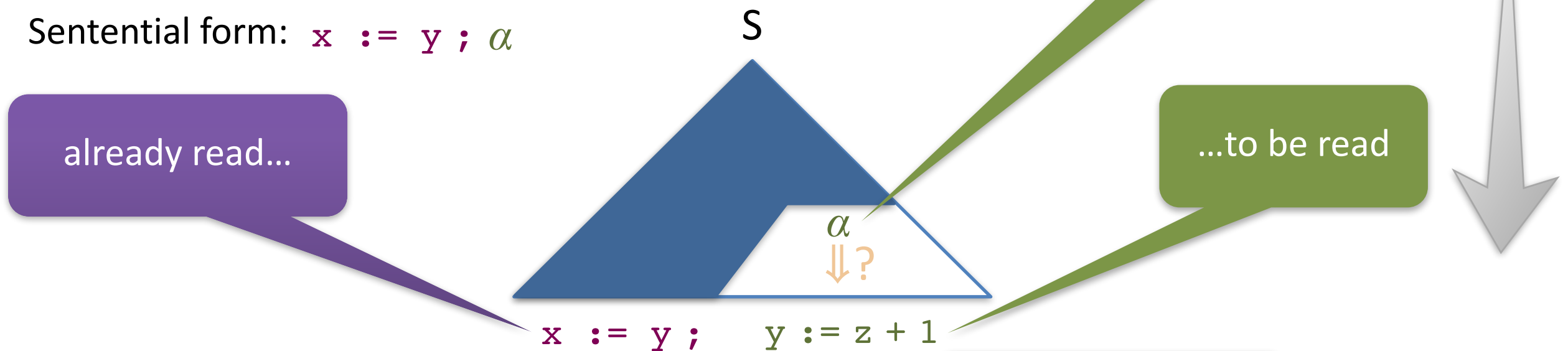
suffix of a sentential form predicted to derive the suffix of the input

...to be read

Efficient Parsers

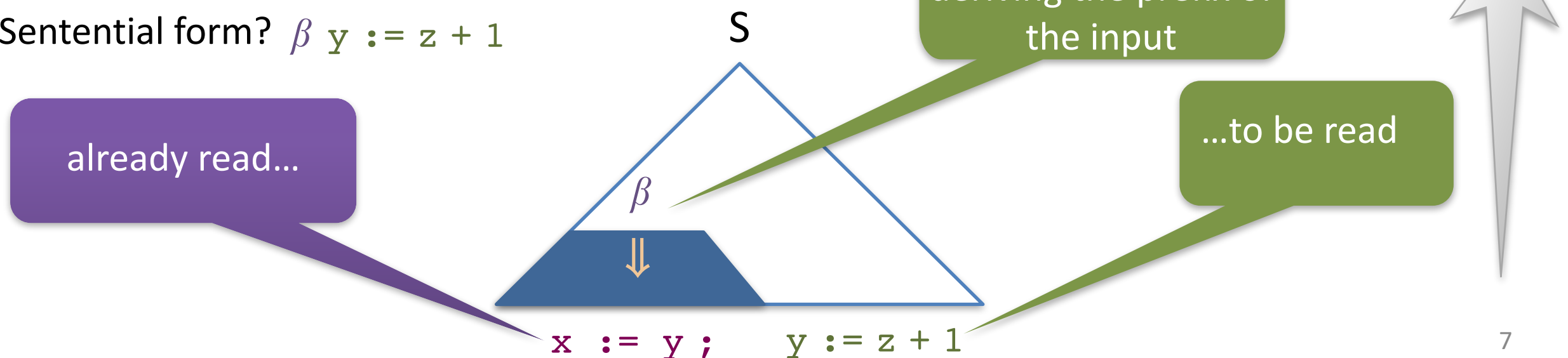
- Top-down (predictive)

Sentential form: $x := y ; \alpha$



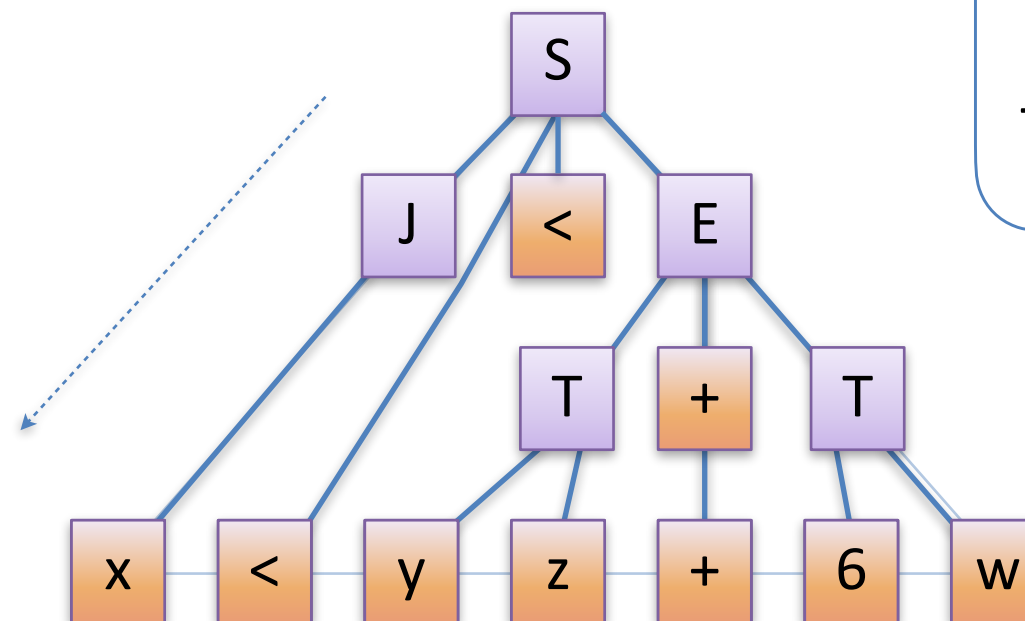
- Bottom-up (shift-reduce)

Sentential form? $\beta y := z + 1$



Efficient Parsers – Illustrated

- Top-down (predictive)



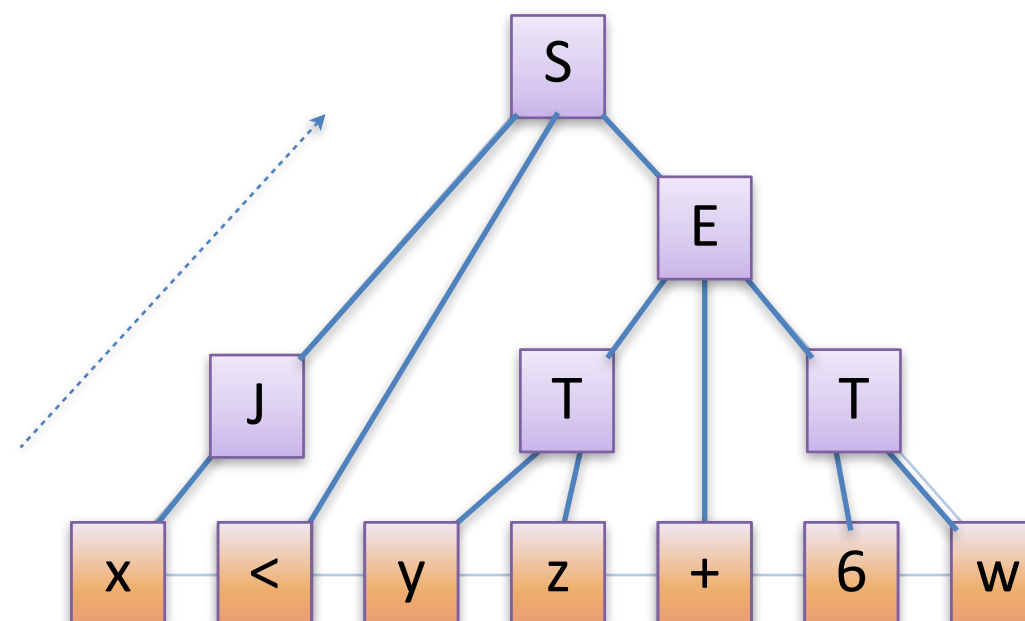
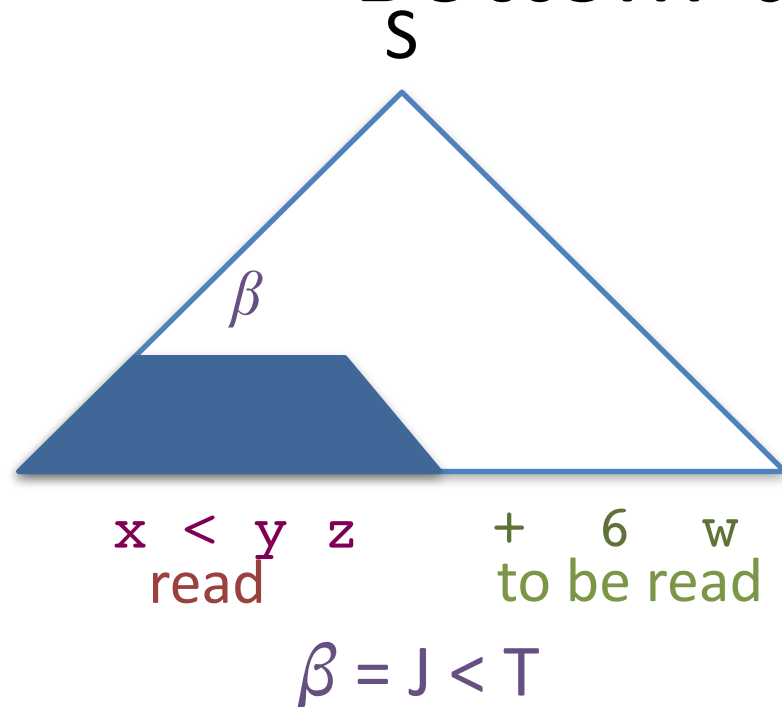
$S \rightarrow J < E$

$J \rightarrow \text{id}$

$E \rightarrow T + T$

$T \rightarrow \text{id id} \mid \text{num id}$

- Bottom-up (shift-reduce)



Terminology: Reductions

- The opposite of **derivation** is called *reduction*
 - Let $A \rightarrow \alpha$ be a production rule
 - Derivation (step): $\beta A \mu \Rightarrow \beta \alpha \mu$
 - Reduction (step): $\beta \alpha \mu \Rightarrow \beta A \mu$
- A *reduction* is a sequence of reduction steps

Bottom-Up Parsing

- Goal: Build a parse tree*
 - Report error if the input is not a legal program
- How:
 - Scan input left-to-right
 - Construct subtrees once all their leaves have been scanned
 - Constructing the first left-most subtree whose children have been constructed
 - **Reduces** the input to the start symbol
 - Finds a rightmost derivation (in reverse order)

* Actually, we want to construct an abstract syntax tree (AST) ; we'll get to it later

Rightmost Derivation

$x := z;$
 $y := x + z$

$S \rightarrow S;S \mid id := E$

$E \rightarrow id \mid num \mid E + E$

S
 $S ; S$

$S ; id := E$

$S ; id := E + E$

$S ; id := E + id$

$S ; id := id + id$

$id := E ; id := id + id$

$id := id ; id := id + id$

$\langle id, "x" \rangle \text{ ASS } \langle id, "z" \rangle ; \langle id, "y" \rangle \text{ ASS } \langle id, "x" \rangle \text{ PLUS } \langle id, "z" \rangle$

$S \rightarrow S;S$

$S \rightarrow id := E$

$E \rightarrow E + E$

$E \rightarrow id$

$E \rightarrow id$

$S \rightarrow id := E$

$E \rightarrow id$

A Reduction in the reverse order of a Rightmost Derivation

```
x := z;
y := x + z
```

$$S \rightarrow S;S \mid id := E$$

$$E \rightarrow id \mid num \mid E + E$$

	S	
	S ; S	
	S ; id := E	
	S ; id := E + E	
	S ; id := E + id	
	S ; id := id + id	
id :=	E ; id := id + id	
id :=	id ; id := id + id	

