

CS335: Bottom-up Parsing

Swarnendu Biswas

Semester 2019-2020-II

CSE, IIT Kanpur

Content influenced by many excellent references, see References slide for acknowledgements.

Rightmost Derivation of $abbcede$

$S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$

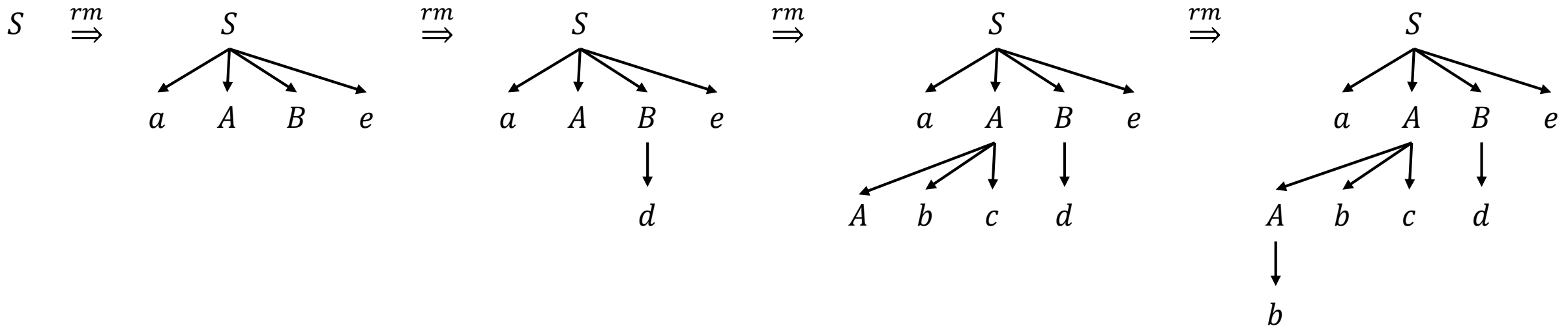
Input string: $abbcede$

$S \rightarrow aABe$

$\rightarrow aAde$

$\rightarrow aAbcde$

$\rightarrow abbcede$



Bottom-up Parsing

Constructs the parse tree starting from the leaves and working up toward the root

$S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$

Input string: <i>abbcede</i>	
$S \rightarrow aABe$	<i>abbcede</i>
$\rightarrow aAde$	$\rightarrow aAbcde$
$\rightarrow aAbcde$	$\rightarrow aAde$
$\rightarrow abbcde$	$\rightarrow aABe$
	$\rightarrow S$

reverse of rightmost derivation

Bottom-up Parsing

$S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$

Input string: $abbcd e$

$abbcd e$

$\rightarrow aAbcd e$

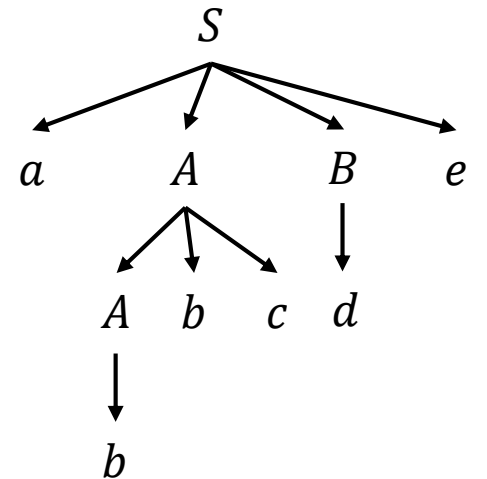
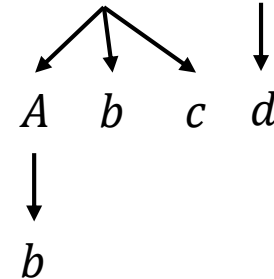
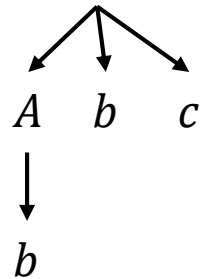
$\rightarrow aAde$

