

Compiler

Parsing: Bottom-Up

© Copyright 2005, Matthew B. Dwyer and Robby. The syllabus and lectures for this course are copyrighted materials and may not be used in other course settings outside University of Nebraska-Lincoln and Kansas State University in their current form or modified form without express written permission of one of the copyright holders. During this course, students are prohibited from selling notes to or being paid for taking notes by any person or commercial firm without the express written permission of one of the copyright holders.

An Introductory Example

- Bottom-up parsers neither fail on left-recursion nor need left-factored grammars
- Hence we can revert to the “natural” grammar for our example:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T^* \text{ int} \mid \text{int} \mid (E)$$

- Consider the string: $\text{int}^* \text{ int} + \text{int}$

The Idea

Bottom-up parsing *reduces* a string to the start symbol by inverting productions:

int * int + int

$T \rightarrow \text{int}$

$T^* \text{ int } + \text{ int}$

$T \rightarrow T^* \text{ int}$

$T + \text{ int}$

$E \rightarrow T$

$E + \text{ int}$

$T \rightarrow \text{int}$

$E + T$

$E \rightarrow E + T$

E

Observation

- Read the productions in this bottom-up parse in reverse (i.e., from bottom to top)
- This is a rightmost derivation!

int * int + int

$T \rightarrow \text{int}$

$T^* \text{ int } + \text{ int}$

$T \rightarrow T^* \text{ int}$

$T + \text{ int}$

$E \rightarrow T$

$E + \text{ int}$

$T \rightarrow \text{int}$

$E + T$

$E \rightarrow E + T$

E

A Bottom-up Parse

int * int + int

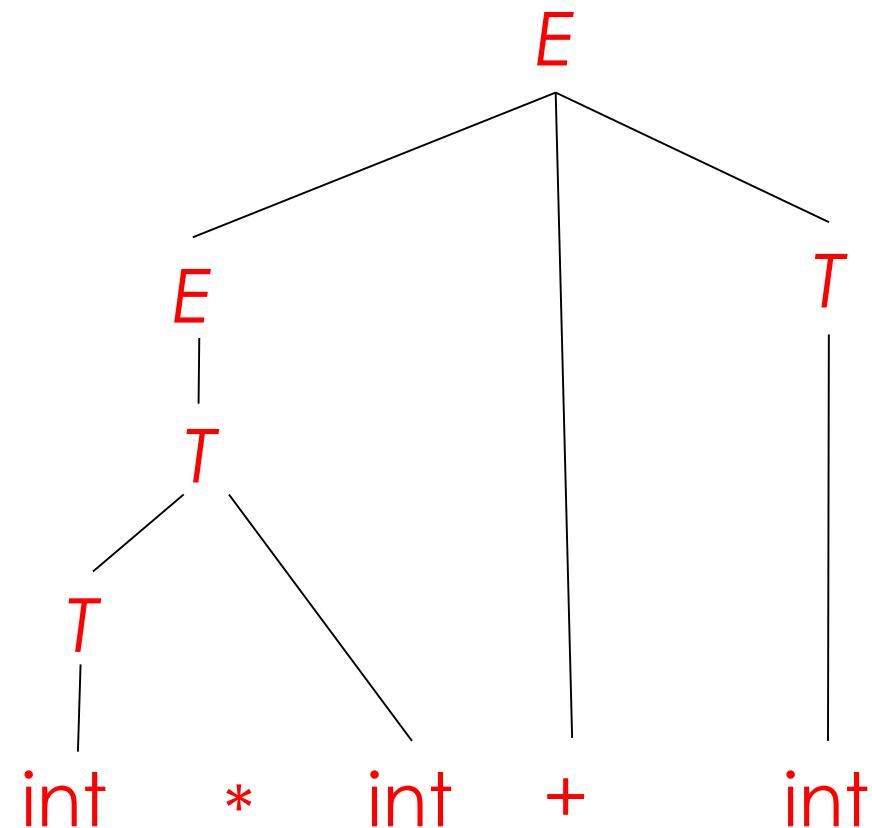
$T^* \text{ int } + \text{ int}$

$T + \text{ int}$

$E + \text{ int}$

$E + T$

E



A bottom-up parser traces a rightmost derivation in reverse

A Bottom-up Parse in Detail (1)

int * int + int

int * int + int

A Bottom-up Parse in Detail (2)

int * int + int

T^* int + int

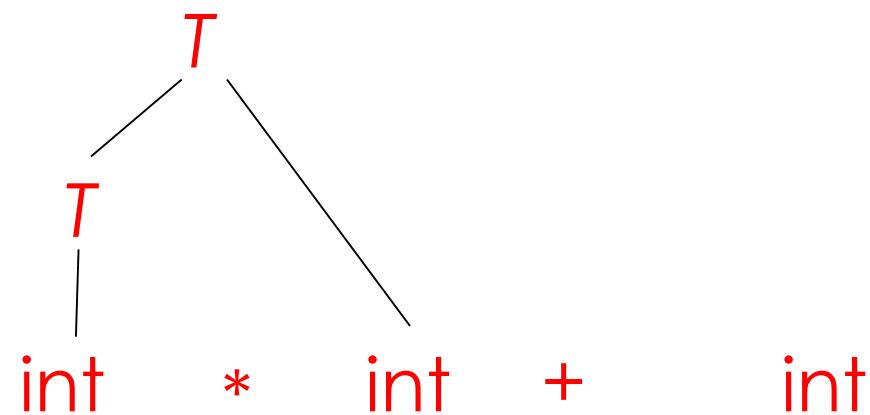


A Bottom-up Parse in Detail (3)

int * int + int

$T^* \text{ int } + \text{ int}$

$T + \text{ int}$



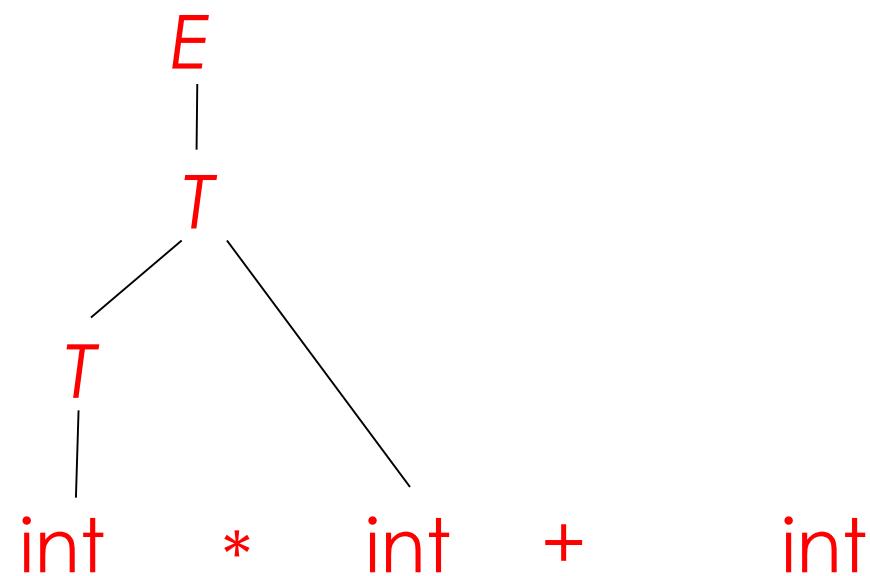
A Bottom-up Parse in Detail (4)

int * int + int

$T^* \text{ int } + \text{ int}$

$T + \text{ int}$

$E + \text{ int}$



A Bottom-up Parse in Detail (5)

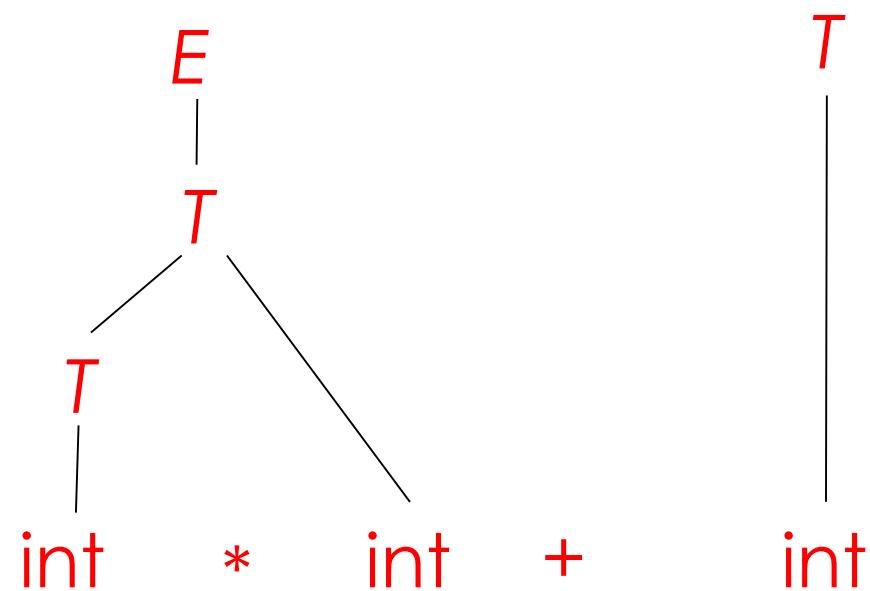
int * int + int

$T^* \text{ int } + \text{ int}$

$T + \text{ int}$

$E + \text{ int}$

$E + T$



A Bottom-up Parse in Detail (6)

int * int + int

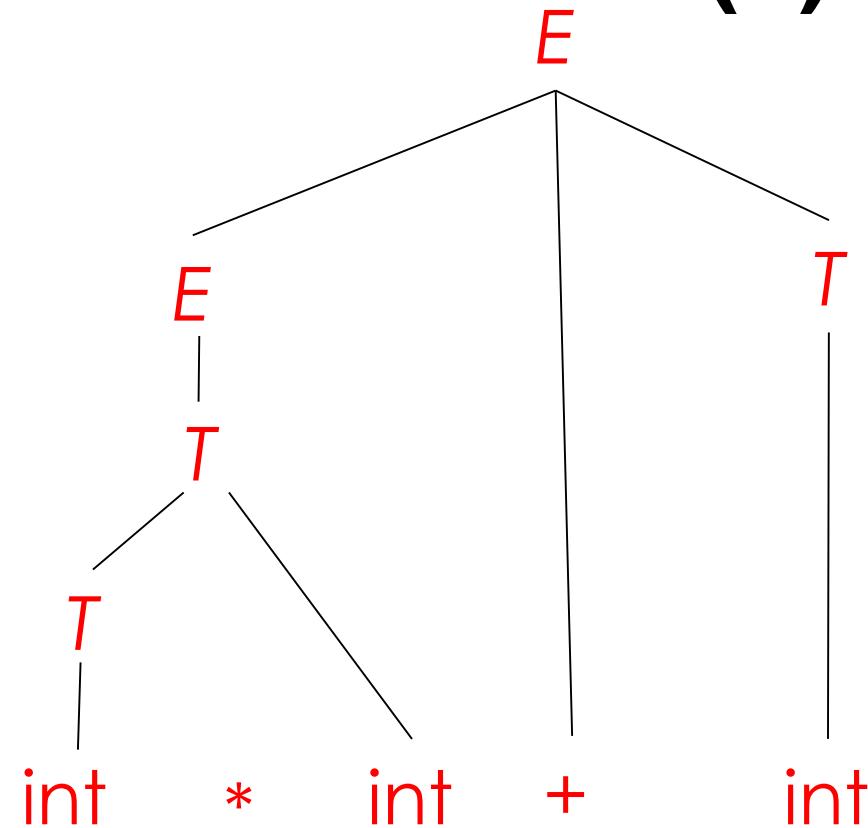
$T^* \text{ int } + \text{ int}$

$T + \text{ int}$

$E + \text{ int}$

$E + T$

E



A Trivial Bottom-Up Parsing Algorithm

Let \mathcal{I} = input string

repeat

 pick a non-empty substring β of \mathcal{I}

 where $X \rightarrow \beta$ is a production

 if no such β , backtrack

 replace one β by X in \mathcal{I}

until $\mathcal{I} = S$ (the start symbol) or all
possibilities are exhausted

For You To Do

Do you see any problems with this algorithm?

Think about performance and completeness!