<id,"x"> ASS <id,"z"> ; <id,"y"> ASS <id,"x"> PLUS <id,"z">

$$S \rightarrow S;S$$

$$S \rightarrow id := E$$

$$E \rightarrow E + E$$

$$E \rightarrow id$$

$$E \rightarrow id$$

$$S \rightarrow id := E$$

$$E \rightarrow id$$

Reduction Example

$$E \rightarrow E * B \mid E + B \mid B$$
$$B \rightarrow 0 \mid 1$$

- Let us number the rules:

$$(1) E \rightarrow E * B$$

$$(2) E \rightarrow E + B$$

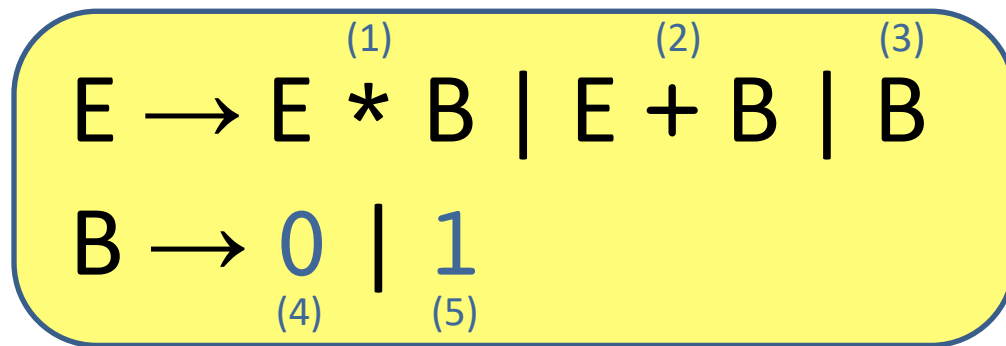
$$(3) E \rightarrow B$$

$$(4) B \rightarrow 0$$

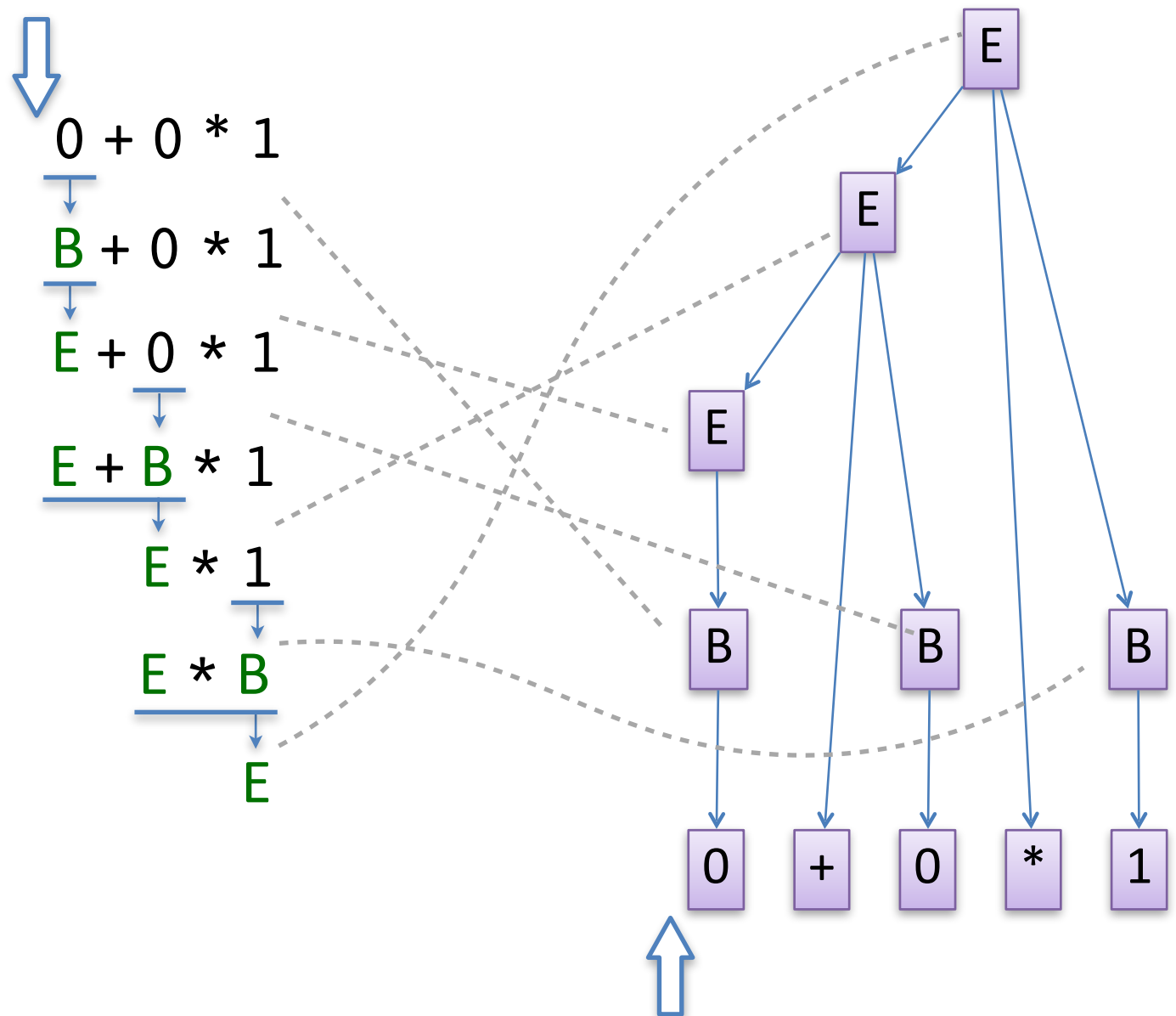
$$(5) B \rightarrow 1$$

Goal: Reduce the Input to the Start Symbol

Building the parse tree during the reduction



Rightmost
derivation

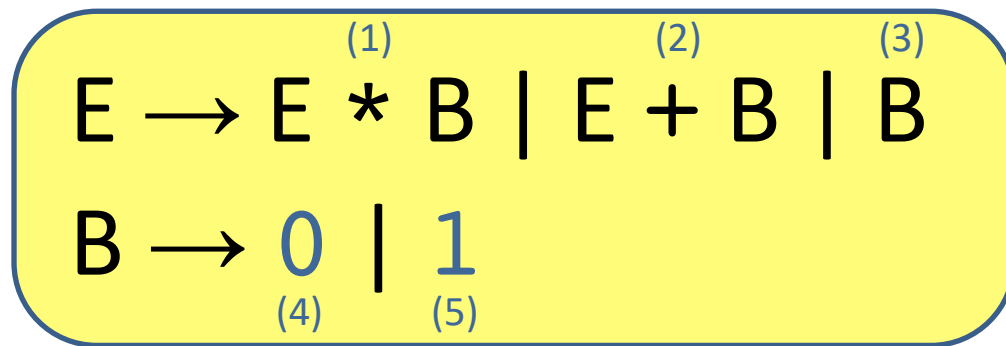


Let $A \rightarrow \alpha$ be a production rule

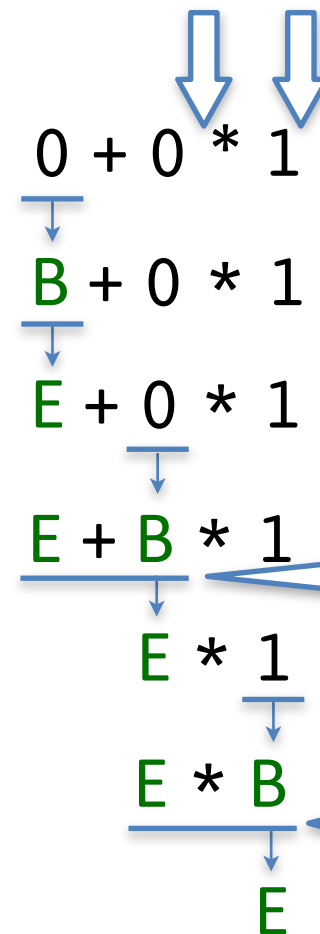
- Derivation: $\beta A \mu \Rightarrow \beta \alpha \mu$
- Reduction: $\beta \alpha \mu \Rightarrow \beta A \mu$

Goal: Reduce the Input to the Start Symbol

Building the parse tree during the reduction

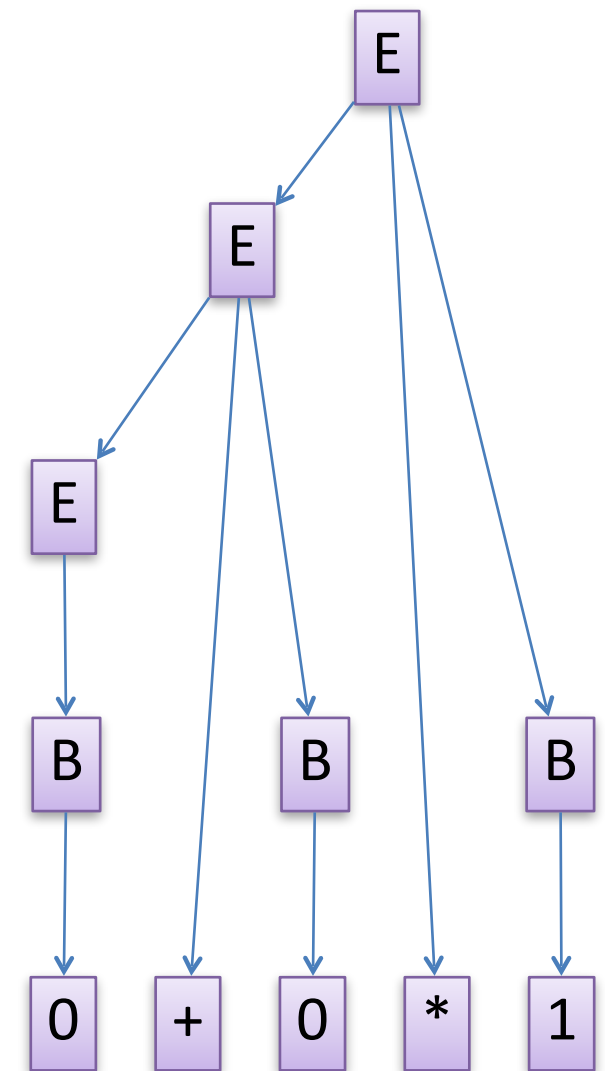


Rightmost
derivation



Using (3)
would be a
mistake ...
(check it out...)

Using (3)
would be a
mistake ...
(check it out...)



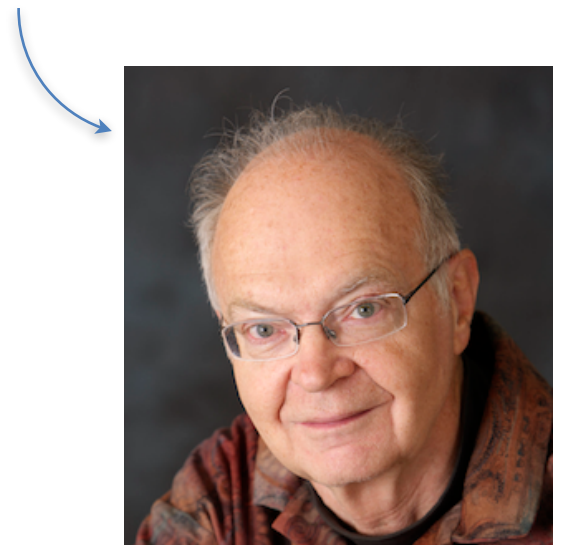
Let $A \rightarrow \alpha$ be a production rule

- Derivation: $\beta A \mu \Rightarrow \beta \alpha \mu$
- Reduction: $\beta \alpha \mu \Rightarrow \beta A \mu$

LR(k) Grammars

- A grammar is in the class LR(k) when it can be derived via:
 - ▶ **Bottom-up** analysis
 - ▶ Scanning the input from **left to right** (L)
 - ▶ Producing the **rightmost derivation** (R)
 - In reverse order
 - ▶ With **lookahead** of k tokens (k)
- A language is said to be LR(k) if it has an LR(k) grammar

Donald Knuth



How does the Magic Work?

- Main idea: Delay decision on which rule to reduce with next to the last possible minute

How does the Magic Work? (by Example)

$S \rightarrow A \mid C$

$A \rightarrow BC \mid BB \mid BbC \mid CB$

$B \rightarrow bb$

$C \rightarrow cc$

Regular Language!

$L = \{ \text{?} \}$

Is the grammar LL(1)?

Is the grammar LL(k) for some k ?

How does the Magic Work?

(by Example)

$S \rightarrow A \mid C$

$A \rightarrow BC \mid BB \mid BbC \mid CB$

$B \rightarrow bb$

$C \rightarrow cc$

Rightmost derivation

$w = bbbcc$

$S \Rightarrow$

$A \Rightarrow$

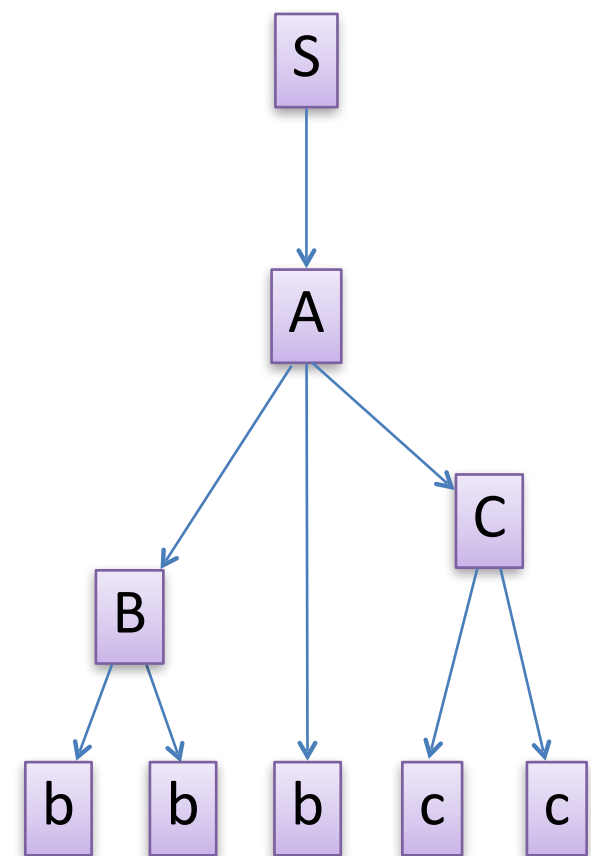
$BbC \Rightarrow$

$Bbcc \Rightarrow$

$bbbcc$

*desired
reduction*

Parse tree



How does the Magic Work?