$\rightarrow aABe$

$\rightarrow S$

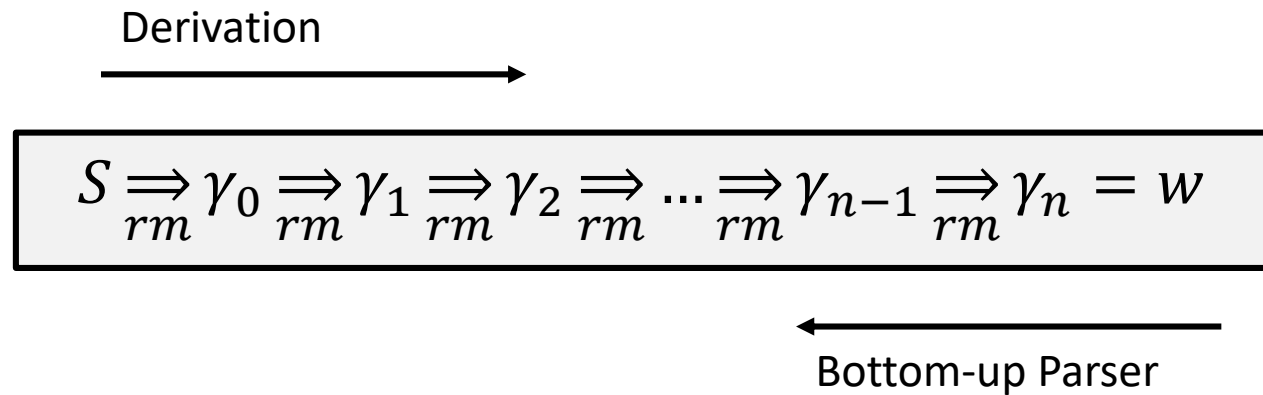
$abbcd e \Rightarrow a \quad A \quad b \quad c \quad d \quad e \Rightarrow a \quad A \quad d \quad e \Rightarrow a \quad A \quad B \quad e \Rightarrow$

\downarrow
 b



Reduction

- Bottom-up parsing **reduces** a string w to the start symbol S
 - At each reduction step, a chosen substring that is the rhs (or body) of a production is replaced by the lhs (or head) nonterminal



Handle

- Handle is a substring that matches the body of a production
 - Reducing the handle is one step in the reverse of the rightmost derivation

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

Right Sentential Form	Handle	Reducing Production
$\text{id}_1 * \text{id}_2$	id_1	$F \rightarrow \text{id}$
$F * \text{id}_2$	F	$T \rightarrow F$
$T * \text{id}_2$	id_2	$F \rightarrow \text{id}$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$

Handle

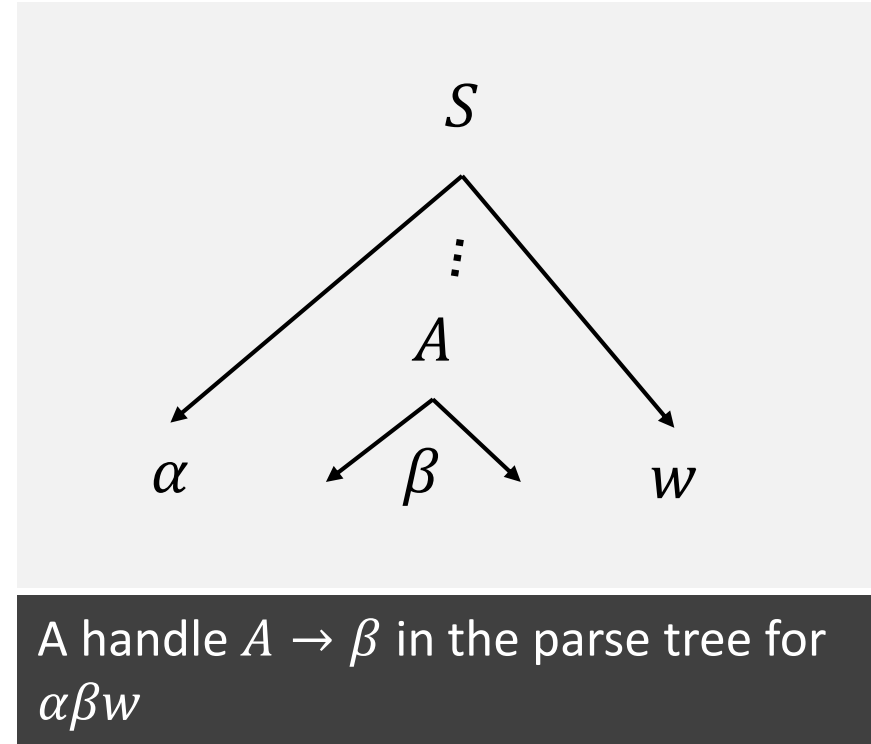
Although T is the body of the production $E \rightarrow T$, T is not a handle in the sentential form $T * \text{id}_2$