Where Do Reductions Happen

Let $\alpha\beta\omega$ be a step of a bottom-up parse

- Assume the next reduction is by $X \rightarrow \beta$
- Then ω is a string of terminals, i.e. $\omega \in T^*$

Why?

Because $\alpha X \omega \Rightarrow \alpha \beta \omega$ is a step in a right-most derivation

Idea

- Split string into two substrings
 - Right substring : as yet unexamined by parsing (a string of terminals)
 - Left substring : has terminals and non-terminals
- The dividing point is marked by a |
 - The | is not part of the string
- Initially, all input is unexamined | $x_1 x_2 \dots x_n$

Shift-Reduce Parsing

Bottom-up parsing uses only two kinds of actions:

Shift

Reduce

Shift

- Move | one place to the right
 - Shifts a terminal to the left string

$$ABC|xyz \Rightarrow ABCx|yz$$

Reduce

- Apply an *inverse production* at the right end of the left string
 - If $A \rightarrow xy$ is a production, then

$$Cbxy|ijk \Rightarrow CbA|ijk$$

- xy is called a *handle*

The Example with Shift-Reduce Parsing

|int * int + int

The Example with Shift-Reduce Parsing

|int * int + int shift
int | * int + int

The Example with Shift-Reduce Parsing

|int * int + int

shift

int | * int + int

reduce $T \rightarrow \text{int}$

T | * int + int

The Example with Shift-Reduce Parsing

|int * int + int

shift

int | * int + int

reduce $T \rightarrow \text{int}$

T | * int + int

shift

T^* | int + int

The Example with Shift-Reduce Parsing

|int * int + int

shift

int | * int + int

reduce $T \rightarrow \text{int}$

T | * int + int

shift

T^* | int + int

shift

T^* int | + int

The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
T * int + int	shift
T^* int + int	shift
T^* int + int	reduce $T \rightarrow T^* \text{int}$
T + int	

The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
T * int + int	shift
T^* int + int	shift
T^* int + int	reduce $T \rightarrow T^* \text{int}$
T + int	reduce $E \rightarrow T$
E + int	

The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
T * int + int	shift
T^* int + int	shift
T^* int + int	reduce $T \rightarrow T^* \text{ int}$
T + int	reduce $E \rightarrow T$
E + int	shift, shift
E + int	

The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
T * int + int	shift
T^* int + int	shift
T^* int + int	reduce $T \rightarrow T^* \text{ int}$
T + int	reduce $E \rightarrow T$
E + int	shift, shift
E + int	reduce $T \rightarrow \text{int}$
E + T	

The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
T * int + int	shift
T^* int + int	shift
T^* int + int	reduce $T \rightarrow T^* \text{ int}$
T + int	reduce $E \rightarrow T$
E + int	shift, shift
E + int	reduce $T \rightarrow \text{int}$
E + T	reduce $E \rightarrow E+T$
E	

The Stack

- Left string can be implemented by a stack
 - Top of the stack is the |
- Shift pushes a terminal on the stack
- Reduce pops 0 or more symbols off of the stack (production rhs) and pushes a non-terminal on the stack (production lhs)
- This is how a Pushdown Automaton works.

Key Issue

- How do we decide when to shift or reduce?
 - Consider step $T^* \text{ int } | + \text{ int}$
 - We could reduce by $T \rightarrow \text{ int}$ giving $T^* T | + \text{ int}$
 - A fatal mistake: No way to reduce to the start symbol E
- This is resolved by various bottom-up parsing algorithms

Conflicts

- But what if there is a choice?
 - If it is legal to shift or reduce, there is a *shift-reduce conflict*
 - If it is legal to reduce by two different productions, there is a *reduce-reduce conflict*
- Generic shift-reduce strategy:
 - If there is a handle on top of the stack, reduce
 - Otherwise, shift

Source of Conflicts

- Ambiguous grammars always cause conflicts
- But beware, so do many non-ambiguous grammars

Conflict Example

Consider our favorite ambiguous grammar:

$$\begin{array}{c} E \quad \rightarrow \\ | \quad E + E \\ | \quad E^* E \\ | \quad (E) \\ | \quad \text{int} \end{array}$$

One Shift-Reduce Parse

int * int + int	shift
...	...
$E^* E \mid + \text{int}$	reduce $E \rightarrow E^* E$
$E \mid + \text{int}$	shift
$E + \mid \text{int}$	shift
$E + \text{int} \mid$	reduce $E \rightarrow \text{int}$
$E + E \mid$	reduce $E \rightarrow E + E$
$E \mid$	

Another Shift-Reduce Parse

int * int + int	shift
...	...
$E^* E + int$	shift
$E^* E + int$	shift
$E^* E + int $	reduce $E \rightarrow int$
$E^* E + E $	reduce $E \rightarrow E + E$
$E^* E $	reduce $E \rightarrow E^* E$
$E $	

Example Notes

- In the second step $E^* E | + \text{int}$ we can either shift or reduce by $E \rightarrow E^* E$
- Choice determines priority of $+$ and $*$
- As noted previously, grammar can be rewritten to enforce precedence
- Precedence declarations are an alternative

Precedence Declarations Revisited

- Precedence declarations cause shift-reduce parsers to resolve conflicts in certain ways
- Declaring “ $*$ ” has greater precedence than $+$ causes parser to reduce at $E^* E \mid + \text{int}$
- More precisely, precedence declaration is used to resolve conflict between reducing a $*$ and shifting a $+$ (and vice versa)

Precedence Declarations Revisited

- The term “precedence declaration” is a bit misleading
- These declarations do not define precedence; they define conflict resolutions
 - Not quite the same thing!

Using Lookahead Strings

- Use lookahead for resolving conflicts
- Problem: *reduction history* remains disregarded
- Example: Consider the **unambiguous** grammar

$$S \rightarrow aA \mid Ba$$
$$A \rightarrow Ba$$
$$B \rightarrow b$$

- For input *ba* we need to reduce *Ba* by $S \rightarrow Ba$.
 - For input *aba* we need to reduce *Ba* by $A \rightarrow Ba$.
 - But the lookahead string is the same (ϵ) in both cases.
- Solution: Encode the *reduction history* in the stack.

LR(1) Parsing

- Use lookahead for resolving conflicts
- DFA for memorizing *reduction history*
- Introduce new start symbol S' and unique production $S' \rightarrow S$.
(Reason: S' does not occur in any other production)

LR(1) Parsing

■ NFA

- States are LR(1)-items ($A \rightarrow \alpha \bullet \beta, u$)
where $A \rightarrow a\beta \in P$, α is parsed prefix of handle
 $u \in T \cup \{\epsilon\}$ is lookahead expected behind $a\beta$

- Transitions:

- $(A \rightarrow \alpha \cdot X\beta, u) \xrightarrow{X} (A \rightarrow \alpha X \cdot \beta, u) \quad X \in T \cup N$
- $(A \rightarrow \alpha \cdot B\beta, u) \xrightarrow{\epsilon} (B \rightarrow \cdot \gamma, v) \quad B \in N$
 $B \rightarrow \gamma \in P$
 $v \in \text{First}(\beta u)$

LR(1) Parsing

■ NFA

- Start state ($S' \rightarrow \bullet S, \varepsilon$)
- Accepting state ($S' \rightarrow S^\bullet, \varepsilon$)

■ DFA

- NFA → DFA conversion (states = sets of LR(1)-items)
- Actions associated with states and lookahead
 - $\text{Act}(q, t) = \text{shift}$ iff $(A \rightarrow \alpha \bullet t \beta, u) \in q, t \in T$
 - $\text{Act}(q, u) = (A \rightarrow \beta)$ iff $(A \rightarrow \beta \bullet, u) \in q, u \in T \cup \{\varepsilon\}$

LR(1) Parsing

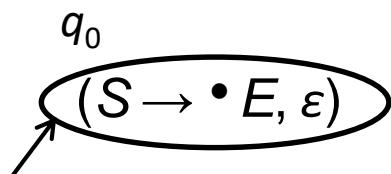
■ Parsing algorithm

- Stack symbols: States (Sets of LR(1)-items)
- Initial stack is $q_0 \mid \omega_c$ where ω_c is the complete input
- Parse steps
 - $\rho q \mid t \omega \rightarrow \rho q \delta(q, t) \mid \omega$ iff $\text{Act}(q, t) = \text{shift}$
 - $\rho q \rho' q' \mid \omega \rightarrow \rho q \delta(q, A) \mid \omega$ iff $\text{Act}(q', u) = (A \rightarrow \beta)$
 $|\beta| = |\rho' q'|$
 $u = \text{First}(\omega)$
- Accept with stack $q_f \mid \varepsilon$ where q_f is accepting state

LR(1) Parsing

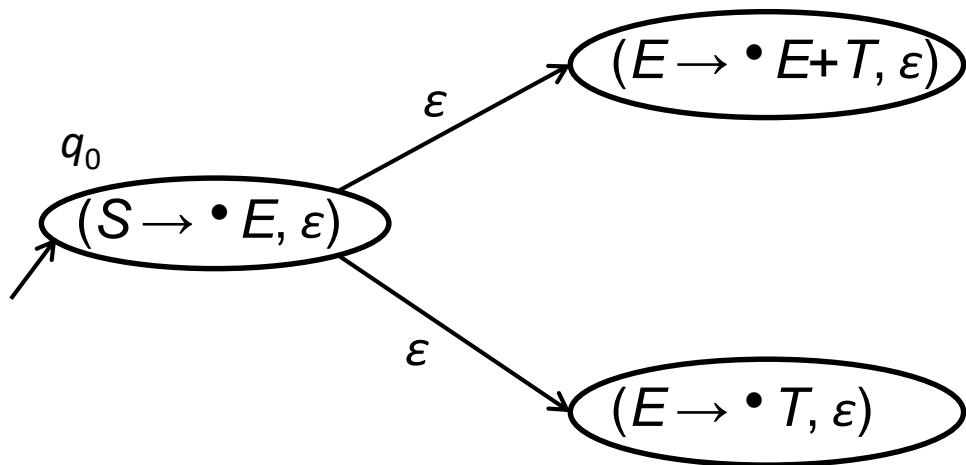
- Example $E \rightarrow E + T \mid T$
 $T \rightarrow T^* \text{ int} \mid \text{int}$
- With new start symbol
 - $S \rightarrow E$
 - $E \rightarrow E + T \mid T$
 - $T \rightarrow T^* \text{ int} \mid \text{int}$