(by Example)

$S \rightarrow A \mid C$

$A \rightarrow BC \mid BB \mid BbC \mid CB$

$B \rightarrow bb$

$C \rightarrow cc$

B

B ← bb

b

B ← bb

b

State

S ← A | C

A ← BC | BB | BbC | CB

B ← bb C ← cc

S

B

b

b

b

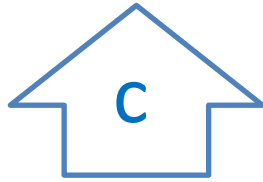
c

c

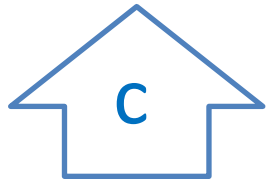
C How does the Magic Work?

(by Example)

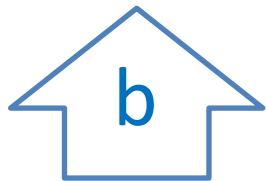
C ← cc



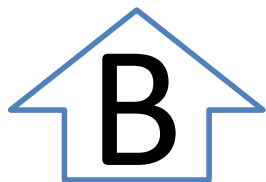
C ← cc



B ← bb A ← BbC C ← cc



A ← BC | BB | BbC
C ← cc B ← bb



S ← A | C
A ← BC | BB | BbC | CB
B ← bb C ← cc

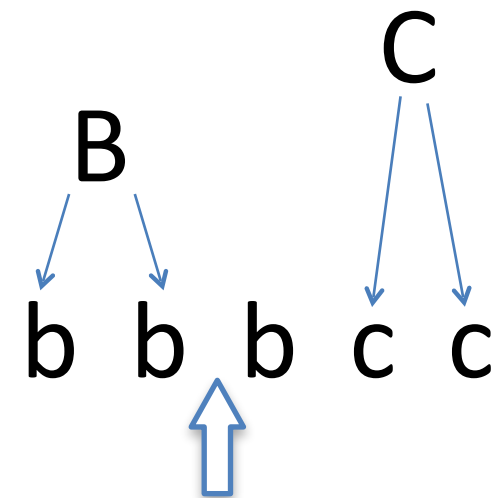


S → A | C

A → BC | BB | BbC | CB

B → bb

C → cc



How does the Magic Work?

(by Example)

A

A ← BbC

C

B ← bb A ← BbC C ← cc

b

A ← BC | BB | BbC
C ← cc B ← bb

B

S ← A | C
A ← BC | BB | BbC | CB
B ← bb C ← cc

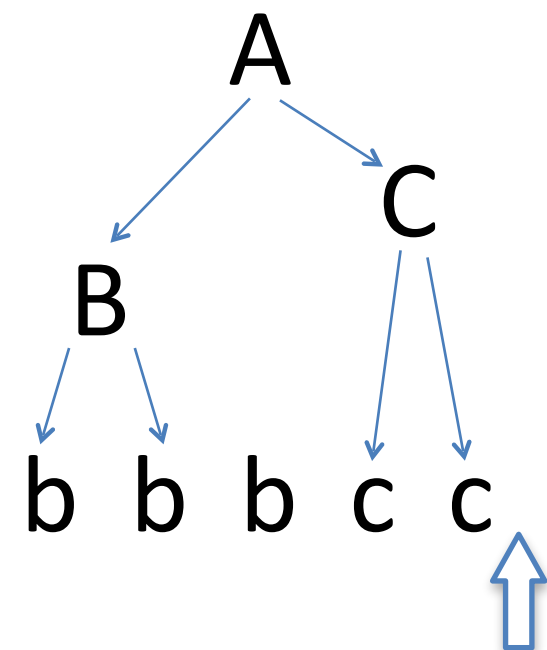
S

S → A | C

A → BC | BB | BbC | CB

B → bb

C → cc



How does the Magic Work?

(by Example)

BTW. what
data structure
did we use?

S

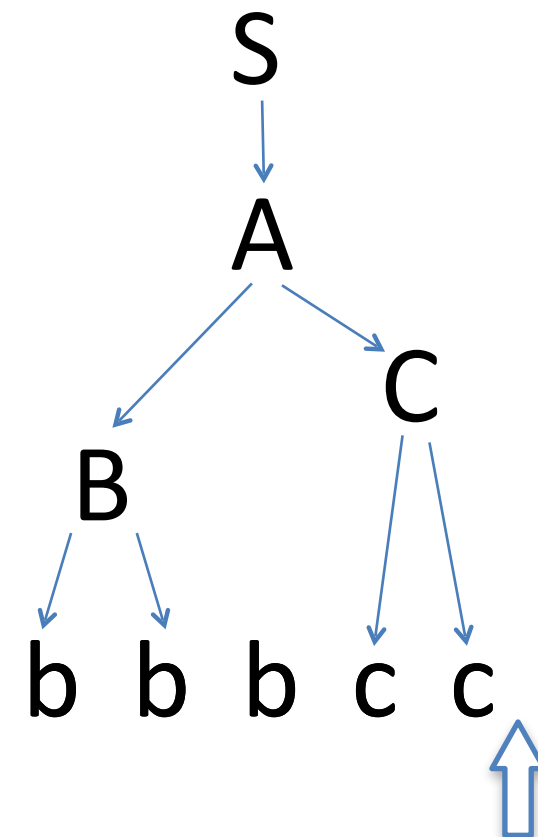
$S \rightarrow A$

A

$S \leftarrow A \mid C$
 $A \leftarrow BC \mid BB \mid BbC \mid CB$
 $B \leftarrow bb \quad C \leftarrow cc$



$S \rightarrow A \mid C$
 $A \rightarrow BC \mid BB \mid BbC \mid CB$
 $B \rightarrow bb$
 $C \rightarrow cc$



Shift & Reduce (LR) Parsers

The parser maintains a **stack** of grammar symbols and states

In each step, we either

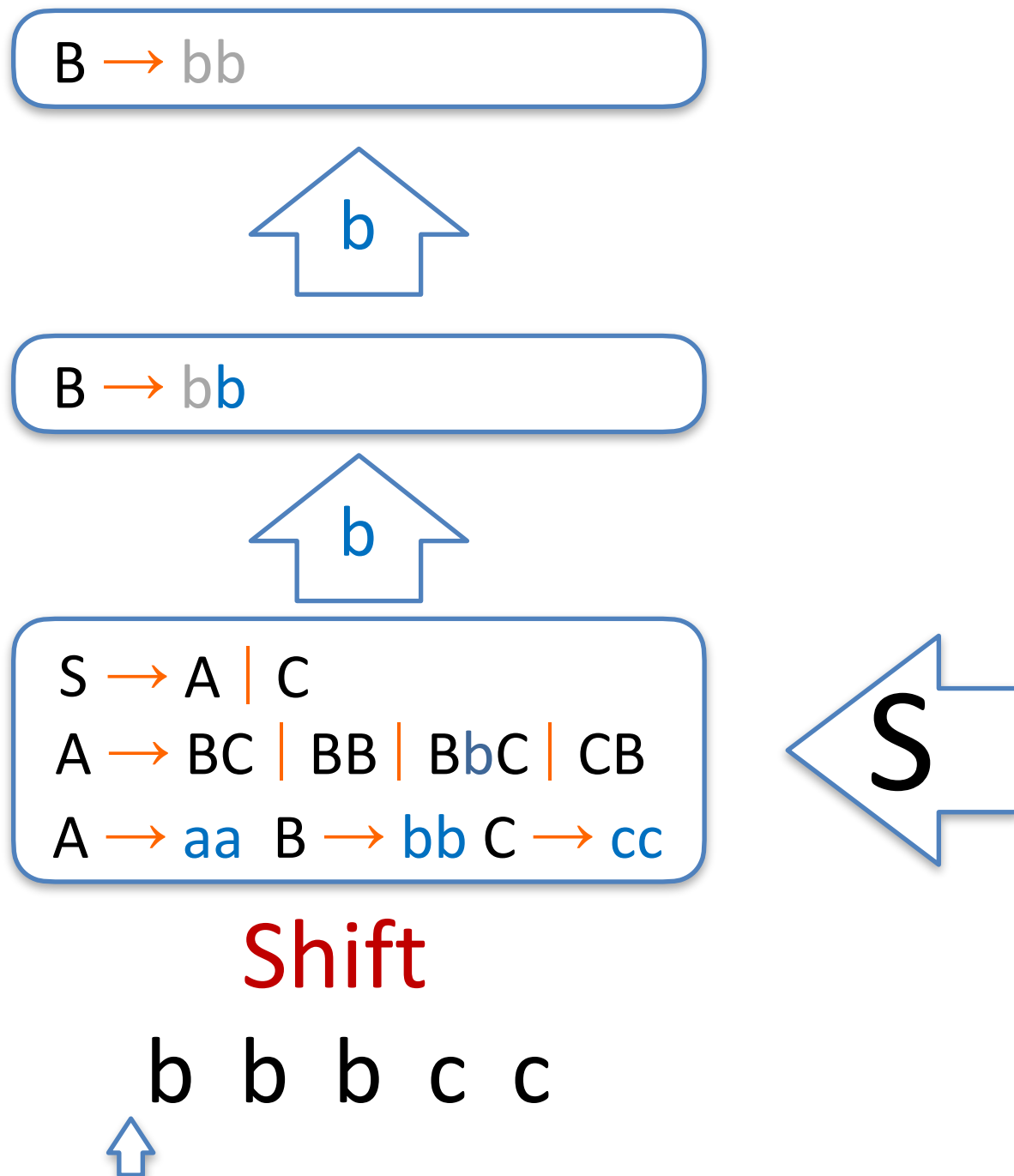
shift moving a symbol from the input to the stack, and
compute and push a new state

or

reduce popping the symbols and states
pertaining to the right-hand side of the
discovered (reduced) rule,
push the non-terminal at the left-hand
side of the rule, and
compute and push a new state

Shift & Reduce (LR) Parsers

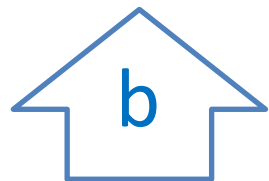
The parser maintains a stack of grammar symbols and states



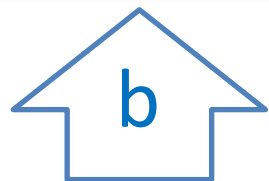
Shift & Reduce (LR) Parsers

The parser maintains a stack of grammar symbols and states

B ← bb



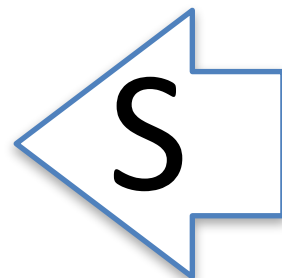
B ← bb



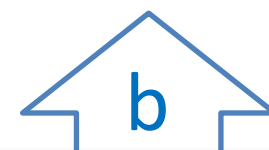
S ← A | C
A ← BC | BB | BbC | CB
A ← aa B ← bb C ← cc

Shift

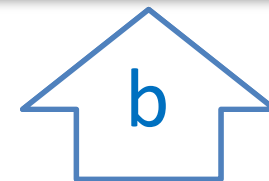
b b b c c



B ← bb



A ← BC | BB | BbC
C ← cc B ← bb



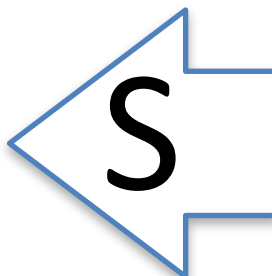
S ← A | C
A ← BC | BB | BbC | CB
A ← aa B ← bb C ← cc

Reduce

b b b c c



State



Important Bottom-Up LR-Parsers

Does an LL(0) parser make sense?