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

Right Sentential Form	Handle	Reducing Production
$\text{id}_1 * \text{id}_2$	id_1	$F \rightarrow \text{id}$
$F * \text{id}_2$	F	$T \rightarrow F$
$T * \text{id}_2$	id_2	$F \rightarrow \text{id}$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$

Handle

- If $S \xRightarrow{*}_{rm} \alpha A w \xRightarrow{rm} \alpha \beta w$, then $A \rightarrow \beta$ is a handle of $\alpha \beta w$
- String w right of a handle must contain only terminals



Handle

If grammar G is unambiguous, then every right sentential form has only one handle

If G is ambiguous, then there can be more than one rightmost derivation of $\alpha\beta w$

Shift-Reduce Parsing

Shift-Reduce Parsing

- Type of bottom-up parsing with two primary actions, shift and reduce
 - Other obvious actions are accept and error
- The input string (i.e., being parsed) consists of two parts
 - Left part is a string of terminals and nonterminals, and is stored in stack
 - Right part is a string of terminals read from an input buffer
 - Bottom of the stack and end of input are represented by \$

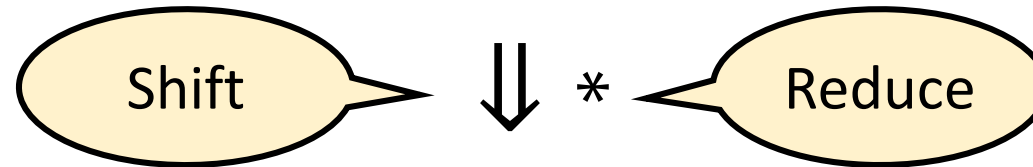
Shift-Reduce Actions

- **Shift:** shift the next input symbol from the right string onto the top of the stack
- **Reduce:** identify a string on top of the stack that is the body of a production, and replace the body with the head

Shift-Reduce Parsing

- Initial

Stack	Input
\$	w\$



- Final goal

Stack	Input
\$S	\$

Shift-Reduce Parsing

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

Stack	Input	Action
\$	$\text{id}_1 * \text{id}_2 \$$	Shift
$\$ \text{id}_1$	$* \text{id}_2 \$$	Reduce by $F \rightarrow \text{id}$
$\$ F$	$* \text{id}_2 \$$	Reduce by $T \rightarrow F$
$\$ T$	$* \text{id}_2 \$$	Shift
$\$ T *$	$\text{id}_2 \$$	Shift
$\$ T * \text{id}_2$	$\$$	Reduce by $F \rightarrow \text{id}$
$\$ T * F$	$\$$	Reduce by $T \rightarrow T * F$
$\$ T$	$\$$	Reduce by $E \rightarrow T$
$\$ E$	$\$$	Accept