Example: LR(1) NFA Construction



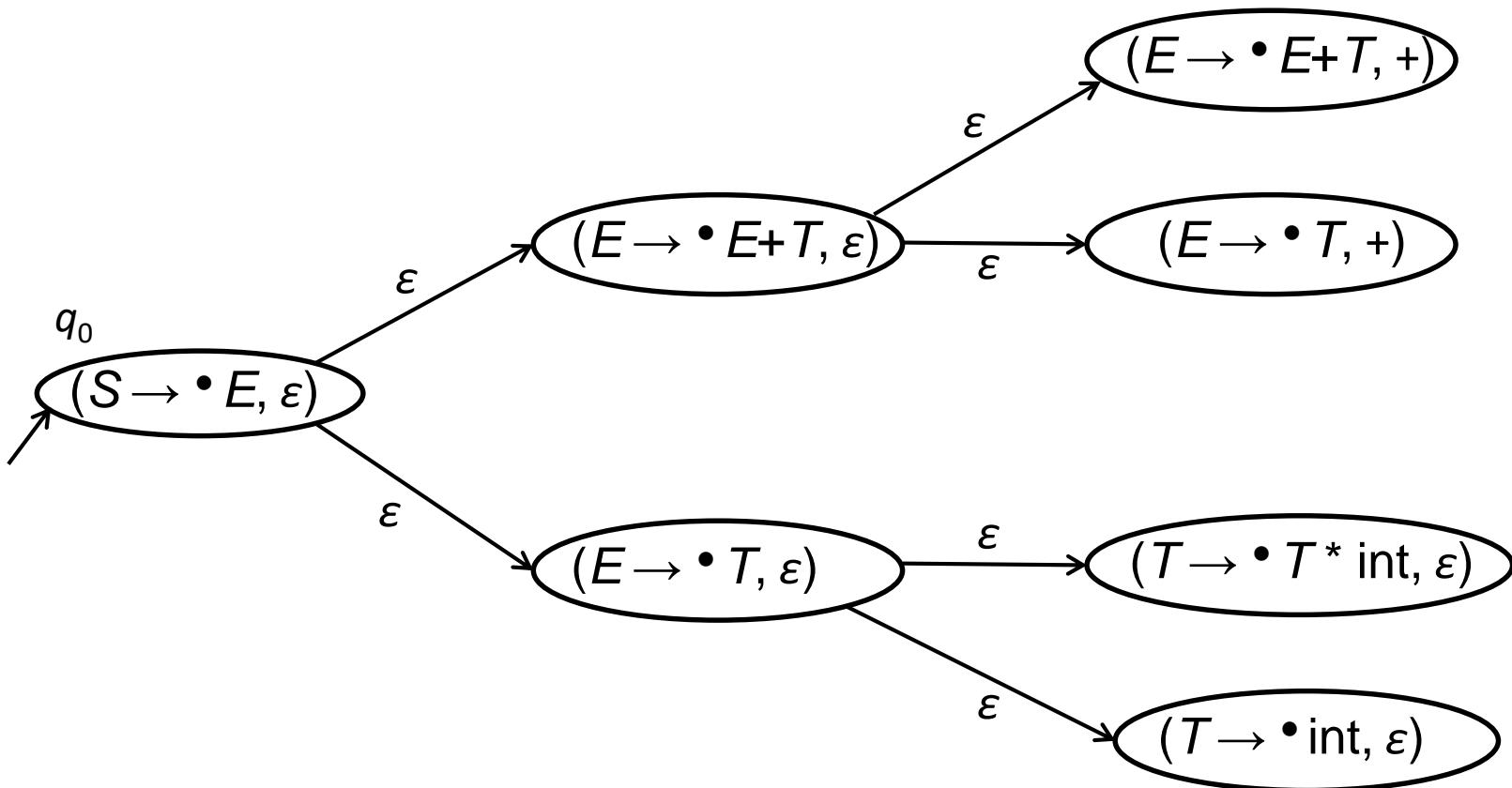
- $(A \rightarrow a^\bullet X\beta, u) \xrightarrow{X} (A \rightarrow aX^\bullet\beta, u) \quad X \in T \cup N$
 - $(A \rightarrow a^\bullet B\beta, u) \xrightarrow{\varepsilon} (B \rightarrow \bullet\gamma, v) \quad B \in N, B \rightarrow \gamma \in P, v \in \text{First}(\beta u)$
- Compiler Parsing Bottom-Up

Example: LR(1) NFA Construction



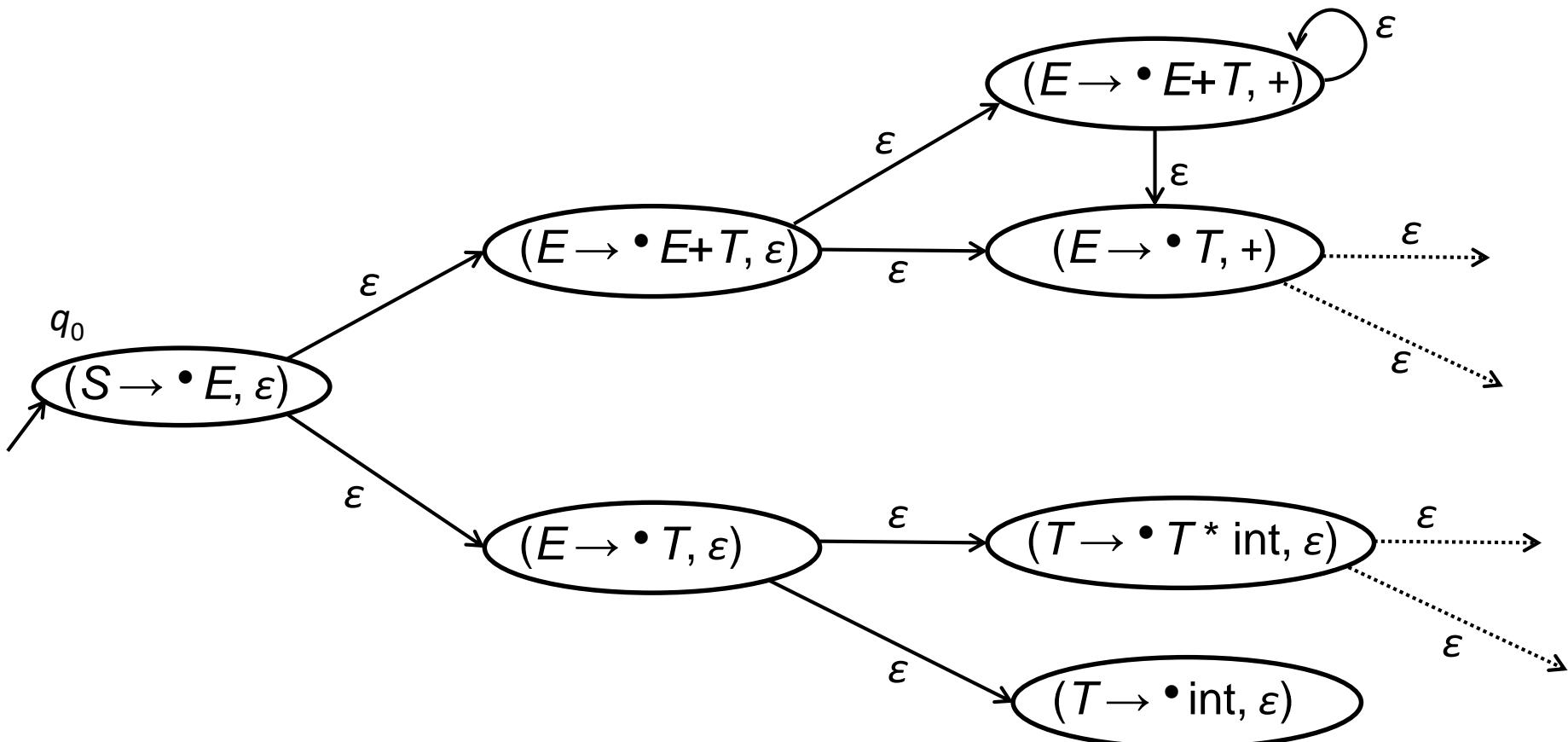
- $(A \rightarrow a^{\bullet} X\beta, u) \xrightarrow{X} (A \rightarrow aX^{\bullet}\beta, u) \quad X \in T \cup N$
 - $(A \rightarrow a^{\bullet} B\beta, u) \xrightarrow{\varepsilon} (B \rightarrow \bullet\gamma, v) \quad B \in N, B \rightarrow \gamma \in P, v \in \text{First}(\beta u)$
- Compiler Parsing Bottom-Up

Example: LR(1) NFA Construction



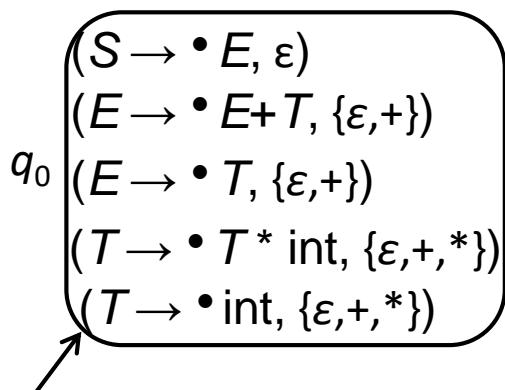
- $(A \rightarrow a^\bullet X\beta, u) \xrightarrow{X} (A \rightarrow aX^\bullet\beta, u) \quad X \in T \cup N$
 - $(A \rightarrow a^\bullet B\beta, u) \xrightarrow{\varepsilon} (B \rightarrow \bullet\gamma, v) \quad B \in N, B \rightarrow \gamma \in P, v \in \text{First}(\beta u)$
- Compiler Parsing Bottom-Up

Example: LR(1) NFA Construction

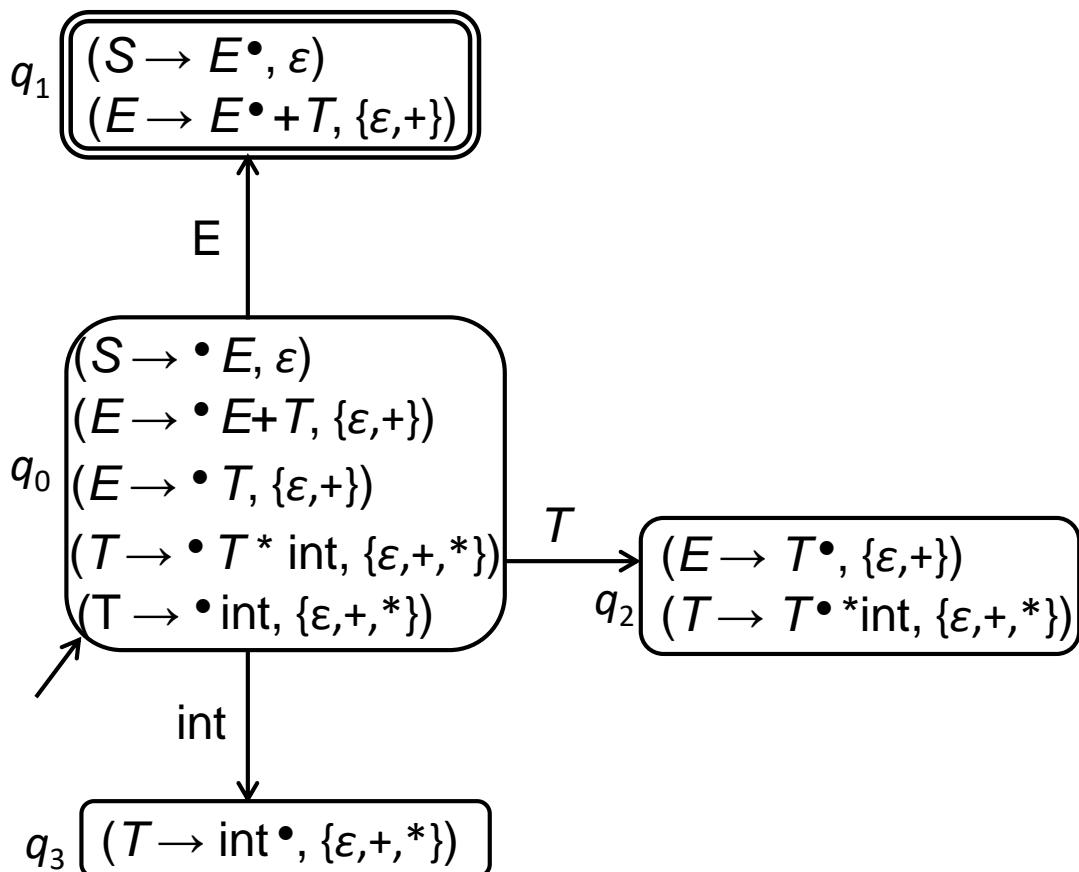


- $(A \rightarrow a^\bullet X\beta, u) \xrightarrow{X} (A \rightarrow aX^\bullet\beta, u) \quad X \in T \cup N$
 - $(A \rightarrow a^\bullet B\beta, u) \xrightarrow{\varepsilon} (B \rightarrow \bullet\gamma, v) \quad B \in N, B \rightarrow \gamma \in P, v \in \text{First}(\beta u)$
- Compiler
Parsing Bottom-Up