- LR(0) – simplest, explains basic ideas
- SLR(1) – simple, explains lookahead
- LR(1) – complicated, very powerful, expensive
- LALR(1) – complicated, powerful enough, used by automatic tools

The lookahead is different from the one in LL(k) grammars:

LL(k): k = number of tokens we look ahead to do a prediction

LR(k): k = number of tokens we look ahead when we do a reduce operation

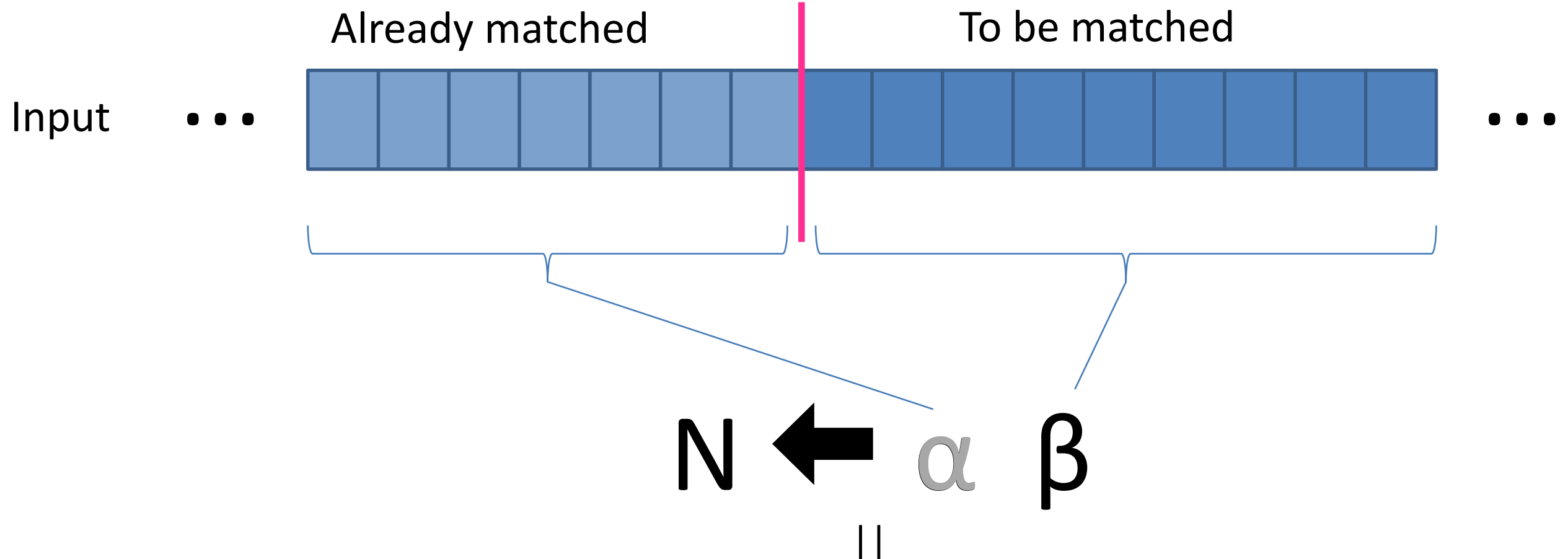
LR(0)

"LR(0) Item"

For a production rule

$N \rightarrow \alpha \beta$

in the grammar,



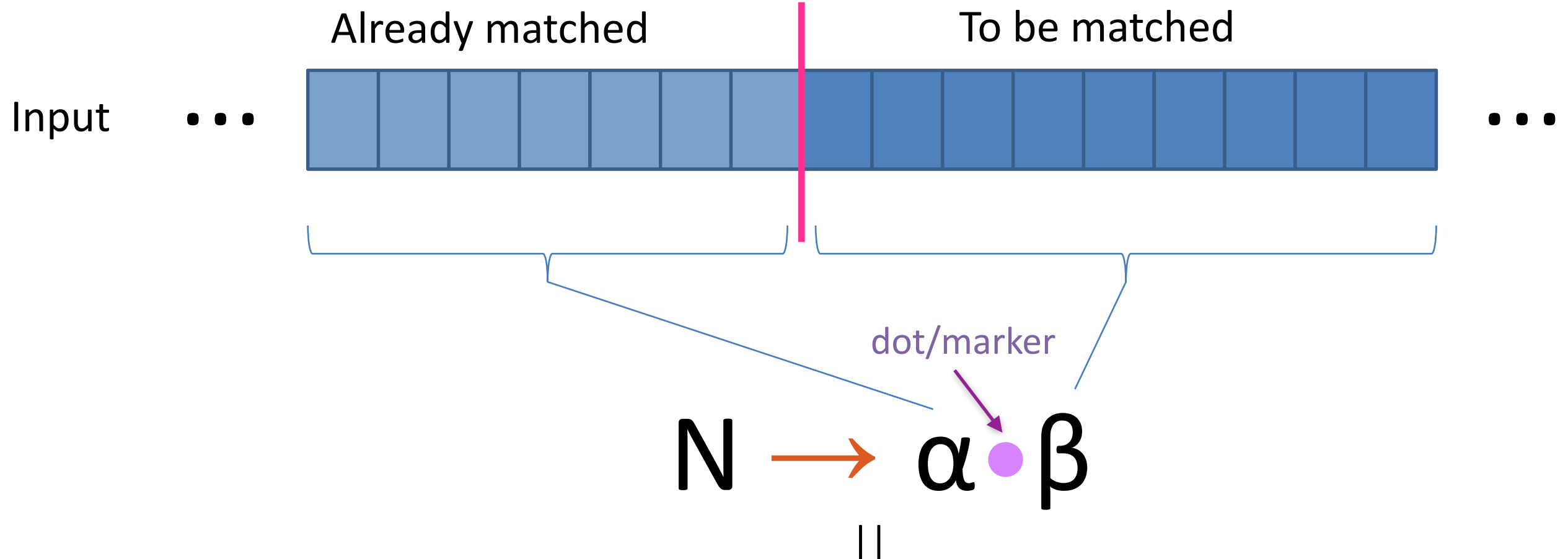
So far we've matched α , expecting to see β

LR(0) Item

For a production rule

$N \rightarrow \alpha \beta$

in the grammar,



So far we've matched α , expecting to see β

LR(0) Item

$S \rightarrow A \mid C$
 $A \rightarrow BC \mid BB \mid \boxed{BbC} \mid CB$
 $B \rightarrow bb$
 $C \rightarrow cc$

$A \rightarrow \bullet BbC$ $A \rightarrow B \bullet bC$ $A \rightarrow Bb \bullet C$

Shift Item

$A \rightarrow BbC \bullet$

Reduce Item

LR(0) Item

$E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

?

Shift Item

?

Reduce Item

LR(0) Item

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

?

Shift Item

?

Reduce Item

Example: Parsing with LR(0) Items

- ▶ **LR(0) state** = set of LR(0) items
- ▶ **LR(0) item** = a derivation rule with a marker
 - ▶ *Separates already determined symbols (past) from the one we hope to determine (future)*

Convention: A single initial rule

Convention: \$ marks end of input

$Z \rightarrow E \$$
 $E \rightarrow T \mid E + T$
 $T \rightarrow i \mid (E)$

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

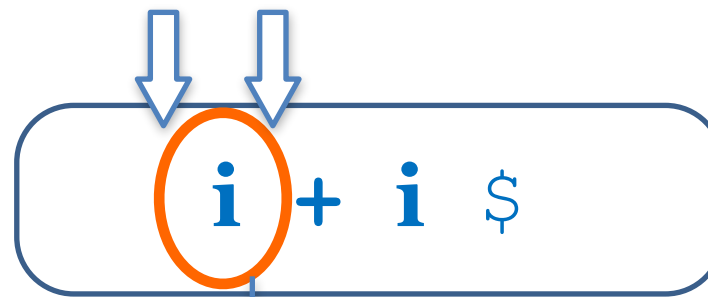
$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

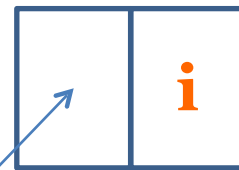
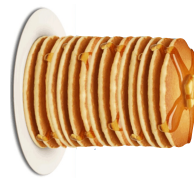
$T \rightarrow \bullet (E)$

Initial state

input



Shift



$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

No need to actually
push states with reduce
items to the stack

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

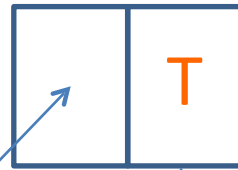
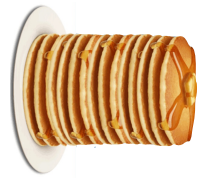
$T \rightarrow \bullet (E)$

$T \rightarrow i \bullet$

Reduce item!

input

$i + i \$$



i

$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

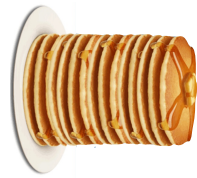
$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input

$i + i \$$



i

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

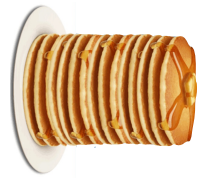
$T \rightarrow \bullet (E)$

$E \rightarrow T \bullet$

Reduce item!

input

i + **i** \$



T
i

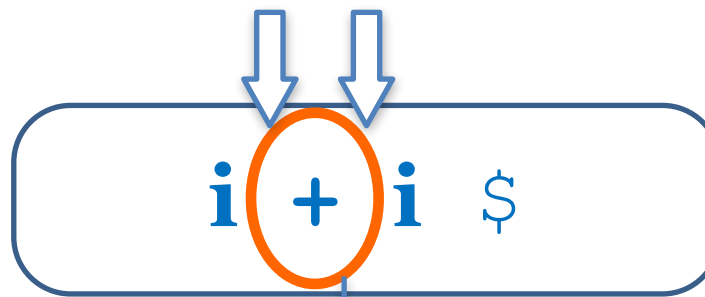
$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \$$

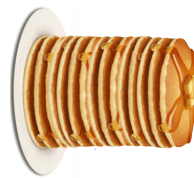
$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input



Shift



T
i

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

$T \rightarrow \bullet (E)$

X

$Z \rightarrow E \bullet \$$

$E \rightarrow E \bullet + T$

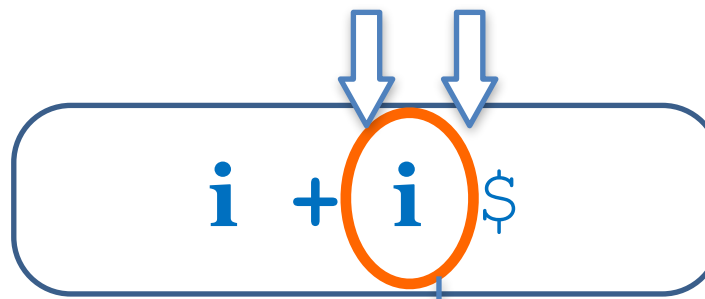
✓

$Z \rightarrow E \$$

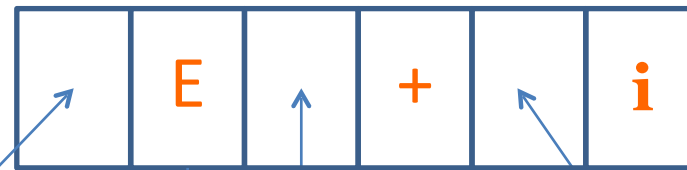
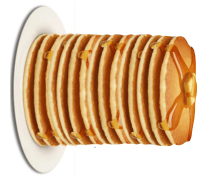
$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input



Shift



T
 i

$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$
 $E \rightarrow E \bullet + T$

$E \rightarrow E + \bullet T$ ✗

 $T \rightarrow \bullet i$ ✓
 $T \rightarrow \bullet (E)$ ✗

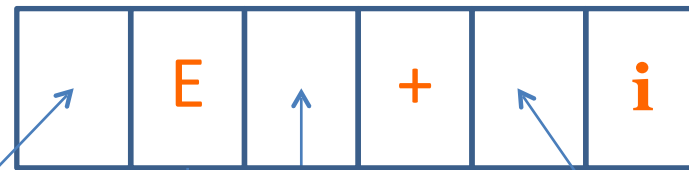
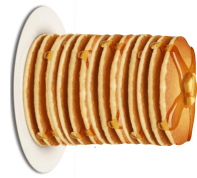
$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input

$i + i \$$



T
—
i

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

Reduce item!

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

$T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$

$E \rightarrow E \bullet + T$

$E \rightarrow E + \bullet T$

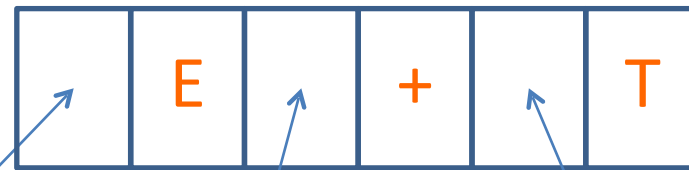
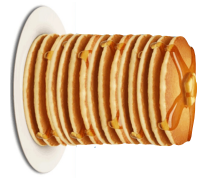
$T \rightarrow \bullet i$

$T \rightarrow \bullet (E)$

$T \rightarrow i \bullet$

input

i + i \$



T

i

i

$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$
 $E \rightarrow E \bullet + T$

$E \rightarrow E + \bullet T$ ✓

 $T \rightarrow \bullet i$ ✗
 $T \rightarrow \bullet (E)$ ✗

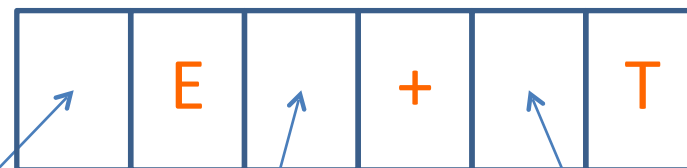
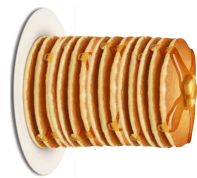
$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input

$i + i \$$



T

i

i

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

Reduce item!