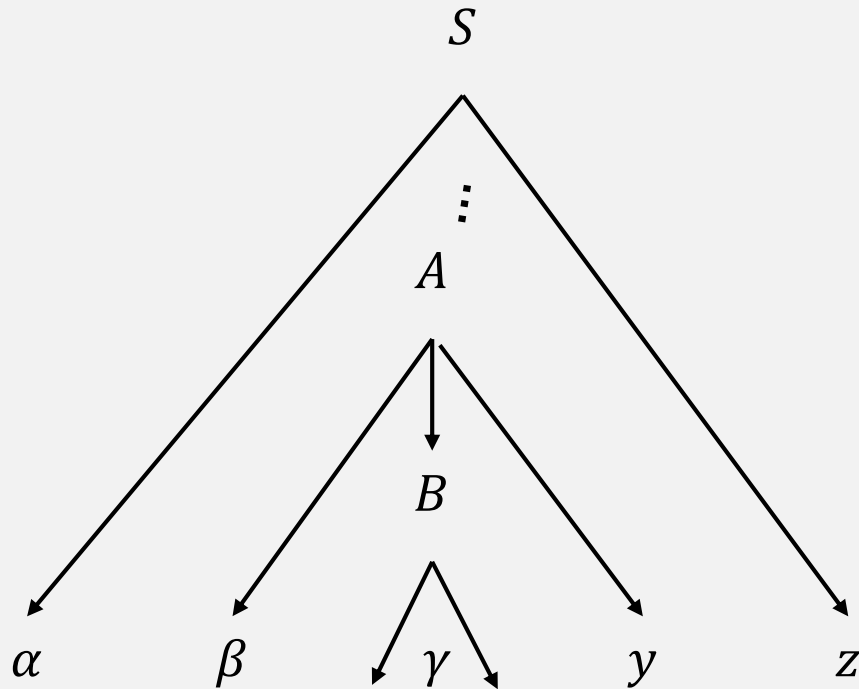
Handle on Top of the Stack

- Is the following scenario possible?

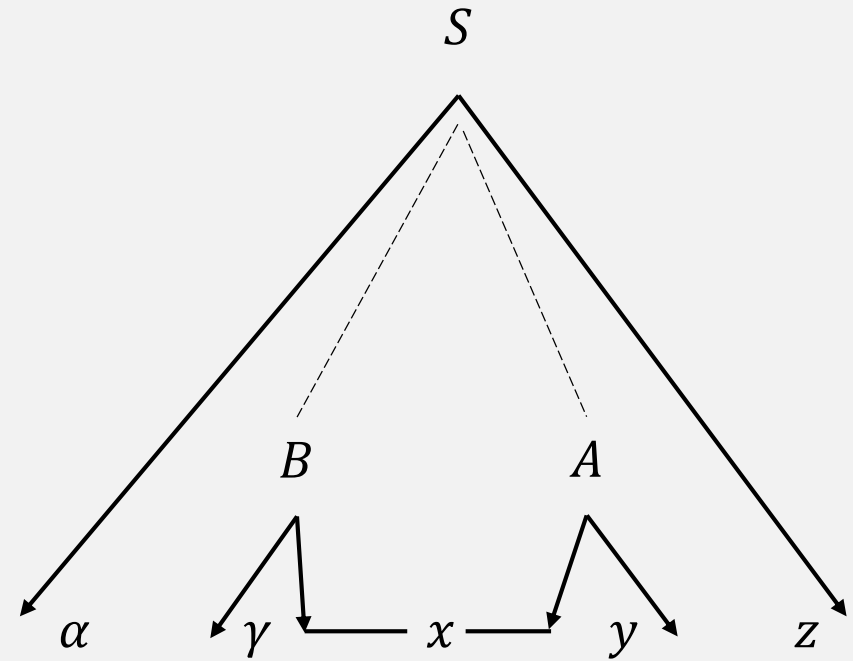
Stack	Input	Action
...		
$\$ \alpha\beta\gamma$	$w\$$	Reduce by $A \rightarrow \gamma$
$\$ \alpha\beta A$	$w\$$	Reduce by $B \rightarrow \beta$
$\$ \alpha B A$	$w\$$	
...		