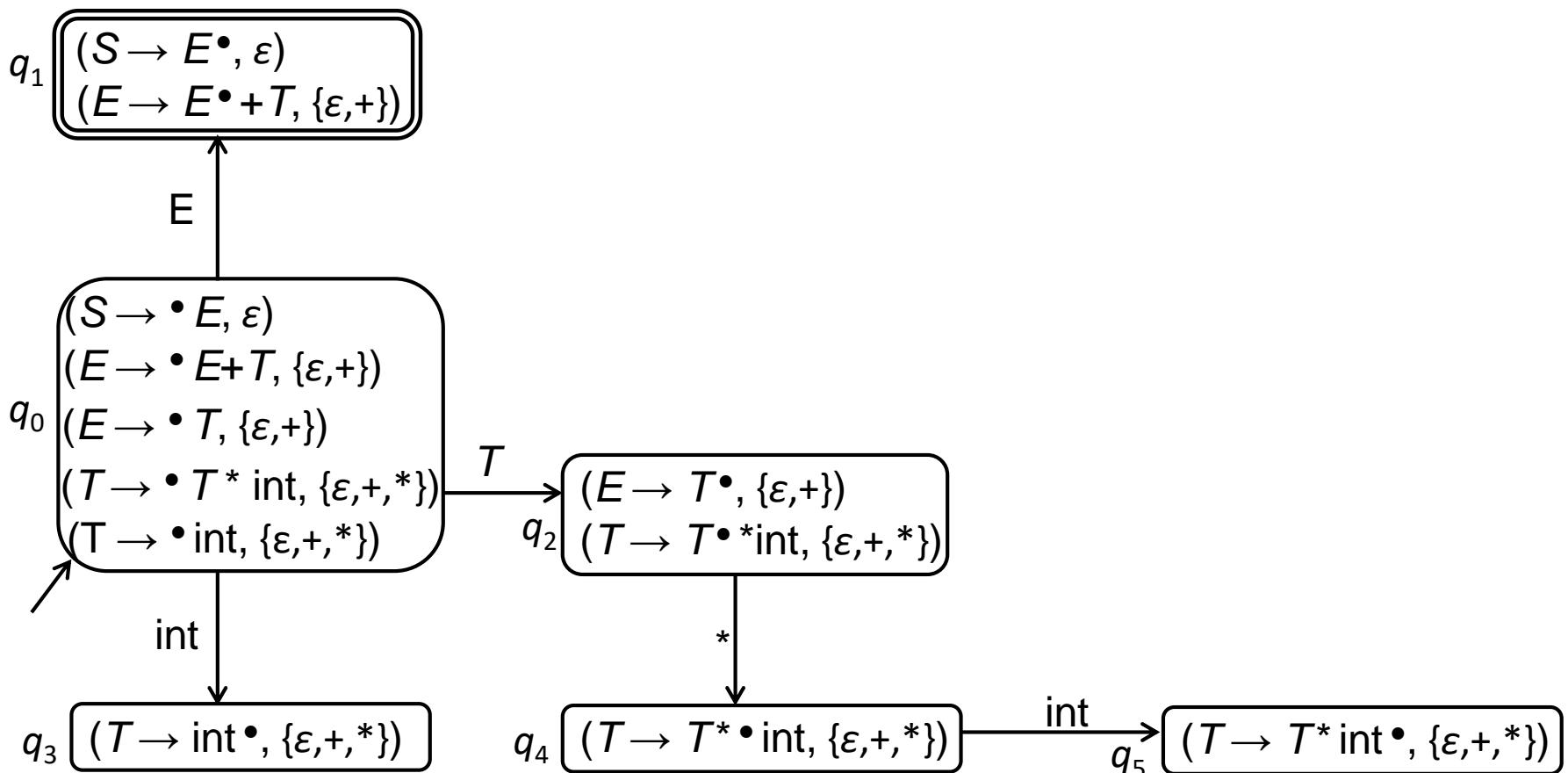
Example: LR(1) DFA Construction on the fly