$Z \rightarrow \bullet E \$$

$E \rightarrow \bullet T$

$E \rightarrow \bullet E + T$

$T \rightarrow \bullet i$

$T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$

$E \rightarrow E \bullet + T$

$E \rightarrow E + \bullet T$

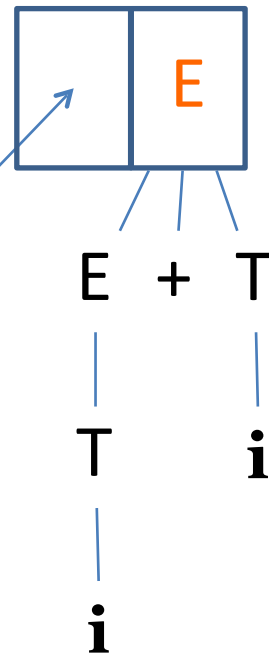
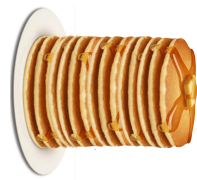
$T \rightarrow \bullet i$

$T \rightarrow \bullet (E)$

$E \rightarrow E + T \bullet$

input

i + i \$



Items in the current state:

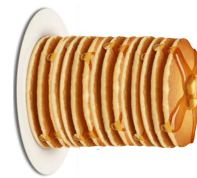
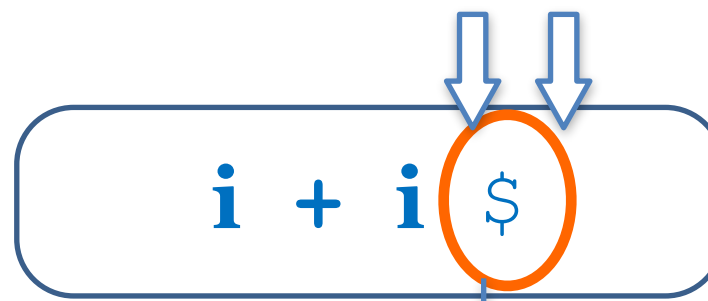
- $Z \rightarrow \bullet E \$$
- $E \rightarrow \bullet T$
- $E \rightarrow \bullet E + T$
- $T \rightarrow \bullet i$
- $T \rightarrow \bullet (E)$

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input



$E + T$
 $T \quad i$
 i

$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$ ✓
 $E \rightarrow E \bullet + T$ ✗

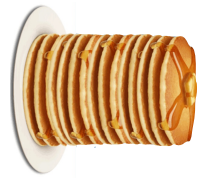
$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

input

i + i \$



E **+** **T**
| |
T **i**
|
i

$Z \rightarrow E \$$

$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

Reduce item!

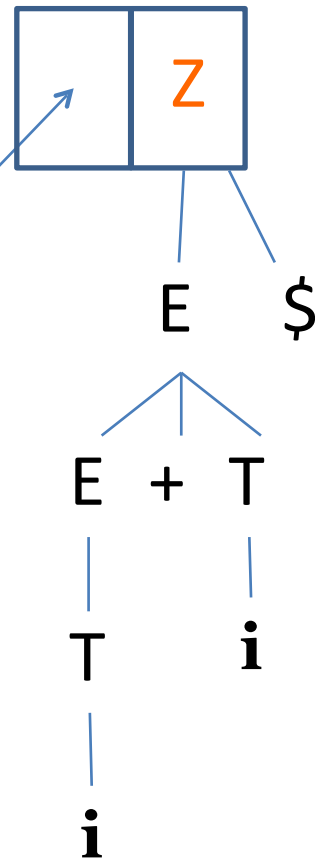
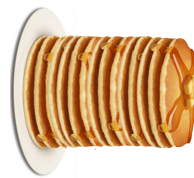
$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \bullet \$$
 $E \rightarrow E \bullet + T$

$Z \rightarrow E \$ \bullet$

input

i + i \$



$Z \rightarrow \bullet E \$$
 $E \rightarrow \bullet T$
 $E \rightarrow \bullet E + T$
 $T \rightarrow \bullet i$
 $T \rightarrow \bullet (E)$

$Z \rightarrow E \$$

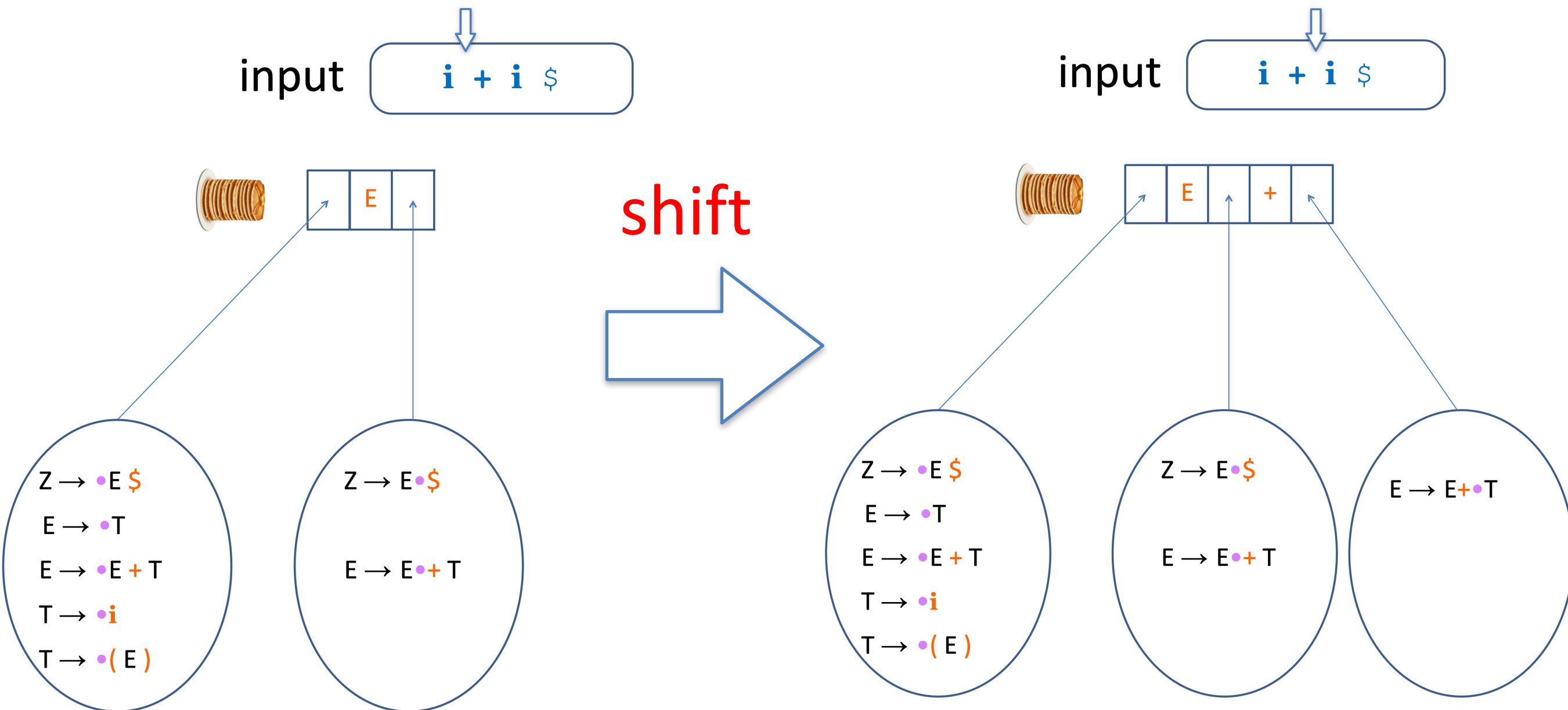
$E \rightarrow T \mid E + T$

$T \rightarrow i \mid (E)$

Reducing the initial rule
means **accept**

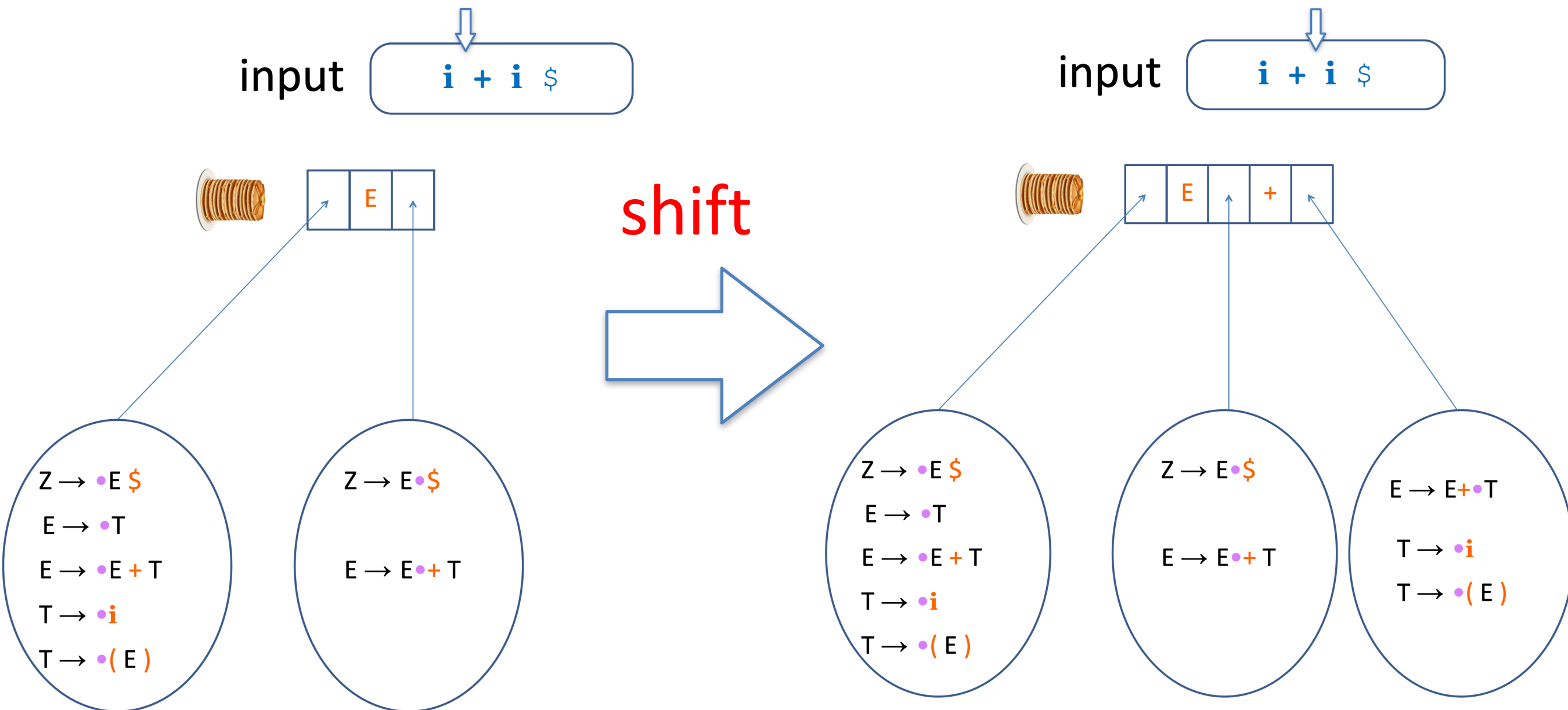
Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?



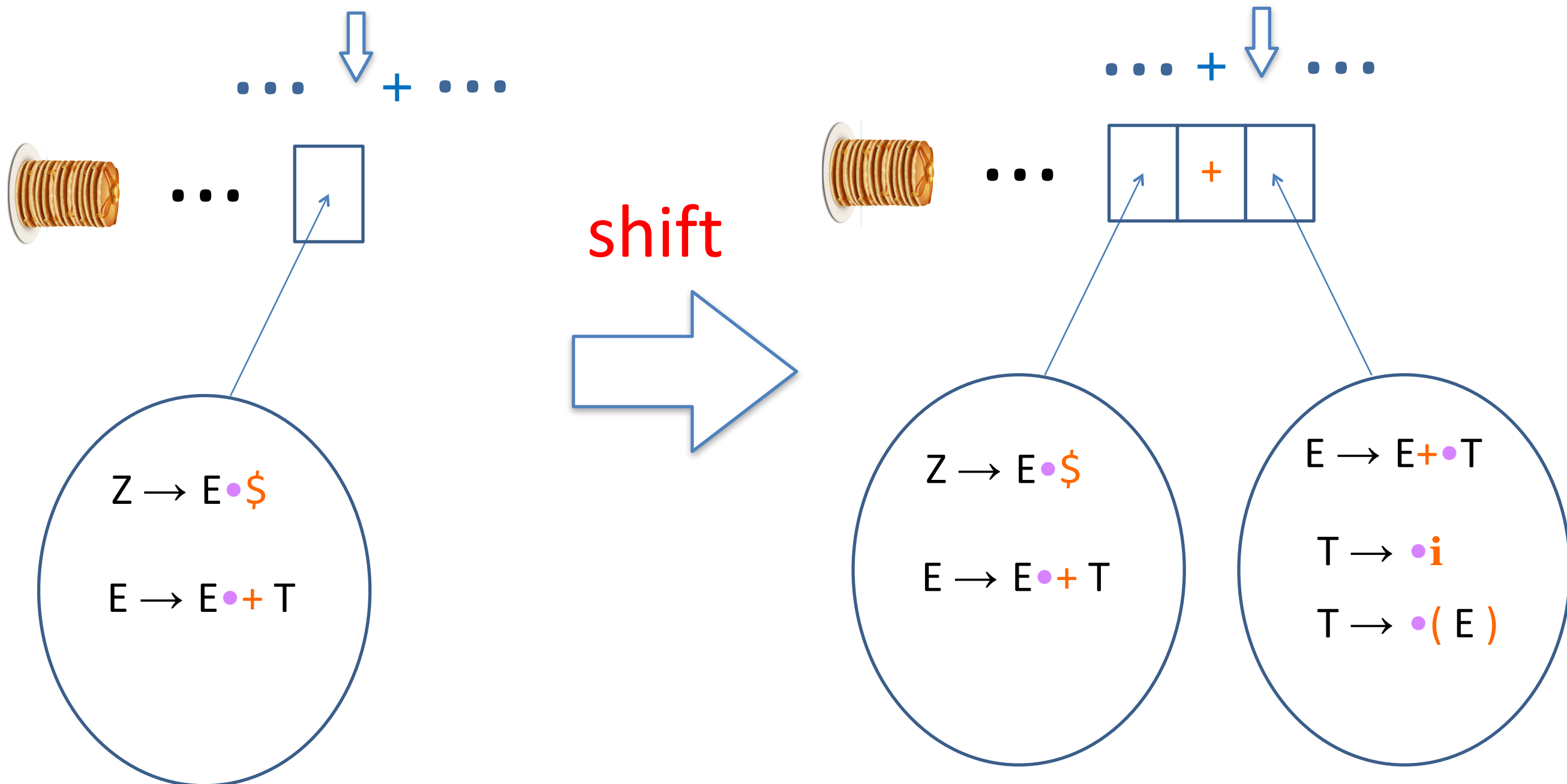
Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?