Possible Choices in Rightmost Derivation

$$1. S \xRightarrow{rm} \alpha Az \xRightarrow{rm} \alpha \beta Byz \xRightarrow{rm} \alpha \beta \gamma yz$$



$$2. S \xRightarrow{rm} \alpha Bx Az \xRightarrow{rm} \alpha Bxyz \xRightarrow{rm} \alpha \gamma xyz$$



Handle on Top of the Stack

- Is the following scenario possible?

Stack	Input	Action
Handle always eventually appears on top of the stack, never inside		
...		

Shift-Reduce Actions

- Shift: shift the next input symbol from the right string onto the top of the stack
- Reduce: identify a string on top of the stack that is the body of a production, and replace the body with the head

How do you decide when to shift and when to reduce?

Steps in Shift-Reduce Parsers

- The general shift-reduce technique is
 - If there is no handle on the stack then shift
 - If there is a handle then reduce
- Bottom up parsing is essentially the process of detecting handles and reducing them
- Different bottom-up parsers differ in the way they detect handles

Challenges in Bottom-up Parsing

Which action do you pick when there is a choice?

- Both shift and reduce are valid, implies a shift-reduce conflict

Which rule to use if reduction is possible by more than one rule?

- Reduce-reduce conflict

Shift-Reduce Conflict

$$E \rightarrow E + T \mid E * E \mid \text{id}$$

id + id * id

Stack	Input	Action
\$	id + id * id\$	Shift
...		
$E + E$	$* \text{id}$ \$	Reduce by $E \rightarrow E + E$
E	$* \text{id}$ \$	Shift
$E *$	id \$	Shift
$E * \text{id}$	\$	Reduce by $E \rightarrow \text{id}$
$E * E$	\$	Reduce by $E \rightarrow E * E$
E	\$	

id + id * id

Stack	Input	Action
\$	id + id * id\$	Shift
...		
$E + E$	$* \text{id}$ \$	Shift
$E + E *$	id \$	Shift
$E + E * \text{id}$	\$	Reduce by $E \rightarrow \text{id}$
$E + E * E$	\$	Reduce by $E \rightarrow E * E$
$E + E$	\$	Reduce by $E \rightarrow E + E$
E	\$	

Shift-Reduce Conflict

Stmt \rightarrow **if** *Expr* **then** *Stmt*
 | **if** *Expr* **then** *Stmt* **else** *Stmt*
 | *other*

Stack	Input	Action
... if <i>Expr</i> then <i>Stmt</i>	else ... \$	

Shift-Reduce Conflict

Stmt → **if** *Expr* **then** *Stmt*
 | **if** *Expr* **then** *Stmt* **else** *Stmt*
 | *other*

Stack	Input	Action
... if <i>Expr</i> then <i>Stmt</i>	else ... \$	

What is a correct thing to do for this grammar – shift or reduce?

Reduce-Reduce Conflict

$M \rightarrow R + R \mid R + c \mid R$ $R \rightarrow c$
--

$c + c$		
Stack	Input	Action
\$	$c + c$ \$	Shift
$\$c$	$+c$ \$	Reduce by $R \rightarrow c$
$\$R$	$+c$ \$	Shift
$\$R +$	c \$	Shift
$\$R + c$	\$	Reduce by $R \rightarrow c$
$\$R + R$	\$	Reduce by $R \rightarrow R + R$
$\$M$	\$	

$c + c$		
Stack	Input	Action
\$	$c + c$ \$	Shift
$\$c$	$+c$ \$	Reduce by $R \rightarrow c$
$\$R$	$+c$ \$	Shift
$\$R +$	c \$	Shift
$\$R + c$	\$	Reduce by $M \rightarrow R + c$
$\$M$	\$	