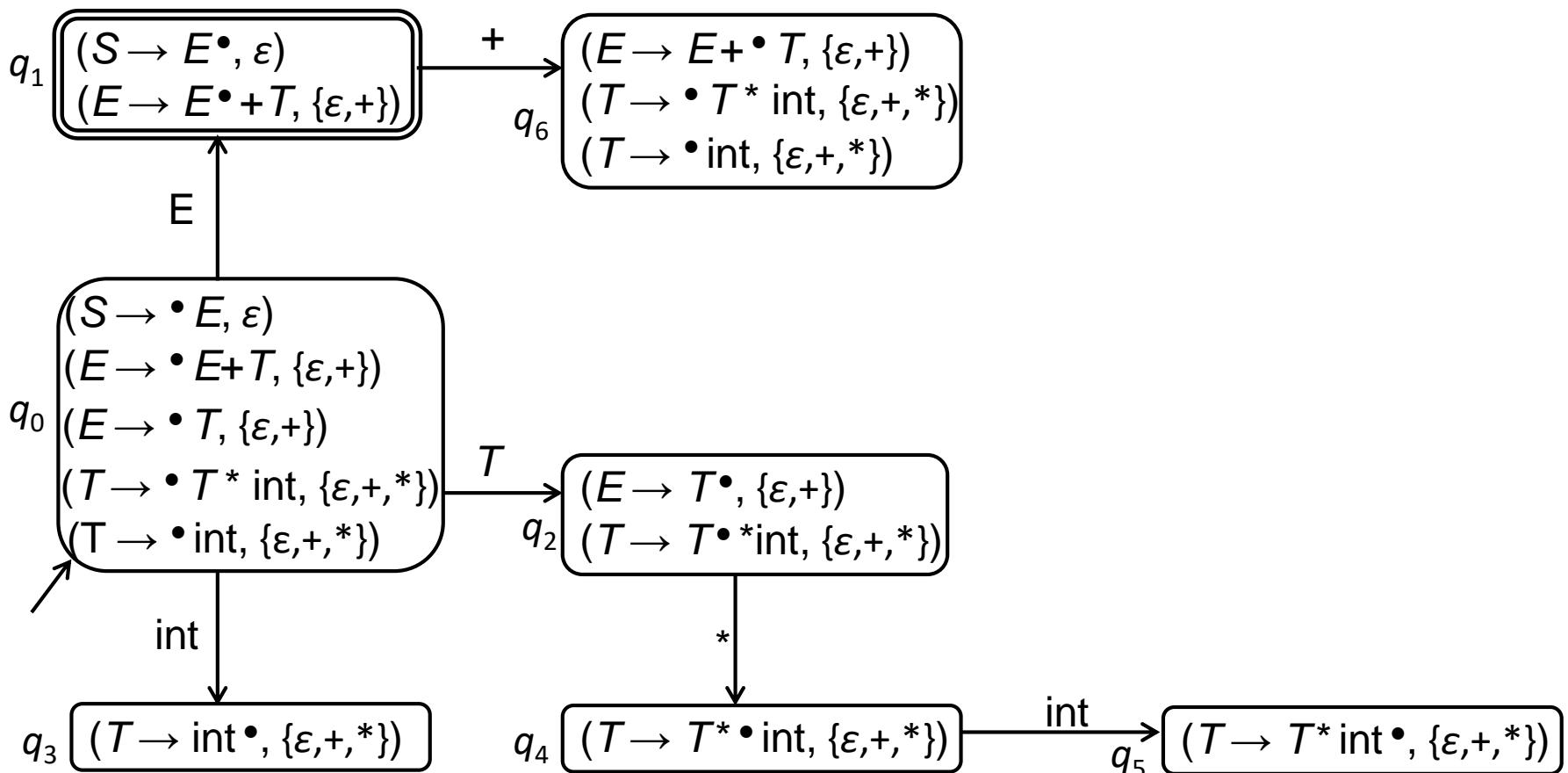
Example: LR(1) DFA Construction



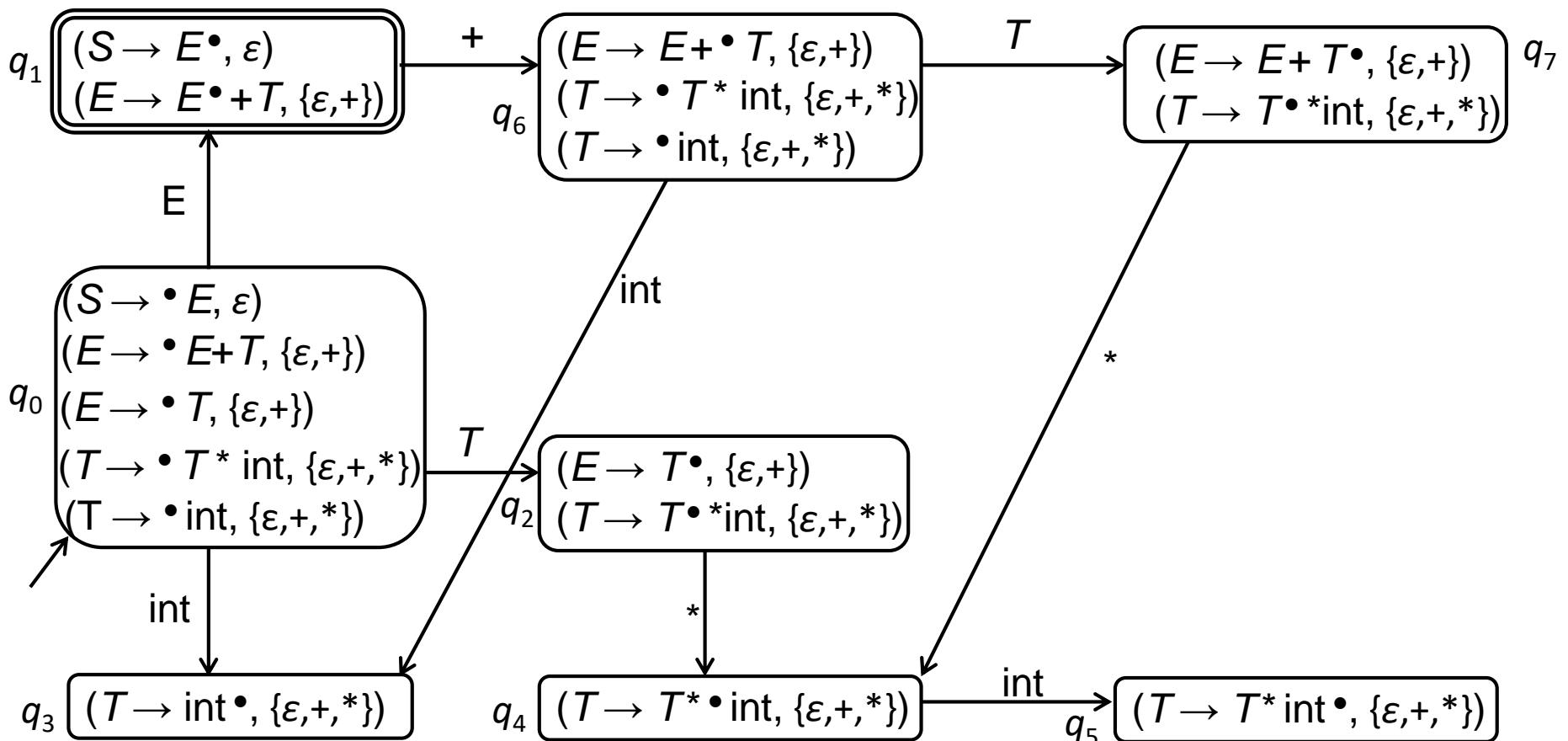
Example: LR(1) DFA Construction



Example: LR(1) DFA Construction



Example: LR(1) DFA Construction



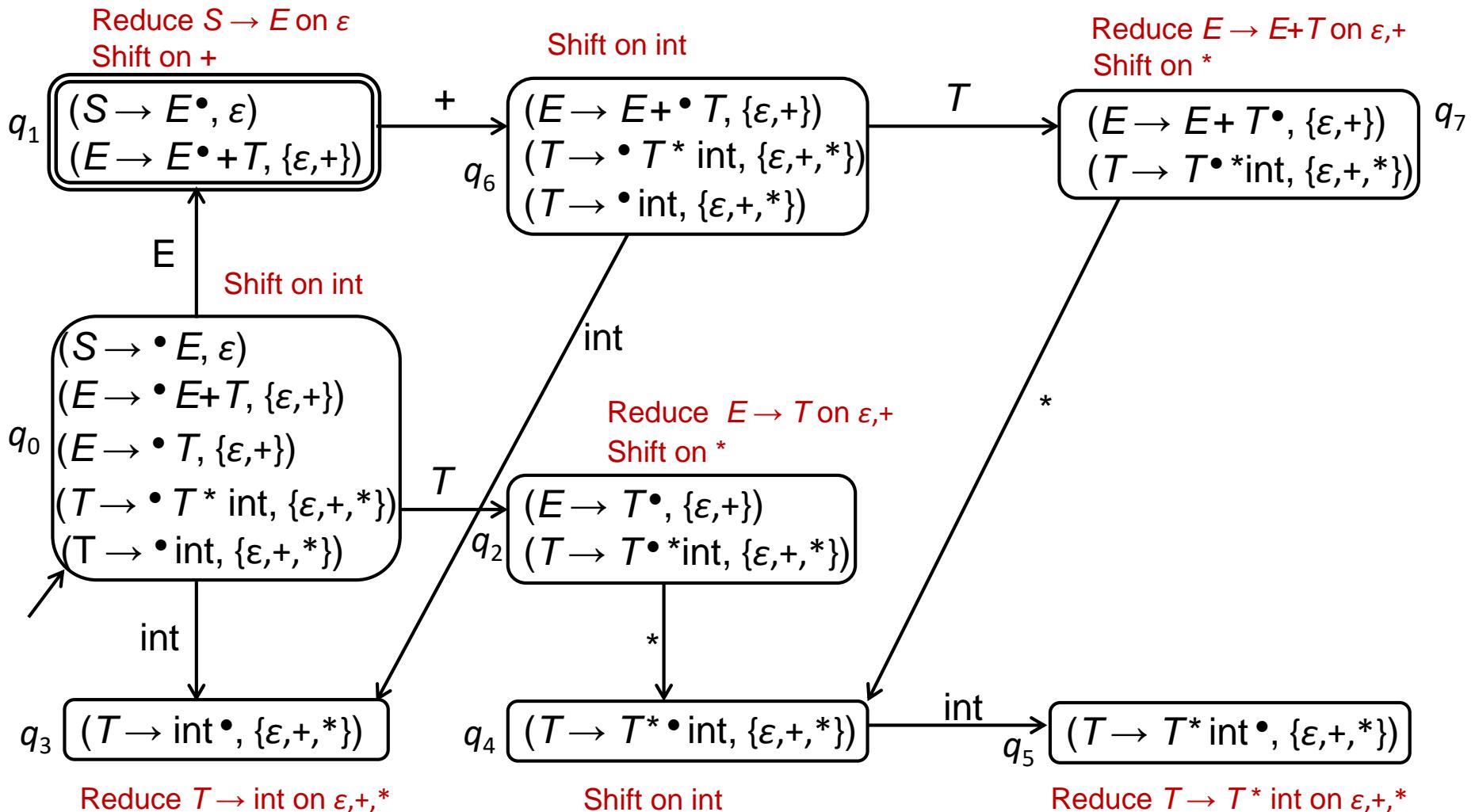
LR(1) Parsing

■ NFA

- Start state ($S' \rightarrow \bullet S, \varepsilon$)
- Accepting state ($S' \rightarrow S^\bullet, \varepsilon$)

■ DFA

- NFA → DFA conversion (states = sets of $LR(1)$ -items)
- Actions associated with states and lookahead
 - $\text{Act}(q, t) = \text{shift}$ iff $(A \rightarrow \alpha \bullet t \beta, u) \in q, t \in T$
 - $\text{Act}(q, u) = (A \rightarrow \beta)$ iff $(A \rightarrow \beta \bullet, u) \in q, u \in T \cup \{\varepsilon\}$



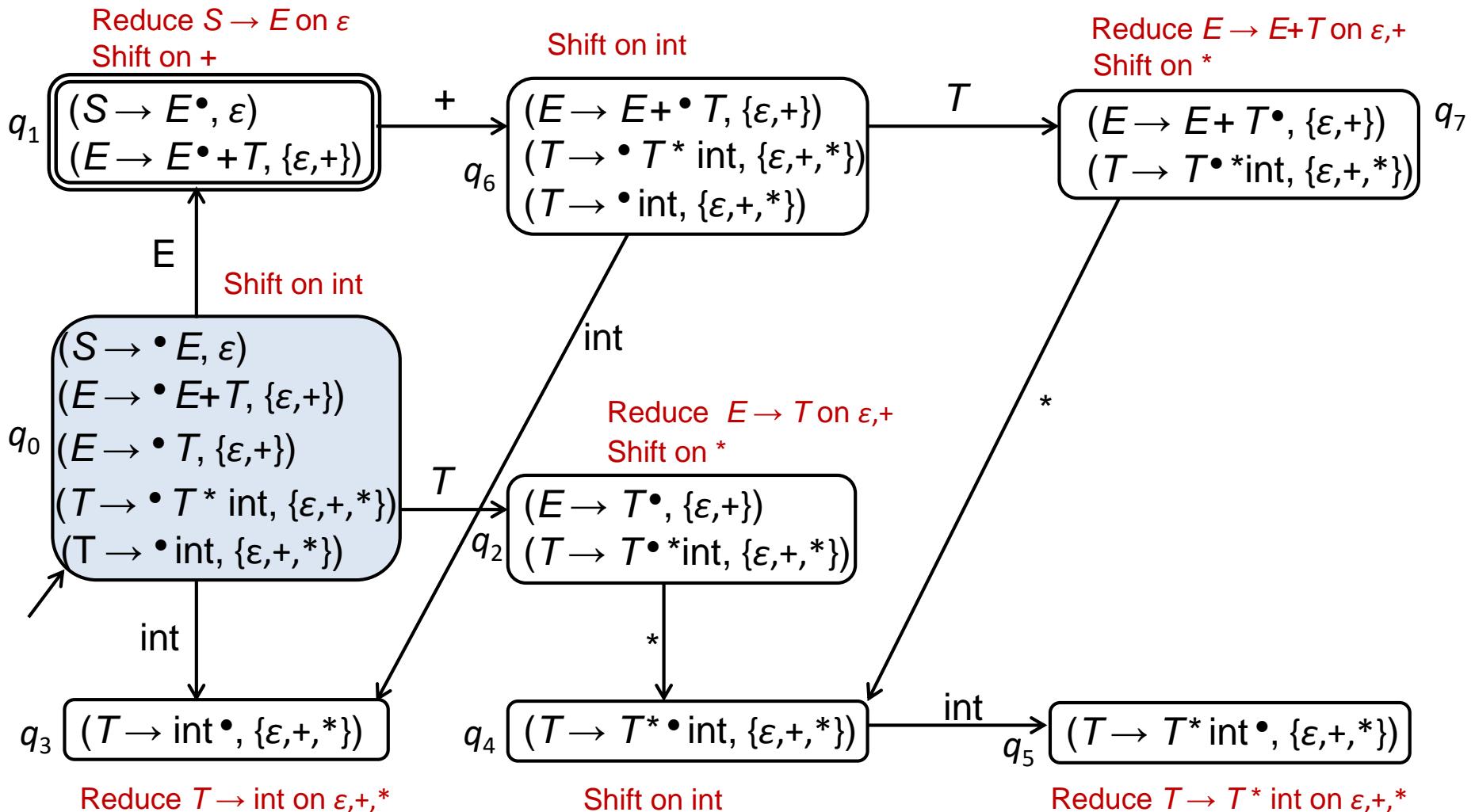
LR(1) Parsing

■ Parsing algorithm

- Stack symbols: States (Sets of LR(1)-items)
- Initial stack is $q_0 \mid \omega_c$ where ω_c is the complete input
- Parse steps
 - $\rho q \mid t \omega \rightarrow \rho q \delta(q, t) \mid \omega$ iff $\text{Act}(q, t) = \text{shift}$
 - $\rho q \rho' q' \mid \omega \rightarrow \rho q \delta(q, A) \mid \omega$ iff $\text{Act}(q', u) = (A \rightarrow \beta)$
 $|\beta| = |\rho' q'|$
 $u = \text{First}(\omega)$
- Accept with stack $q_f \mid \varepsilon$ where q_f is accepting state

Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$

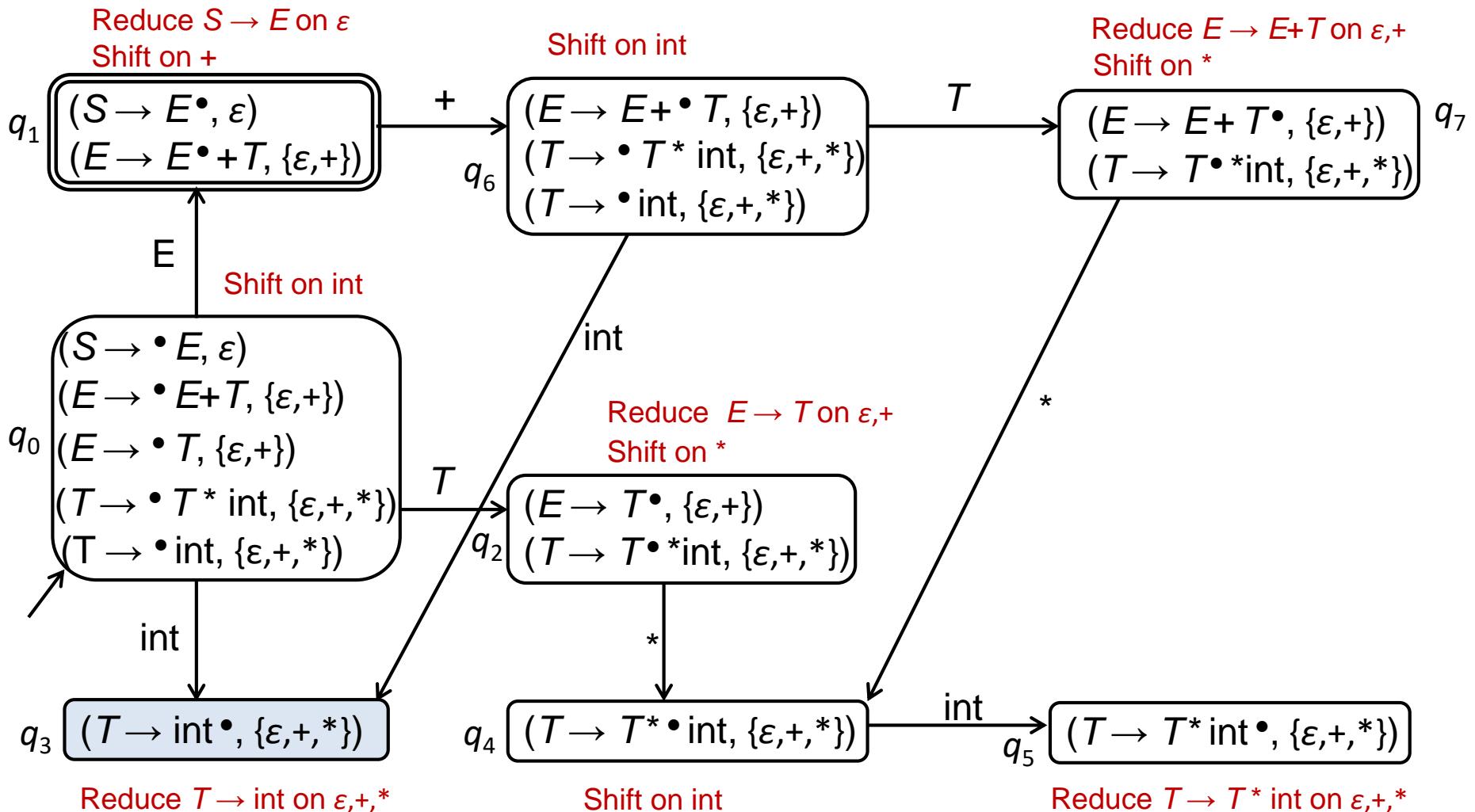


Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$

Shift

$q_0 q_3 \mid * \text{int} + \text{int}$



Example LR(1)-Parse

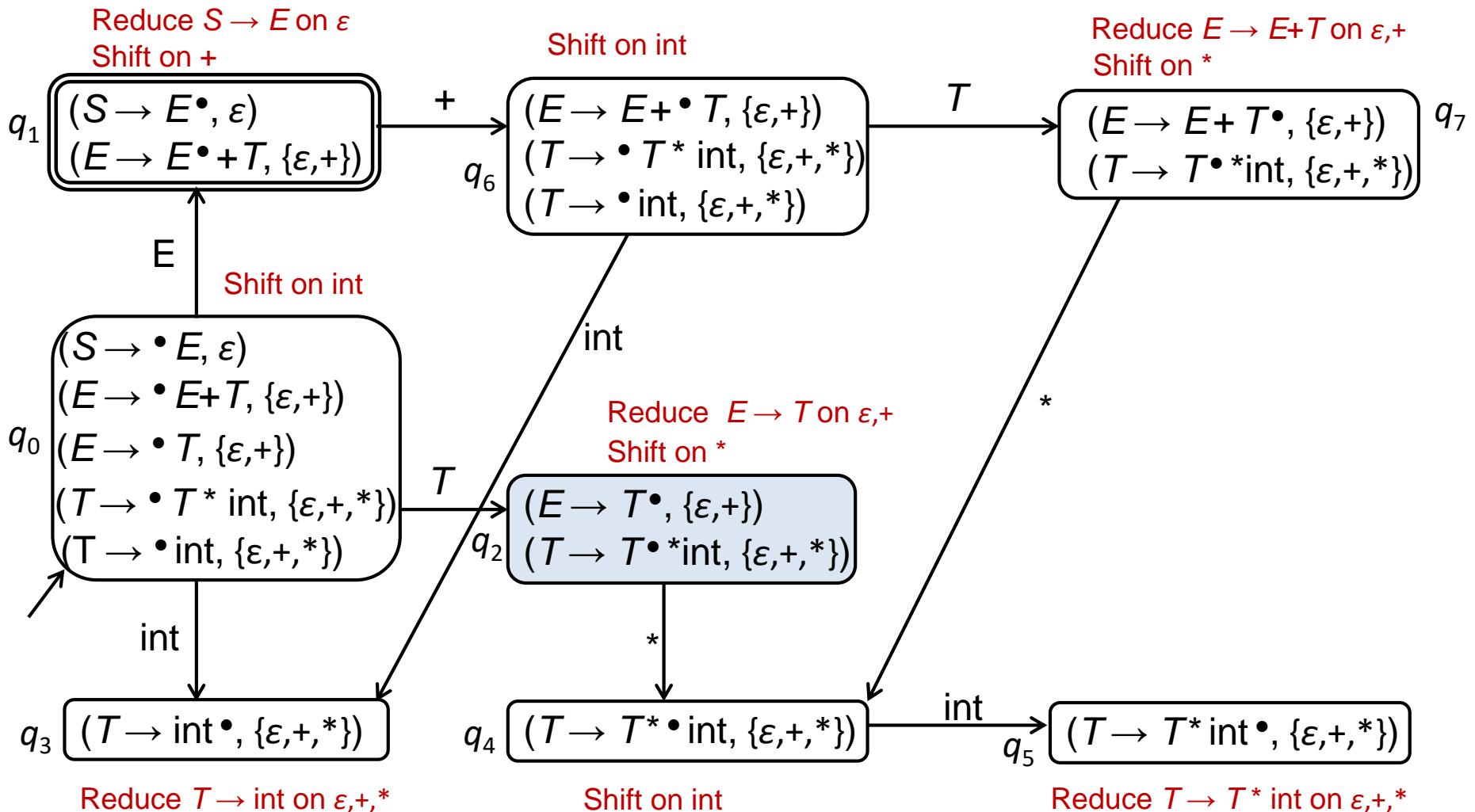
$q_0 \mid \text{int} * \text{int} + \text{int}$

Shift

$q_0 q_3 \mid * \text{int} + \text{int}$

Reduce $T \rightarrow \text{int}$

$q_0 q_2 \mid * \text{int} + \text{int}$



Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$

Shift

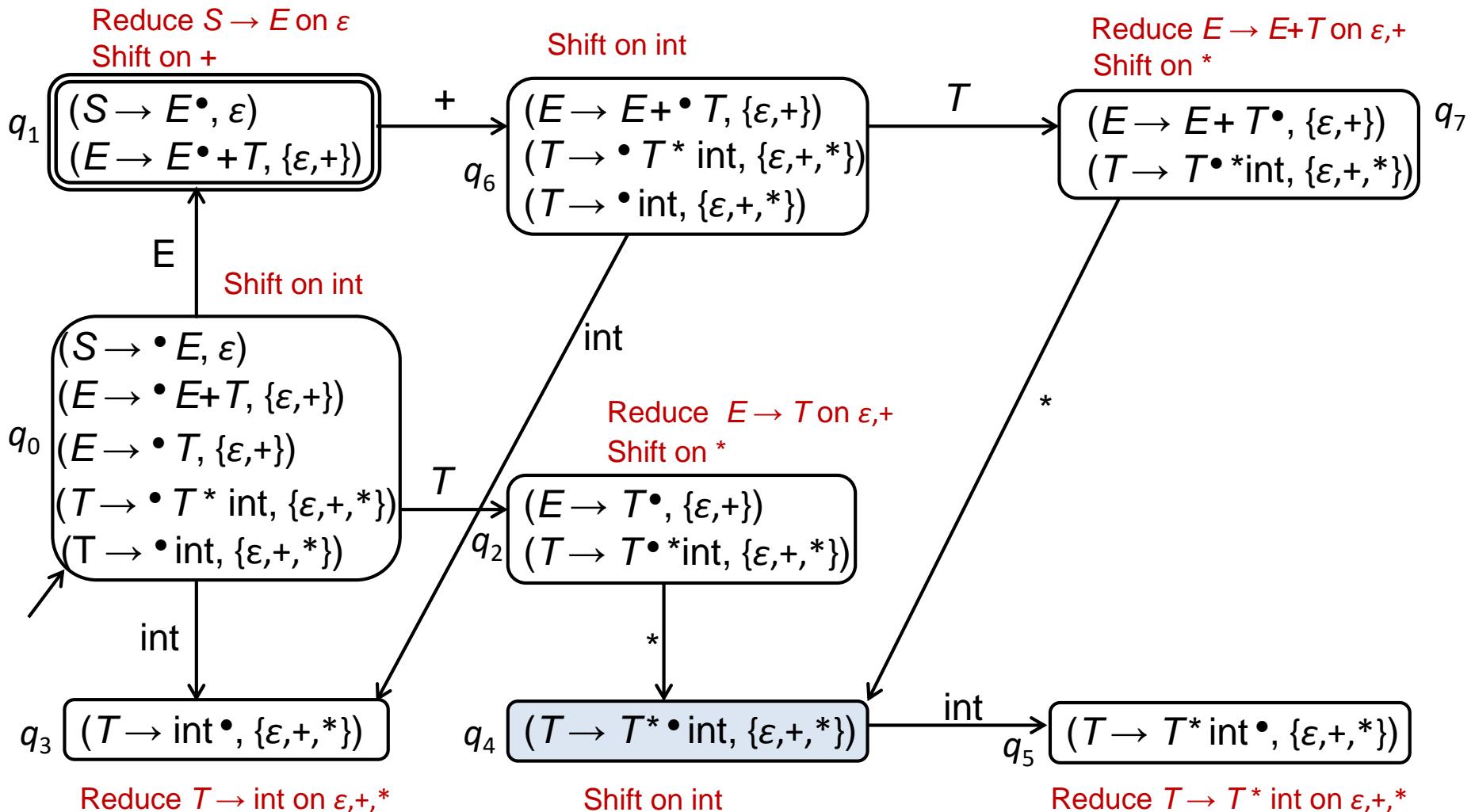
$q_0 q_3 \mid * \text{int} + \text{int}$