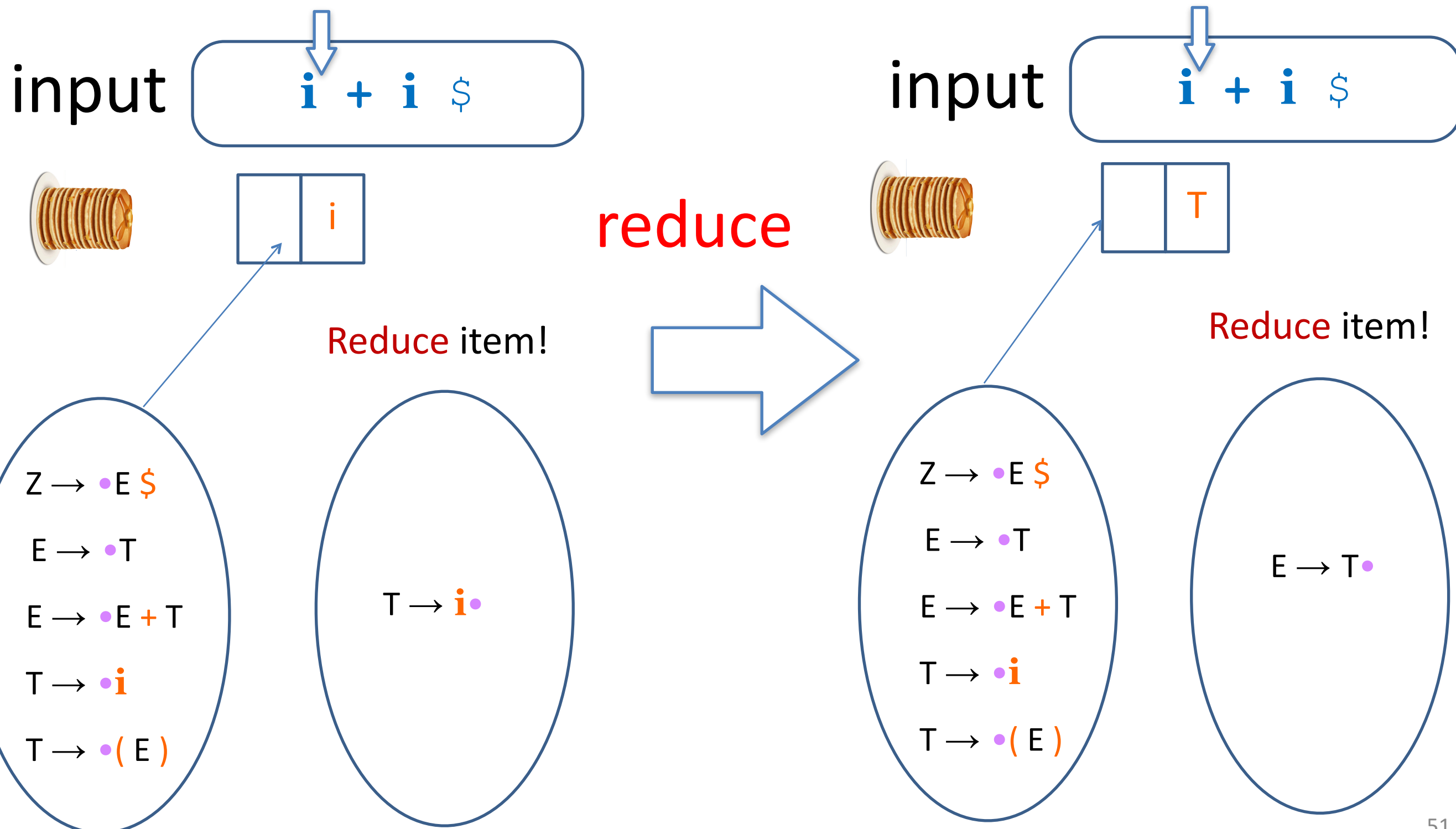
Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?



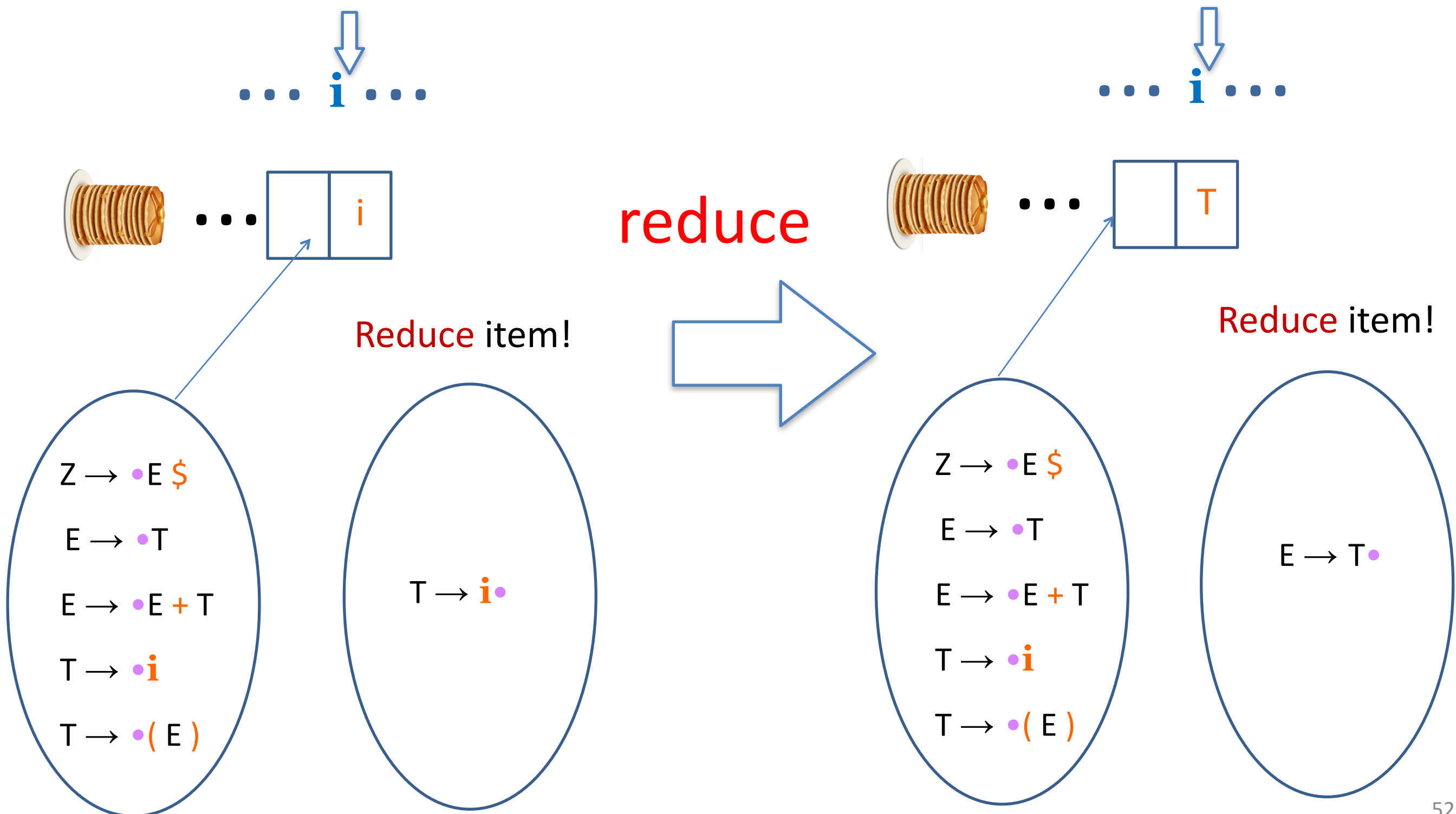
Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?



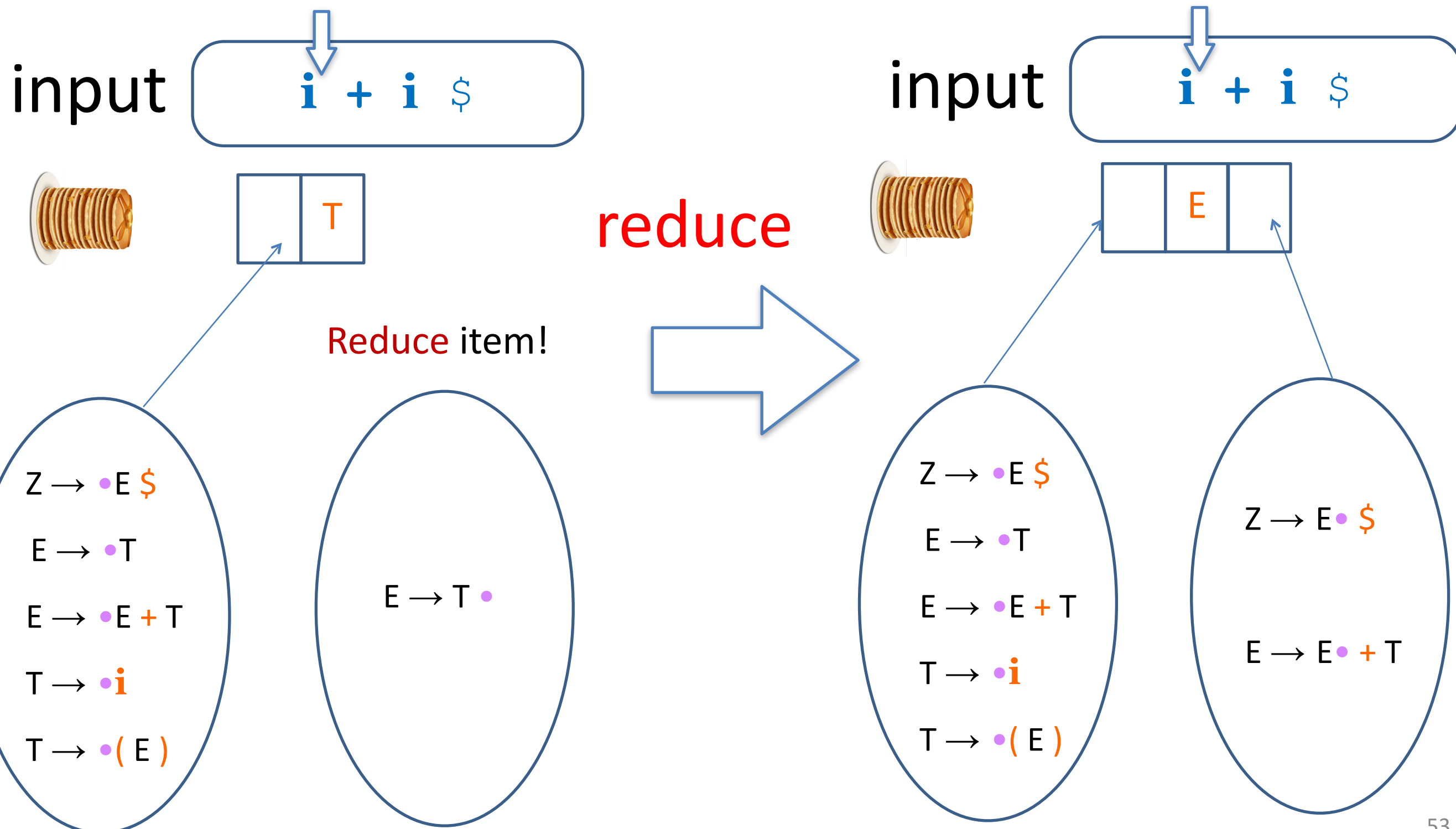
Computing State Transitions

- Do we need to compute the state transitions during parsing?
 - Or can we do it ahead?