LR Parsing

LR(k) Parsing

- Popular bottom-up parsing scheme
 - L is for left-to-right scan of input
 - R is for reverse of rightmost derivation
 - k is the number of lookahead symbols
- LR parsers are table-driven, like the nonrecursive LL parser
- LR grammar is one for which we can construct an LR parsing table

Popularity of LR Parsing

- LR parsers can recognize all programming language constructs for which context-free grammars can be written
- LR-parsing is most general nonbacktracking shift-reduce parsing method
- LR parsing works for a superset of grammars that can be parsed with predictive or LL parsers
 - LL(k) parsing must predict which production to use having seen only the first k tokens of the right-hand side
 - LR(k) parsing postpones the decision until it has seen input tokens corresponding to the entire right-hand side of the production

Popularity of LR Parsing

Can recognize all language constructs with CFGs

Most general nonbacktracking shift-reduce parsing method

Works for a superset of grammars parsed with predictive or LL parsers



Popularity of LR Parsing

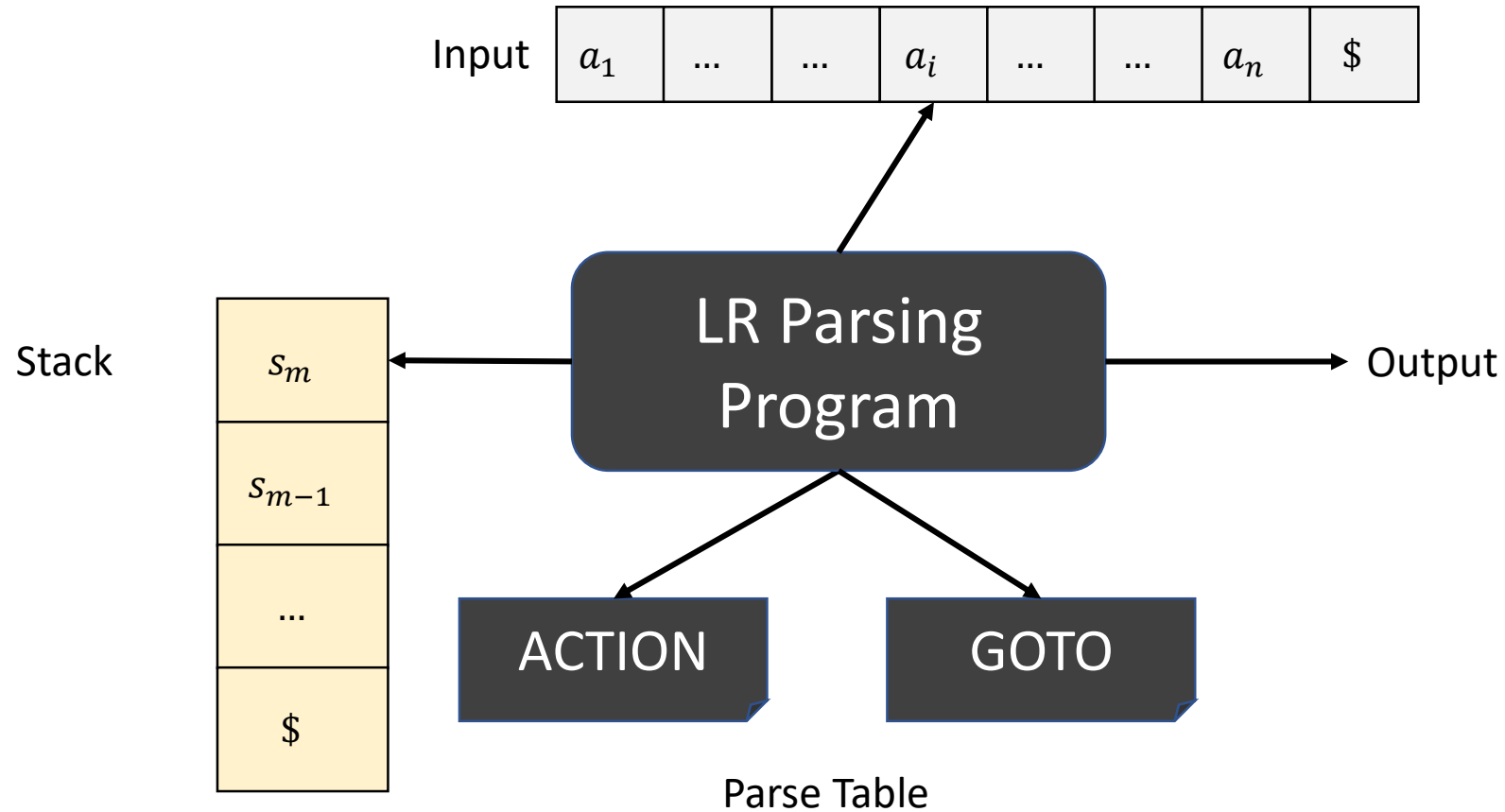
Can recognize all language constructs with CFGs