Reduce $T \rightarrow \text{int}$

$q_0 q_2 \mid * \text{int} + \text{int}$

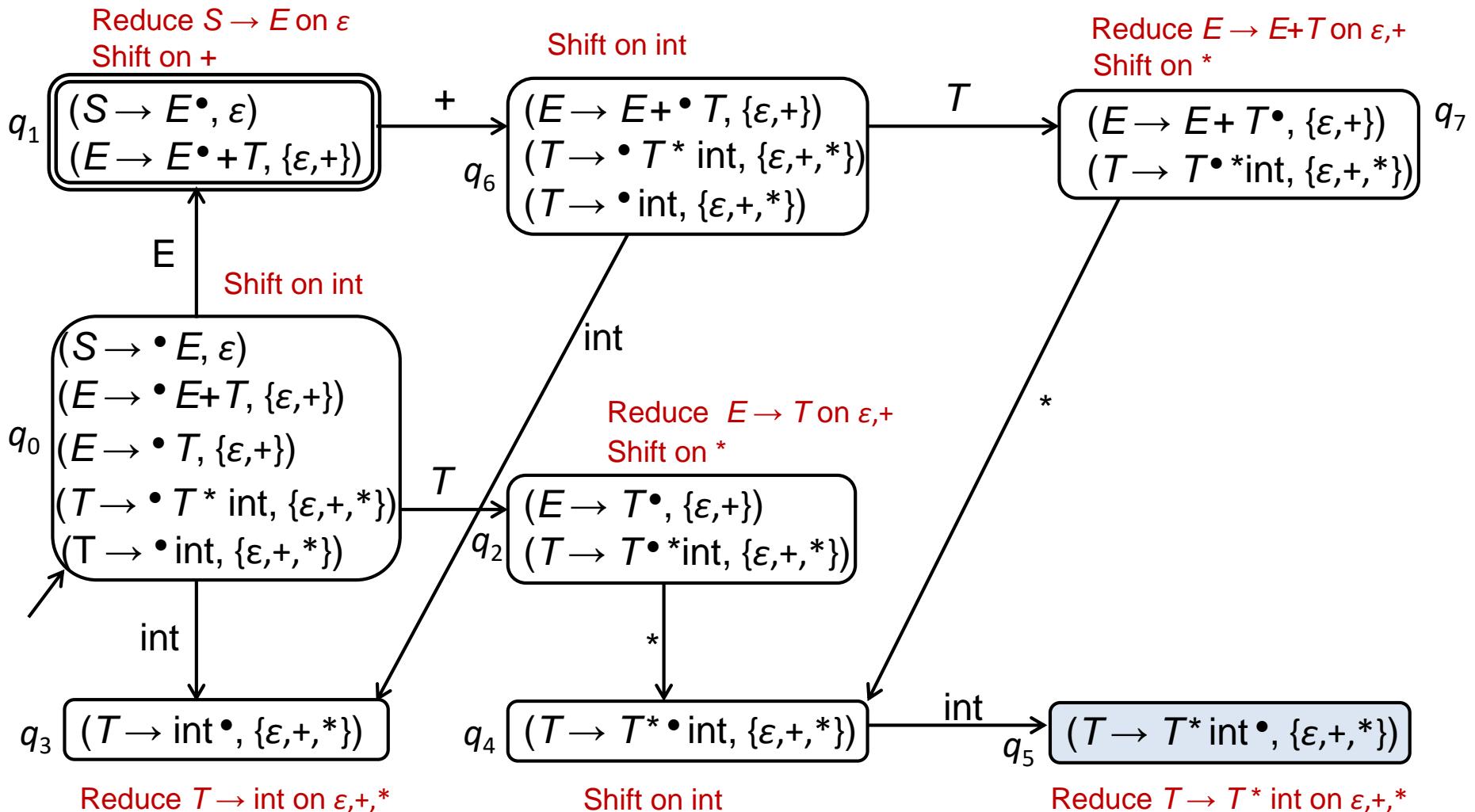
Shift

$q_0 q_2 q_4 \mid \text{int} + \text{int}$



Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	



Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$

Shift

$q_0 q_3 \mid * \text{int} + \text{int}$

Reduce $T \rightarrow \text{int}$

$q_0 q_2 \mid * \text{int} + \text{int}$

Shift

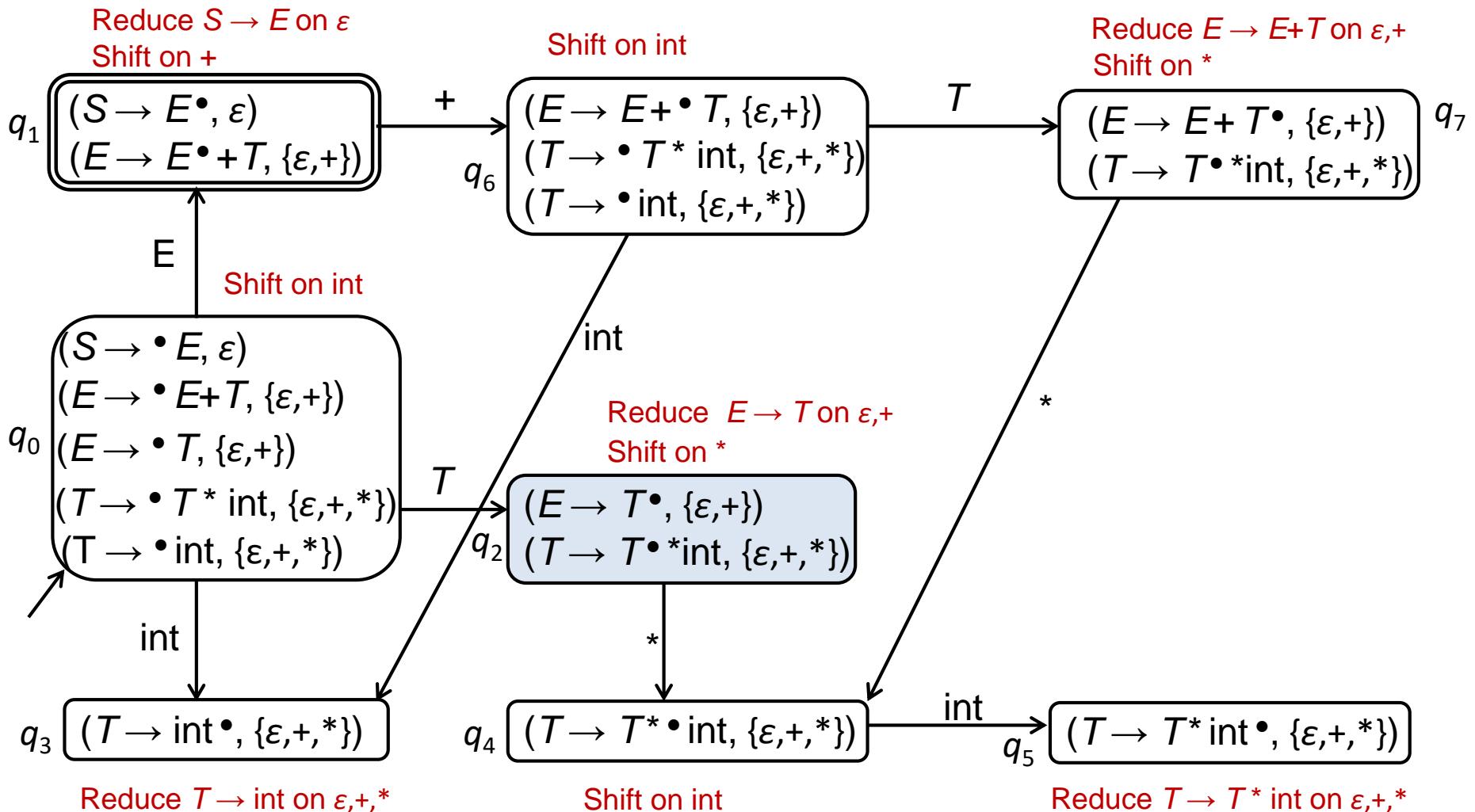
$q_0 q_2 q_4 \mid \text{int} + \text{int}$

Shift

$q_0 q_2 q_4 q_5 \mid + \text{int}$

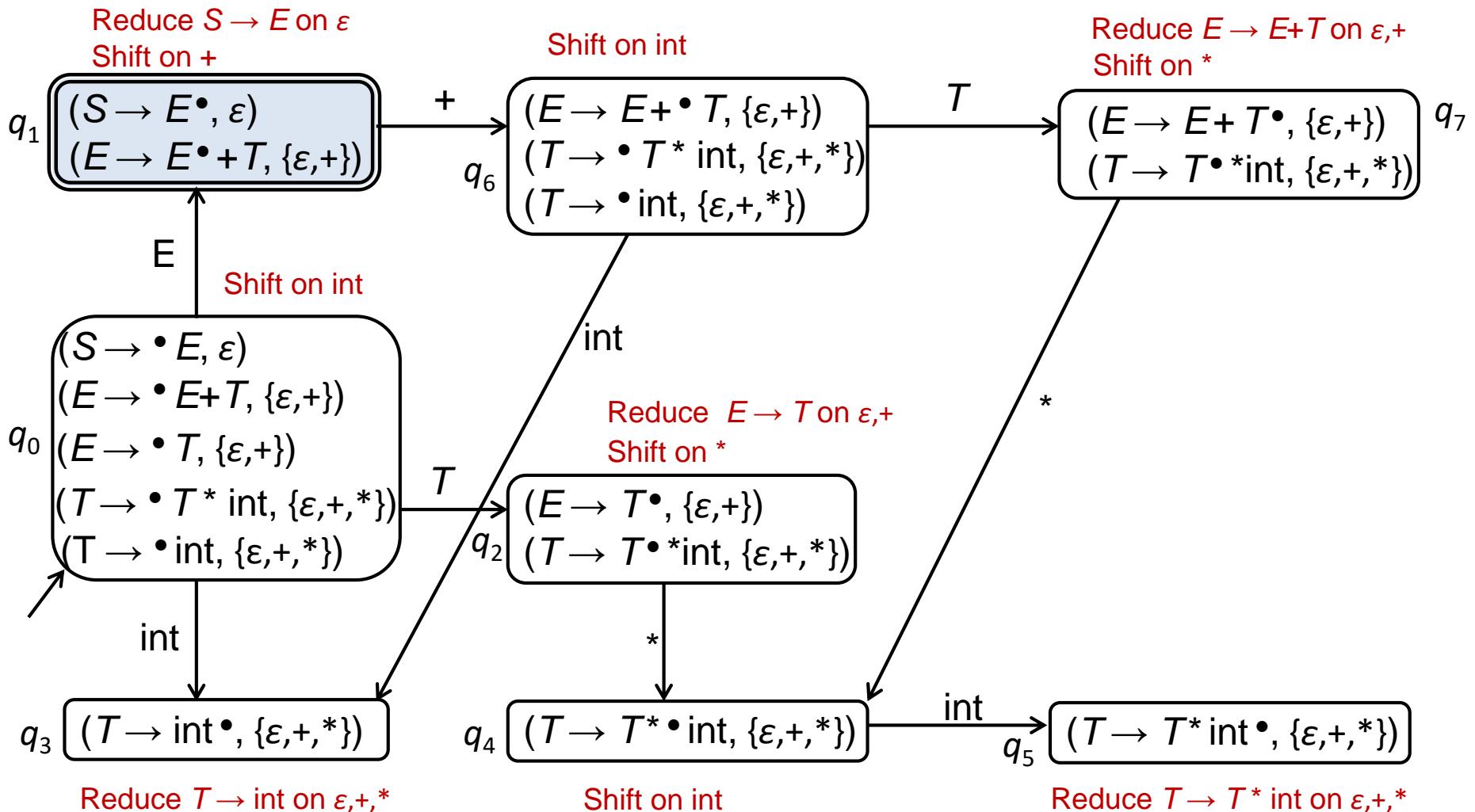
Reduce $T \rightarrow T^* \text{int}$

$q_0 q_2 \mid + \text{int}$



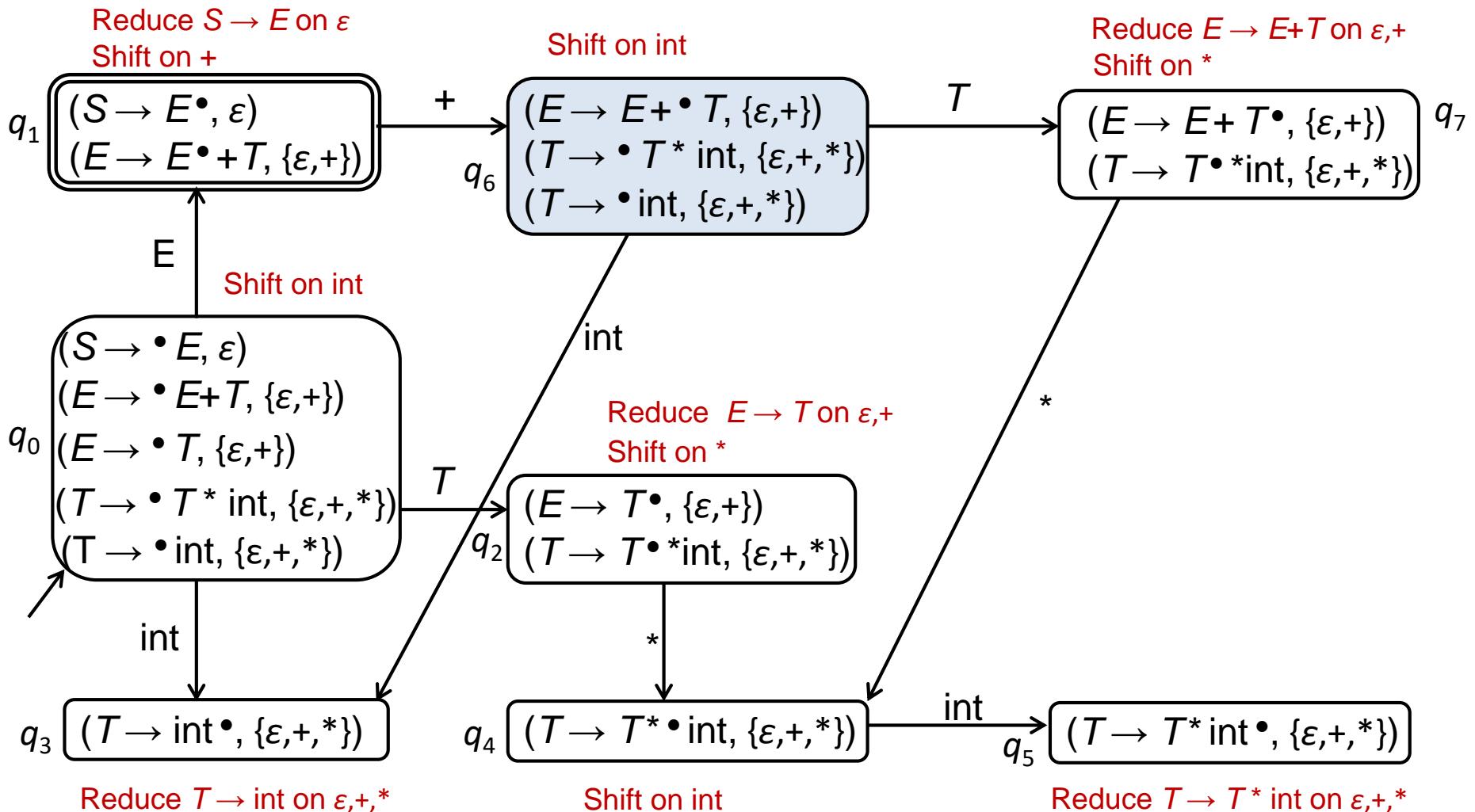
Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	