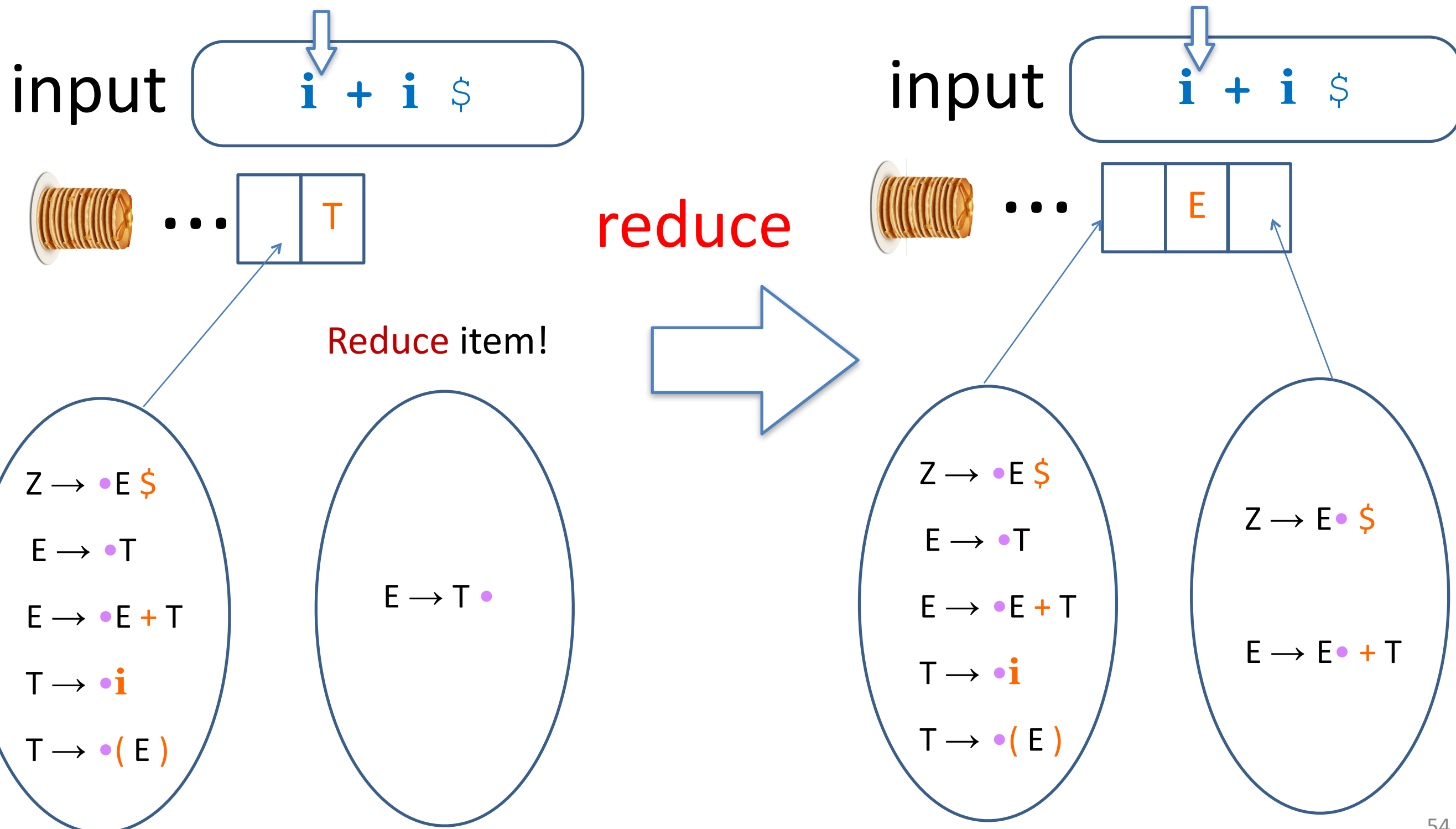
Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?



Computing State Transitions

- Do we need to compute the state transitions during parsing?
- Or can we do it ahead?

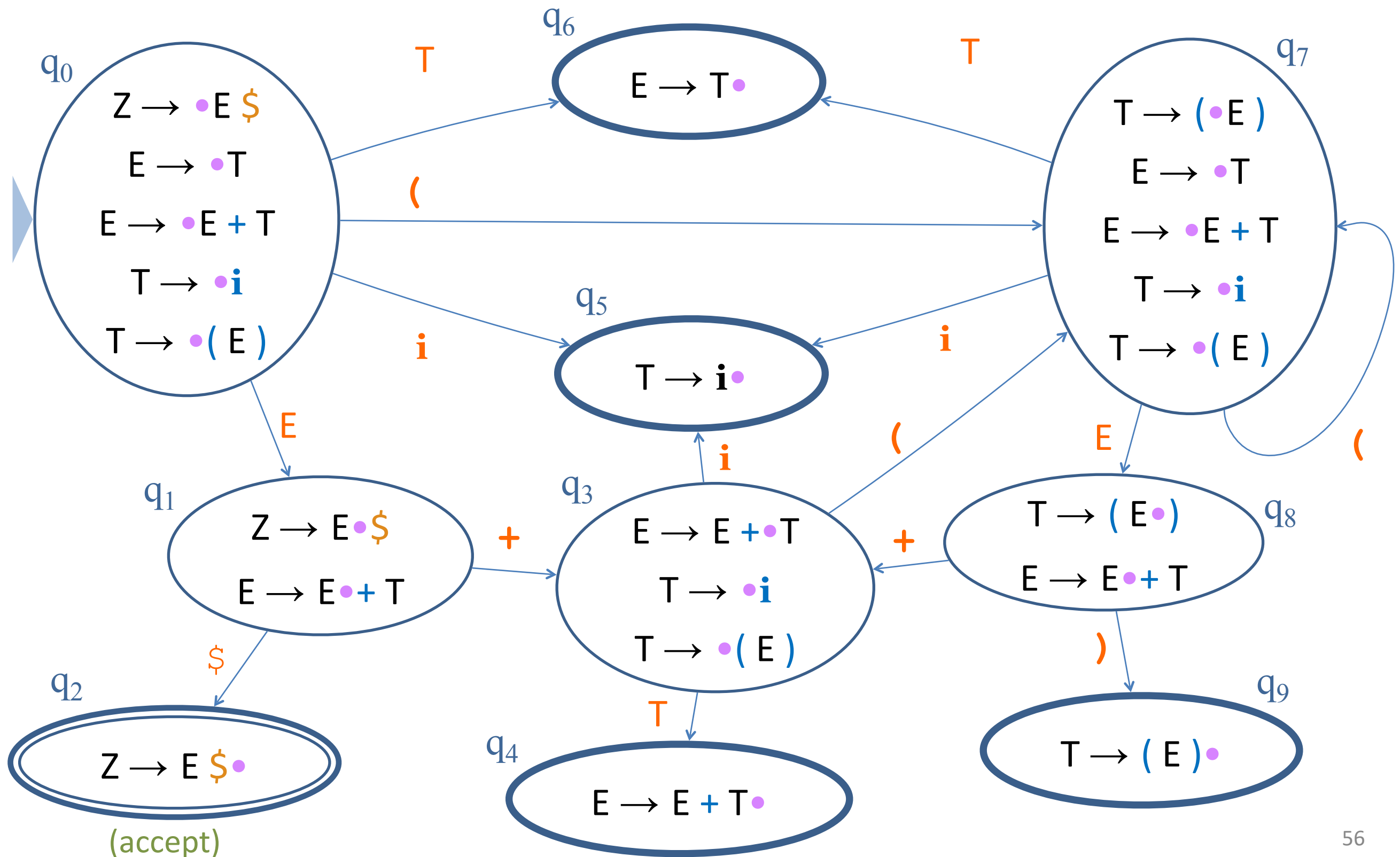


Computing State Transitions

- Do we need to compute the state transitions during parsing?
 - Or can we do it ahead of time?



LR(0) Automaton: All Possible Transitions



How to construct the LR(0) automaton? (A bit more formal)

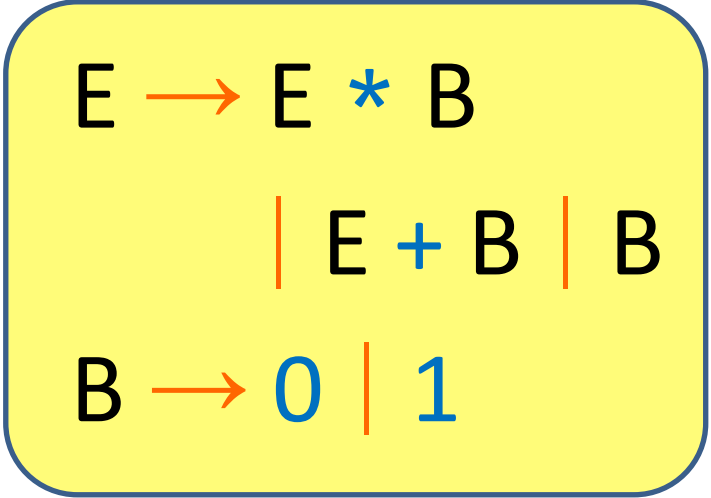
LR(0) Automaton Construction

- An LR(0) state consists of several LR(0) items
- For example, if we identified a string that is reduced to E, then we may be in one of the following LR items:

$$E \rightarrow E \bullet + B \quad \text{or} \quad E \rightarrow E \bullet * B$$

- Therefore one state would be:

$$q = \{E \rightarrow E \bullet + B, E \rightarrow E \bullet * B\}$$



A yellow rounded rectangle containing three LR(0) items. The first item is $E \rightarrow E * B$ with an orange arrow pointing to the asterisk. The second item is $E \rightarrow E + B$ with a blue plus sign, preceded and followed by vertical orange lines. The third item is $B \rightarrow 0 \mid 1$ with an orange arrow pointing to the 0, followed by a vertical orange line and then the 1.

$$\begin{array}{l} E \rightarrow E * B \\ \quad | \quad E + B \quad | \quad B \\ B \rightarrow 0 \mid 1 \end{array}$$

- But if the current state includes $E \rightarrow E + \bullet B$, then we must allow B to be derived too — **Closure!**

Construct the Closure

- Proposition: a **closure set** of LR(0) items has the following property — if the set contains an item of the form

$$A \rightarrow \alpha \bullet B \beta$$

where B is a non terminal, then it must *also* contain an item

$$B \rightarrow \bullet \delta$$

for *each rule* of the form $B \rightarrow \delta$ in the grammar.

- Building the closure set for a given item set is a recursive fixpoint computation, as δ may also begin with a variable

Closure: an example

$E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

grammar

$C = \{ E \rightarrow E + \bullet B \}$

set C of LR items

- The closure of the set C is

$$\text{clos}(C) = \{ E \rightarrow E + \bullet B , \\ B \rightarrow \bullet 0 , \\ B \rightarrow \bullet 1 \}$$

- This will become another parser LR(0) state

Closure: an example

$E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

grammar

$C = \{ E \rightarrow \bullet E + B \}$

set C of LR items

- The closure of the set C is

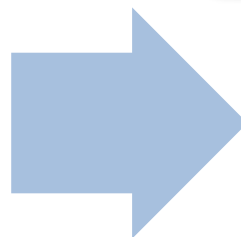
$\text{clos}(C) = \{$

Extended Grammar

- Goal: simple termination condition
 - ▶ Assume that the initial variable only appears in a single rule.
This guarantees that the last reduction can be (easily) detected
 - ▶ Any grammar can be (easily) extended to have such a structure

Example: the grammar

(1) $E \rightarrow E * B$
(2) $E \rightarrow E + B$
(3) $E \rightarrow B$
(4) $B \rightarrow 0$
(5) $B \rightarrow 1$



Can be extended into

Initial
rule

(0) $S \rightarrow E \$$
(1) $E \rightarrow E * B$
(2) $E \rightarrow E + B$
(3) $E \rightarrow B$
(4) $B \rightarrow 0$
(5) $B \rightarrow 1$

end-of-input
marker

The Initial State

- To build the automaton, we go through all possible states that can be seen during reduction
 - ▶ Each state represents a (closure) set of LR(0) items
- The initial state q_0 is the closure of the initial rule
 - ▶ In our example, the initial rule is $S \rightarrow \bullet E \$$, so:

$$q_0 = \text{clos}(\{S \rightarrow \bullet E \$\}) =$$

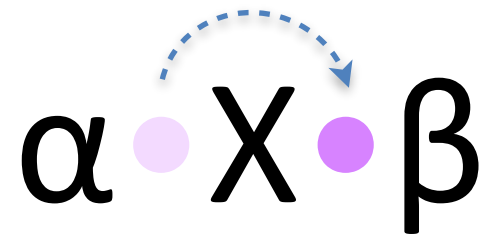
$$\{ S \rightarrow \bullet E \$, E \rightarrow \bullet E * B, E \rightarrow \bullet E + B, \\ E \rightarrow \bullet B, B \rightarrow \bullet 0, B \rightarrow \bullet 1 \}$$

$S \rightarrow E \$$
 $E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

- Next, we build all possible states that can be reached by following a *single symbol* (token or variable)

The Next States

- For every symbol (terminal or variable) X , and every possible state (closure set) q ,
 - Find all items in the set of q in which the dot is before an X .
(We denote this set by $q|X$)
 - Move the dot ahead of the X in all items in $q|X$
 - Find the closure of the obtained set:
this is the state into which we move from q upon seeing X



- Formally, the **next set** is defined on a set C and next symbol X :

- $C|X = \{ (N \rightarrow \alpha \bullet X \beta) \in C \}$
- $\text{step}(C, X) = \{ N \rightarrow \alpha X \bullet \beta \mid (N \rightarrow \alpha \bullet X \beta) \in C \}$
- $\text{nextSet}(C, X) = \text{clos}(\text{step}(C, X))$