Most general nonbacktracking shift-reduce parsing method

Works for a superset of grammars parsed with predictive or LL parsers

- LL(k) parsing predicts which production to use having seen only the first k tokens of the right-hand side
- LR(k) parsing can decide after it has seen input tokens corresponding to the entire right-hand side of the production

Block Diagram of LR Parser



LR(0) Item

- An LR(0) item (also called item) of a grammar G is a production of G with a dot at some position in the body
- An item indicates how much of a production we have seen
 - Symbols on the left of “.” are already on the stack
 - Symbols on the right of “.” are expected in the input

Production	Items
$A \rightarrow XYZ$	$A \rightarrow .XYZ$
	$A \rightarrow X.YZ$
	$A \rightarrow XY.Z$
	$A \rightarrow XYZ.$

$A \rightarrow .XYZ$ indicates that we expect a string derivable from XYZ next on the input

LR(0) Automaton

- An LR parser makes shift-reduce decisions by maintaining states
- States represent sets of LR(0) items
- The canonical LR(0) collection is used for constructing a DFA for parsing

Closure Operation

- Let I be a set of items for a grammar G
- $\text{Closure}(I)$ is constructed by
 1. Add every item in I to $\text{Closure}(I)$
 2. If $A \rightarrow \alpha.B\beta$ is in $\text{Closure}(I)$ and $B \rightarrow \gamma$ is a rule, then add $B \rightarrow \gamma$ to $\text{Closure}(I)$ if not already added
 3. Repeat until no more new items can be added to $\text{Closure}(I)$

Example of Closure

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

Suppose $I = \{E' \rightarrow \cdot E\}$, compute $\text{Closure}(I)$

Example of Closure

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

Suppose $I = \{E' \rightarrow \cdot E\}$

Closure(I) = {
 $E' \rightarrow \cdot E$,
 $E \rightarrow \cdot E + T$,
 $E \rightarrow \cdot T$,
 $T \rightarrow \cdot T * F$,
 $T \rightarrow \cdot F$,
 $F \rightarrow \cdot (E)$,
 $F \rightarrow \cdot \mathbf{id}$
}

Goto Operation

- Suppose I is a set of items and X is a grammar symbol
- $\text{Goto}(I, X)$ is the closure of set all items $[A \rightarrow \alpha X. \beta]$ such that $[A \rightarrow \alpha. X \beta]$ is in I
 - If I is a set of items for some valid prefix α then $\text{Goto}(I, X)$ is set of valid items for prefix αX
- Intuitively, $\text{Goto}()$ defines the transitions in the LR(0) automaton

Example of Goto

- Suppose $I = \{E' \rightarrow E., E \rightarrow E. + T\}$
- Compute $\text{Goto}(I, +)$