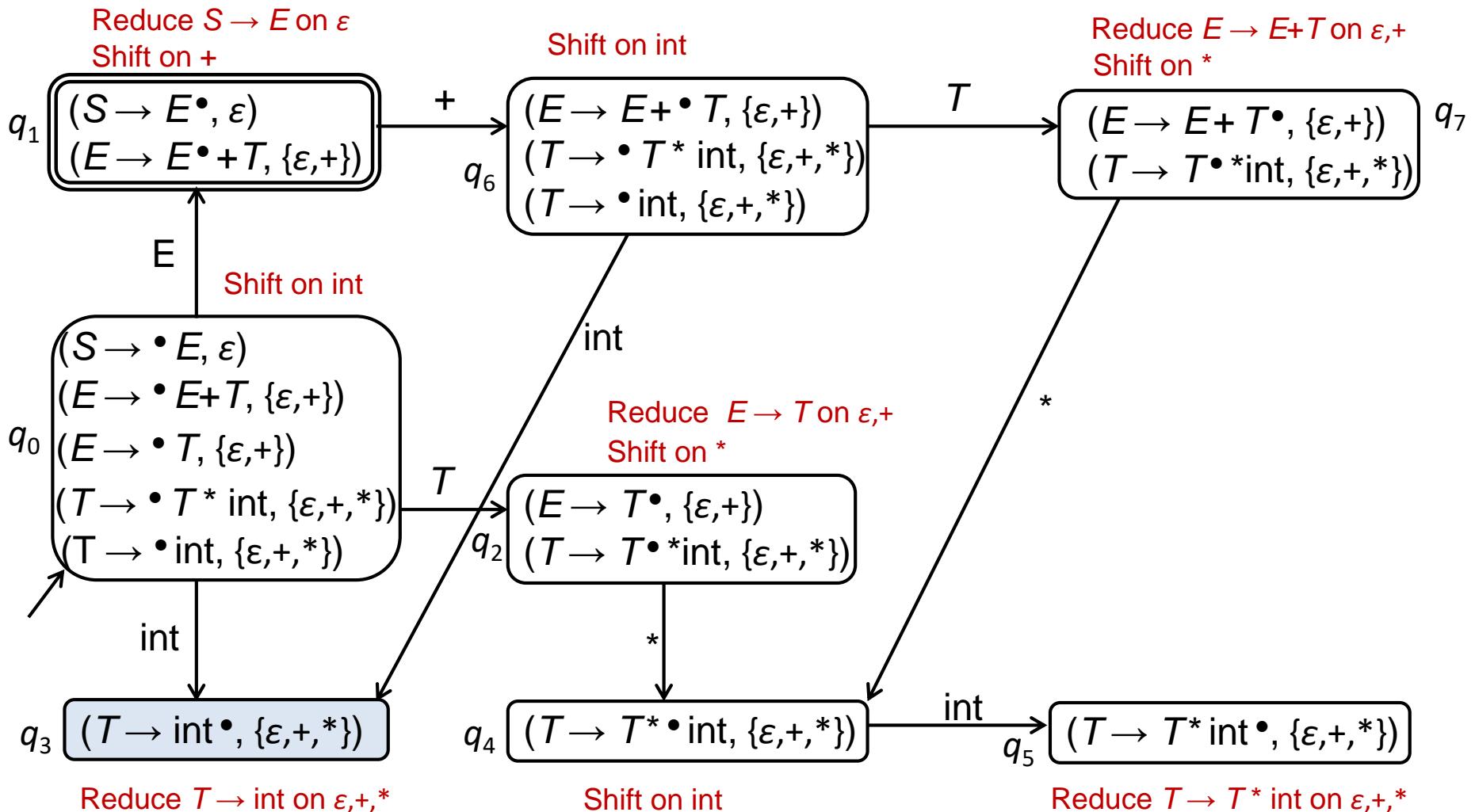
Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	Shift
$q_0 q_1 q_6 \mid \text{int}$	



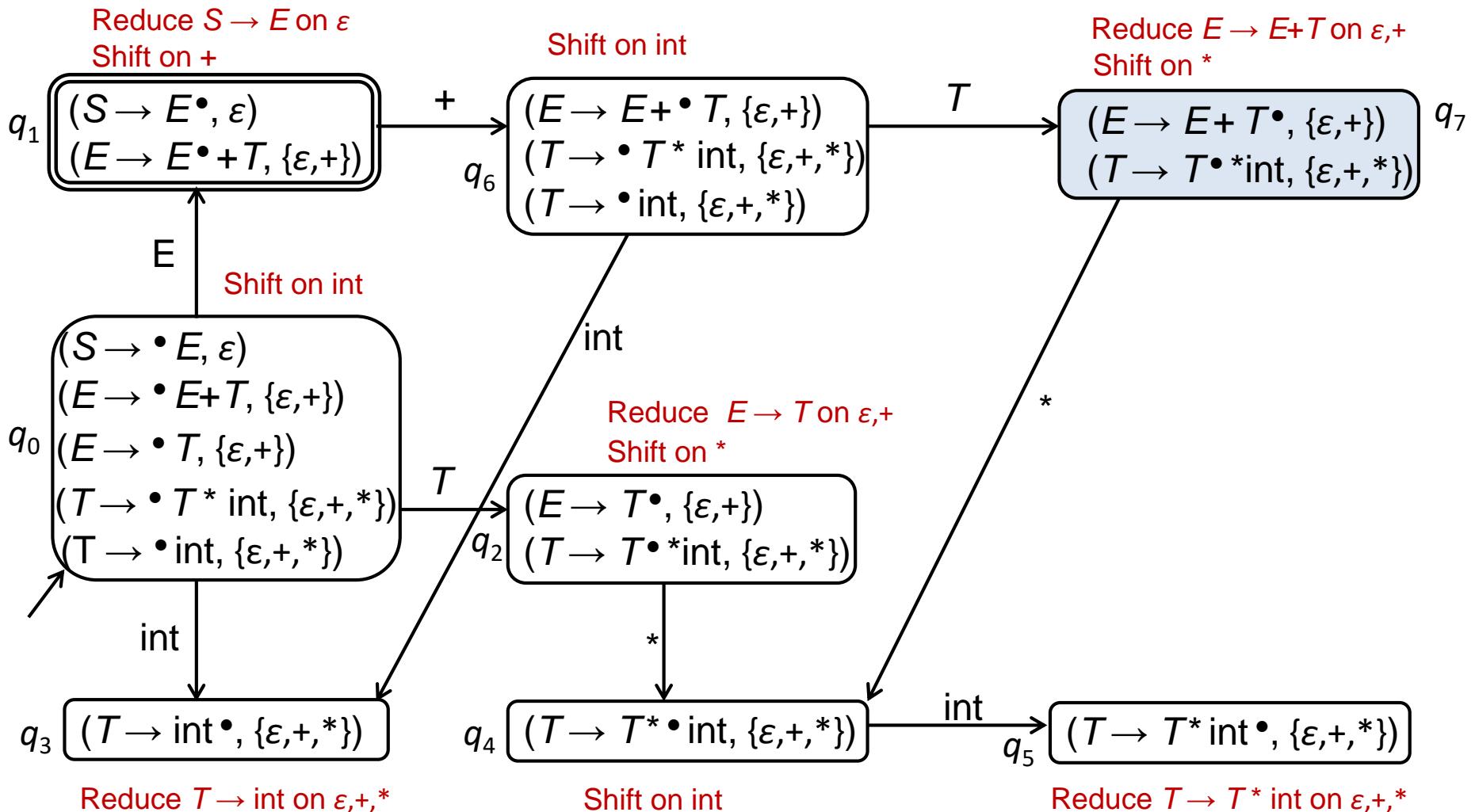
Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	Shift
$q_0 q_1 q_6 \mid \text{int}$	Shift
$q_0 q_1 q_6 q_3 \mid$	



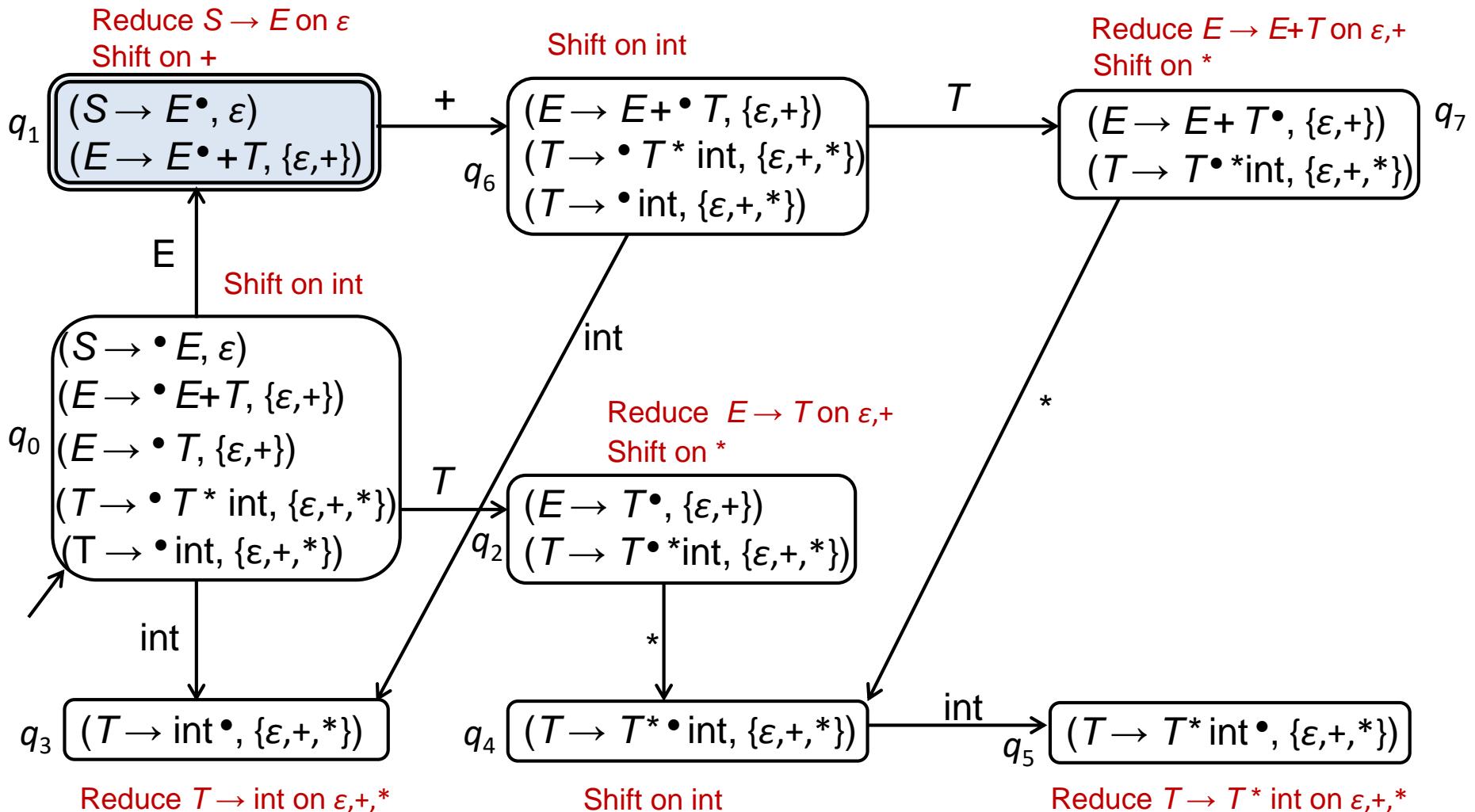
Example LR(1)-Parse

$q_0 \mid \text{int} * \text{int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	Shift
$q_0 q_1 q_6 \mid \text{int}$	Shift
$q_0 q_1 q_6 q_3 \mid$	Reduce $T \rightarrow \text{int}$
$q_0 q_1 q_6 q_7 \mid$	



Example LR(1)-Parse

$q_0 \mid \text{int}^* \text{ int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{ int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{ int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{ int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	Shift
$q_0 q_1 q_6 \mid \text{int}$	Shift
$q_0 q_1 q_6 q_3 \mid$	Reduce $T \rightarrow \text{int}$
$q_0 q_1 q_6 q_7 \mid$	Reduce $E \rightarrow E + T$
$q_0 q_1 \mid$	



Example LR(1)-Parse

$q_0 \mid \text{int}^* \text{ int} + \text{int}$	Shift
$q_0 q_3 \mid * \text{ int} + \text{int}$	Reduce $T \rightarrow \text{int}$
$q_0 q_2 \mid * \text{ int} + \text{int}$	Shift
$q_0 q_2 q_4 \mid \text{int} + \text{int}$	Shift
$q_0 q_2 q_4 q_5 \mid + \text{int}$	Reduce $T \rightarrow T^* \text{ int}$
$q_0 q_2 \mid + \text{int}$	Reduce $E \rightarrow T$
$q_0 q_1 \mid + \text{int}$	Shift
$q_0 q_1 q_6 \mid \text{int}$	Shift
$q_0 q_1 q_6 q_3 \mid$	Reduce $T \rightarrow \text{int}$
$q_0 q_1 q_6 q_7 \mid$	Reduce $E \rightarrow E + T$
$q_0 q_1 \mid$	Accept

LR Parsing

- Can be extended by increasing lookahead
 - $k > 1$ is not relevant in practice
 - $\text{LR}(k) = \text{LR}(1) = \text{Det. CFL}$ (for languages!)
- Can be simplified by disregarding the Follow-Sets
 - LR-items have the form $(A \rightarrow \alpha \bullet \beta)$
 - Actions must not depend on lookahead
 - $\text{LR}(0) = \text{Prefix-Free Det. CFL}$ (for languages!)

LR Parsing (2)

- Can be simplified by reducing state space
 - LALR(1): Collapse LR(1) states that only differ in lookahead components (*unify those Follow-sets*)
 - SLR(1): Compute LR(0) states and add lookahead information by using Follow-construction (*ignoring the context from former steps*).
 - $\text{LR}(0) \subset \text{SLR}(1) \subset \text{LALR}(1) \subset \text{LR}(1)$
- Lots of Tools: E.g. Yacc (Bison), CUP, LPG, ...