The Next States

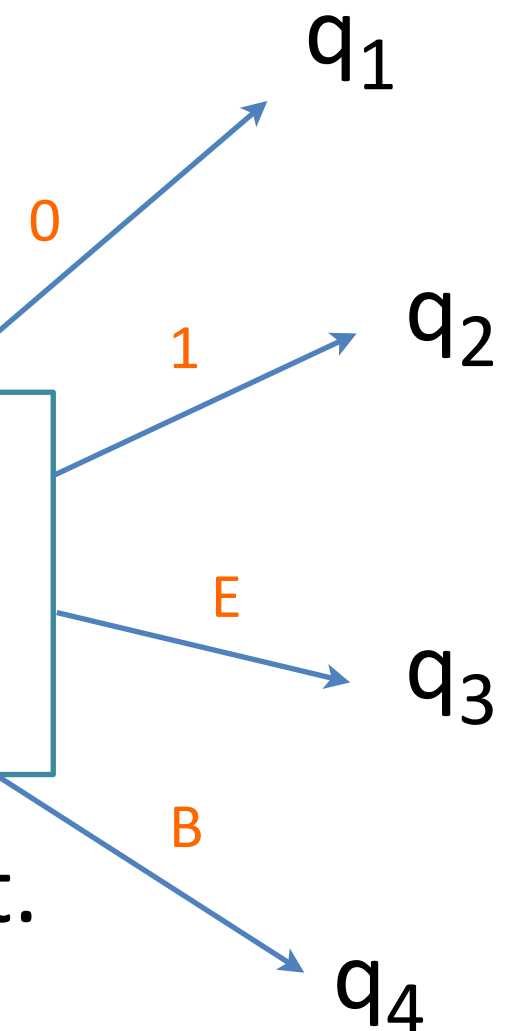
$S \rightarrow E \$$
 $E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

Recall that in our example

$$q_0 = \text{clos}(\{S \rightarrow \bullet E\}) =$$

$\{S \rightarrow \bullet E \$, E \rightarrow \bullet E * B, E \rightarrow \bullet E + B,$
 $E \rightarrow \bullet B, B \rightarrow \bullet 0, B \rightarrow \bullet 1\}$

Let us check which states are reachable from it.



nextSet(q_0, \cdot) in the example

$S \rightarrow E \$$
 $E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

$S \rightarrow \bullet E \$$
 $E \rightarrow \bullet E * B$
 $E \rightarrow \bullet E + B$
 $E \rightarrow \bullet B$
 $B \rightarrow \bullet 0 \mid B \rightarrow \bullet 1$

$$q_0 | 0 = \{B \rightarrow \bullet 0\}$$

$$q_0 | 1 = \{B \rightarrow \bullet 1\}$$

q_1

$$q_1 = \text{clos}(\{B \rightarrow 0 \bullet\}) = \{B \rightarrow 0 \bullet\}$$

q_2

$$q_2 = \text{clos}(\{B \rightarrow 1 \bullet\}) = \{B \rightarrow 1 \bullet\}$$

q_0

0

1

E

B

$$q_0 | E = \{S \rightarrow \bullet E \$, E \rightarrow \bullet E * B, E \rightarrow \bullet E + B\}$$

$$q_0 | B = \{E \rightarrow \bullet B\}$$

q_3

$$q_3 = \text{clos}(\{S \rightarrow E \bullet \$, E \rightarrow E \bullet * B, E \rightarrow E \bullet + B\}) = \{S \rightarrow E \bullet \$, E \rightarrow E \bullet * B, E \rightarrow E \bullet + B\}$$

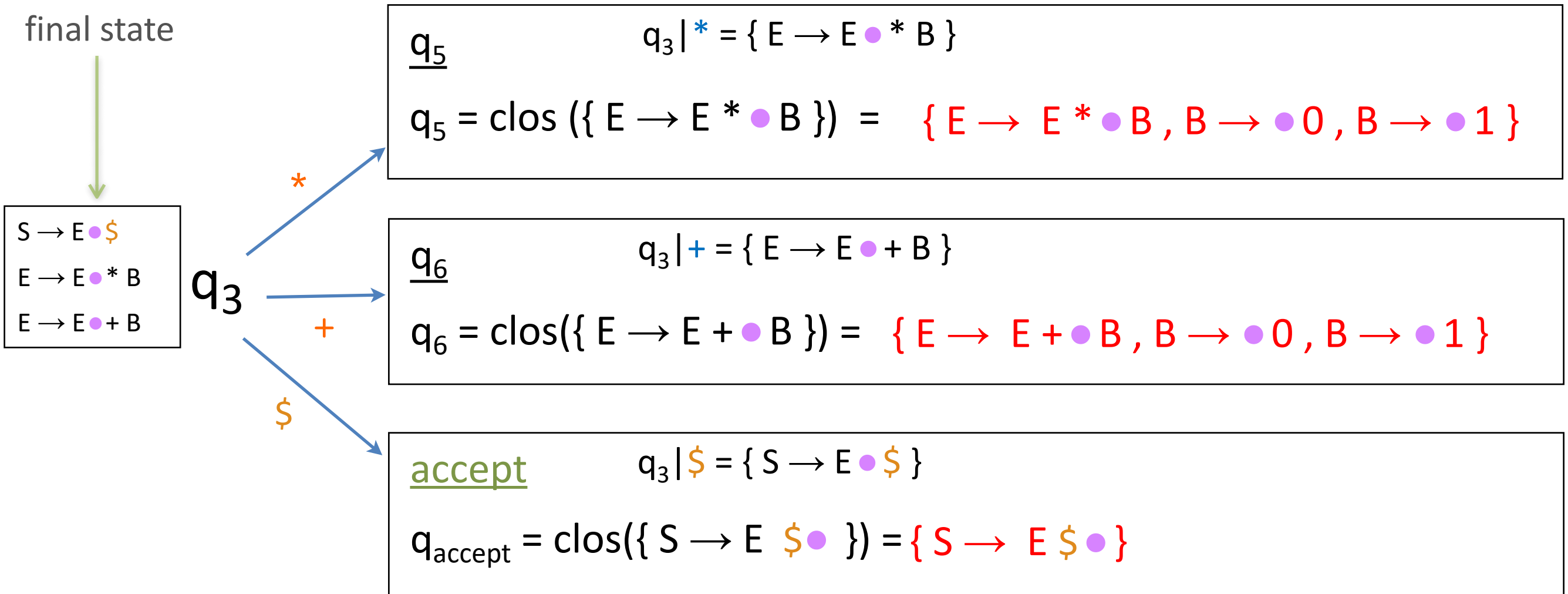
q_4

$$q_4 = \text{clos}(\{E \rightarrow B \bullet\}) = \{E \rightarrow B \bullet\}$$

From these new states there are more reachable states

$S \rightarrow E \$$
 $E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

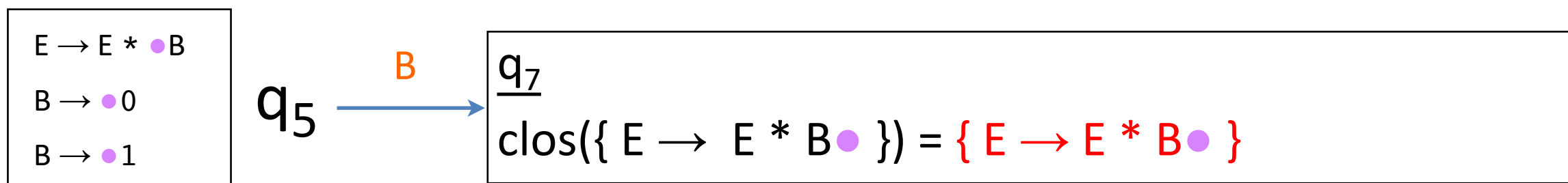
- From q_1, q_2, q_4 , there are no steps because the dot is at the end of every item in their sets.
- From state q_3 we can reach the following three states —



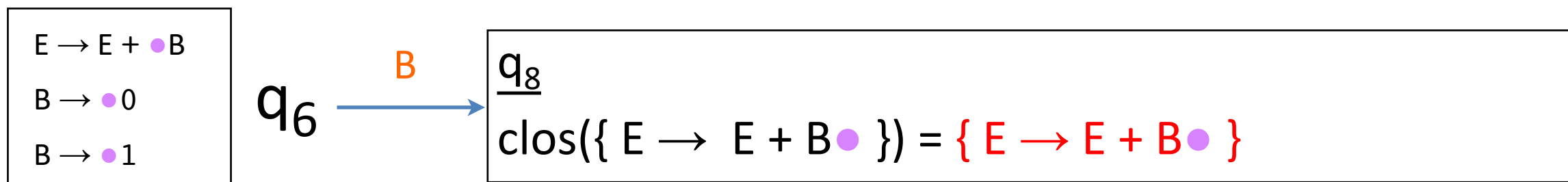
Finally

$S \rightarrow E \$$
 $E \rightarrow E * B \mid E + B \mid B$
 $B \rightarrow 0 \mid 1$

- From q_5 we can proceed with $X=0$, or $X=1$, or $X=B$.
- For $X=0$ we reach q_1 again; for $X=1$ we reach q_2 .
- For $X=B$ we get q_7 :

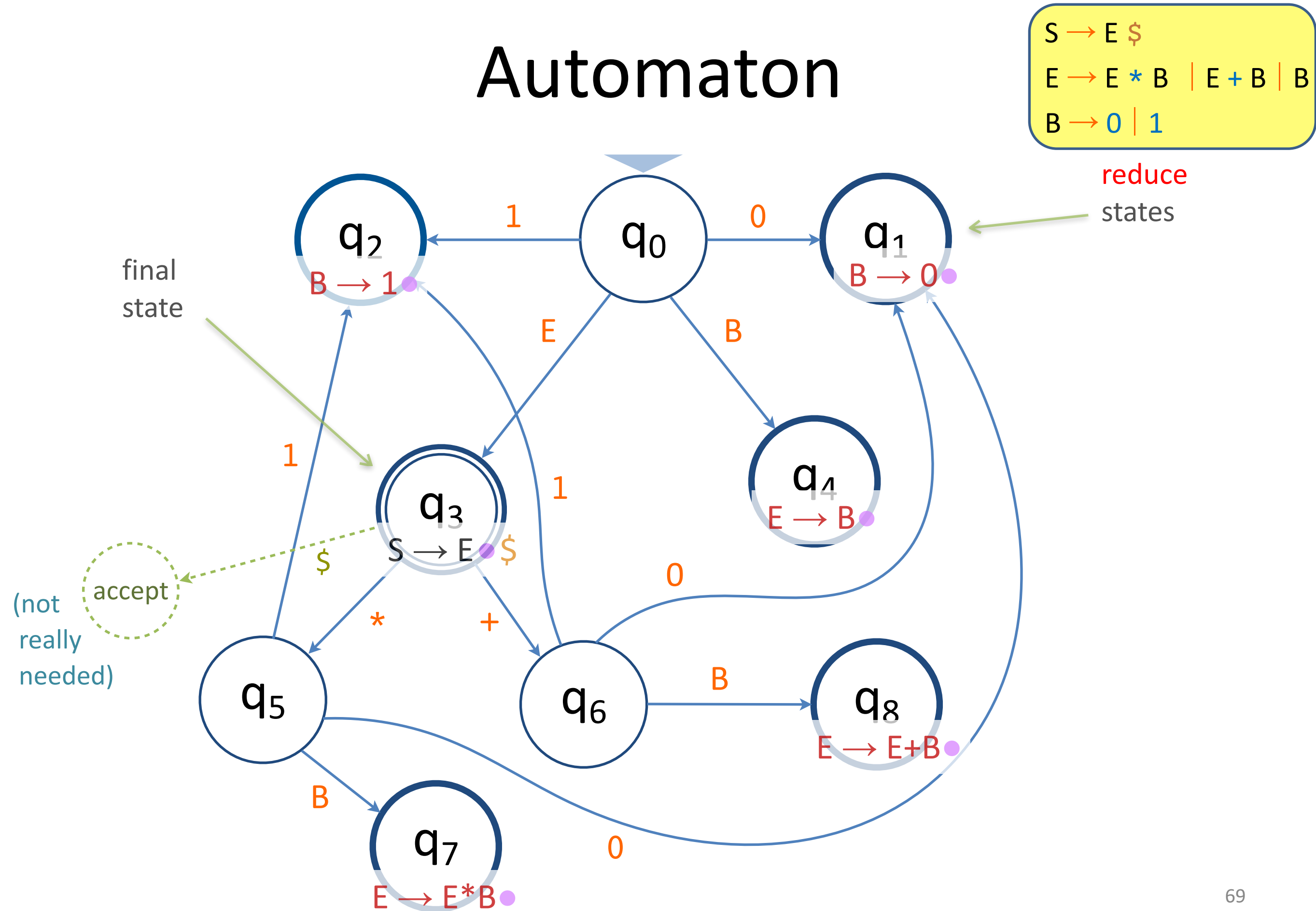


- Similarly, from q_6 with $X=B$ we get q_8 :



- These two states have no further steps. (Why?)

Automaton



Another LR(0) Automaton

$Z \rightarrow E \$$
 $E \rightarrow T \mid E + T$
 $T \rightarrow i \mid (E)$