Example of Goto

$$\begin{array}{l} E' \rightarrow E \\ E \rightarrow E + T | T \\ T \rightarrow T * F | F \\ F \rightarrow (E) | \mathbf{id} \end{array}$$

- Suppose $I = \{E' \rightarrow E., E \rightarrow E. + T\}$

$$\begin{aligned} \text{Goto}(I, +) = \{ & \\ & E \rightarrow E+.T, \\ & E \rightarrow E+.T, \\ & T \rightarrow .T * F, \\ & T \rightarrow .F, \\ & F \rightarrow .(E), \\ & F \rightarrow .\mathbf{id} \\ & \} \end{aligned}$$

Canonical Collection of Sets of LR(0) Items

$C = \text{Closure}(\{S' \rightarrow S\})$

repeat

 for each set of items I in C

 for each grammar symbol X

 if ($\text{Goto}(I, X)$ is not empty and not in C)

 add $\text{GOTO}(I, X)$ to C

until no new sets of items are added to C

Canonical Collection of Sets of LR(0) Items

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

- Compute the canonical collection for the expression grammar

Canonical Collection of Sets of LR(0) Items

$$I_0 = \text{Closure}(E' \rightarrow E) = \{ \\ E' \rightarrow \cdot E, \\ E \rightarrow \cdot E + T, \\ E \rightarrow \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot \text{id}, \\ \}$$

$$I_1 = \text{Goto}(I_0, E) = \{ \\ E' \rightarrow E \cdot, \\ E \rightarrow E \cdot + T \\ \}$$

$$I_2 = \text{Goto}(I_0, T) = \{ \\ E \rightarrow T \cdot, \\ T \rightarrow T \cdot * F \\ \}$$

$$I_3 = \text{Goto}(I_0, F) = \{ \\ T \rightarrow F \cdot \\ \}$$

$$I_5 = \text{Goto}(I_0, \text{id}) = \{ \\ F \rightarrow \text{id} \cdot \\ \}$$

$$I_4 = \text{Goto}(I_0, "(") = \{ \\ F \rightarrow (\cdot E), \\ E \rightarrow \cdot E + T, \\ E \rightarrow \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot \text{id}, \\ \}$$

$$I_7 = \text{Goto}(I_2, *) = \{ \\ T \rightarrow T \cdot * F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot \text{id} \\ \}$$

Canonical Collection of Sets of LR(0) Items

$$I_6 = \text{Goto}(I_1, +) = \{$$
$$E \rightarrow E+.T,$$
$$T \rightarrow .T * F,$$
$$T \rightarrow .F,$$
$$F \rightarrow .(E),$$
$$F \rightarrow .\text{id},$$
$$\}$$
$$I_8 = \text{Goto}(I_4, E) = \{$$
$$E \rightarrow E.+T,$$
$$F \rightarrow (E.).$$
$$\}$$
$$I_9 = \text{Goto}(I_6, T) = \{$$
$$E \rightarrow E + T.,$$
$$T \rightarrow T.* F$$
$$\}$$
$$I_{10} = \text{Goto}(I_7, F) = \{$$
$$T \rightarrow T * F.,$$
$$\}$$
$$I_{11} = \text{Goto}(I_8,)) = \{$$
$$F \rightarrow (E).$$
$$\}$$
$$I_2 = \text{Goto}(I_4, T)$$
$$I_3 = \text{Goto}(I_4, F)$$
$$I_4 = \text{Goto}(I_4, "(")$$
$$I_5 = \text{Goto}(I_4, \text{id})$$
$$I_3 = \text{Goto}(I_6, F)$$
$$I_4 = \text{Goto}(I_6, "(")$$
$$I_5 = \text{Goto}(I_6, \text{id})$$
$$I_4 = \text{Goto}(I_7, "(")$$
$$I_5 = \text{Goto}(I_7, \text{id})$$
$$I_6 = \text{Goto}(I_8, +)$$
$$I_7 = \text{Goto}(I_9, *)$$