



Programming Language Syntax: Top-down Parsing

Read: Scott, Chapter 2.3.2 and 2.3.3

Lecture Outline

- Top-down parsing (also called LL parsing)
 - LL(1) parsing table
 - FIRST, FOLLOW, and PREDICT sets
 - LL(1) grammars
- Bottom-up parsing (also called LR parsing)
 - A brief overview, no detail

Recursive Descent

$\overline{start} \rightarrow expr \$\$$

$expr \rightarrow term\ term_tail$

$term \rightarrow id\ factor_tail$

$term_tail \rightarrow +\ term\ term_tail \mid \epsilon$

$factor_tail \rightarrow *\ id\ factor_tail \mid \epsilon$

start()

case lookahead() of

id: **expr()**; match(\$\$)

(\$\$ - end-of-input marker)

otherwise PARSE_ERROR

expr()

case lookahead() of

id: **term()**; **term_tail()**

otherwise PARSE_ERROR

term_tail()

case lookahead() of

+: match(' '); **term()**; **term_tail()**

\$\$: skip

otherwise: PARSE_ERROR

Predicting production $term_tail \rightarrow +\ term\ term_tail$

Predicting epsilon production $term_tail \rightarrow \epsilon$

Recursive Descent

$start \rightarrow expr \$\$$

$expr \rightarrow term term_tail$

$term \rightarrow id factor_tail$

$term_tail \rightarrow + term term_tail \mid \epsilon$

$factor_tail \rightarrow * id factor_tail \mid \epsilon$

term()

case lookahead() of

id: match('id'); factor_tail()

otherwise: PARSE_ERROR

factor_tail()

case lookahead() of

*****: match('*'); match('id'); factor_tail();

+, \$\$: skip

otherwise PARSE_ERROR

Predicting production $factor_tail \rightarrow *id factor_tail$

Predicting production $factor_tail \rightarrow \epsilon$

Recursive Descent Parsing

- Parse id + id * id \$\$

Start()

expr()

term()

id factor-tail() // return

term-tail()

+ term()

id factor-tail()

* id factor-tail()

...

LL(1) Parsing Table

- But how does the parser “predict”?
 - E.g., how does the parser know to expand a *factor_tail* by *factor_tail* $\rightarrow \epsilon$ on + and \$\$?
- It uses the LL(1) parsing table
 - One dimension is nonterminal to expand
 - Other dimension is lookahead token
 - We are interested in **one** token of lookahead
 - Entry “nonterminal on token” contains the production to apply or contains nothing

LL(1) Parsing Table

- One dimension: nonterminal to expand
- Other dimension: lookahead token

	a	
A	α	

- E.g., entry “nonterminal A on terminal \mathbf{a} ” contains production $A \rightarrow \alpha$
- Meaning: when parser is at nonterminal A and lookahead token is \mathbf{a} , then parser expands A by production $A \rightarrow \alpha$

LL(1) Parsing Table

Goal: Algorithm that constructs this LL(1) table for a given grammar.

$start \rightarrow expr \$\$$

$expr \rightarrow term \ term_tail$

$term \rightarrow \underline{id} \ factor_tail$

$term_tail \rightarrow + \ term \ term_tail \mid \epsilon$

$factor_tail \rightarrow * \ id \ factor_tail \mid \epsilon$

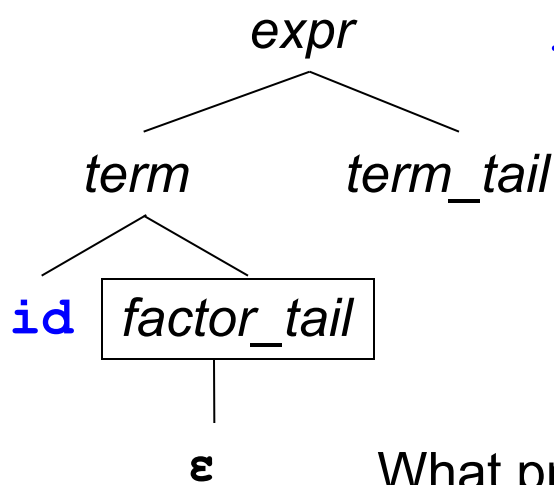
	<u>id</u>	+	*	\$\$
<u>start</u>	$expr \$\$$	—	—	—
expr	$term \ term_tail$	—	—	—
term_tail	—	$+ \ term \ term_tail$	—	ϵ
term	$\underline{id} \ factor_tail$	—	—	—
factor_tail	—	ϵ	$* \ id \ factor_tail$	ϵ

Intuition

■ Top-down parsing

- Parse tree is built from the top to the leaves
- Always expand the leftmost nonterminal

$expr \rightarrow term\ term_tail$
 $term_tail \rightarrow +\ term\ term_tail \mid \epsilon$
 $term \rightarrow id\ factor_tail$
 $factor_tail \rightarrow * id\ factor_tail \mid \epsilon$



id $+$ id $+$ $id * id$

$factor_tail \rightarrow * id factor_tail$ ~~X~~
 $factor_tail \rightarrow \epsilon$ ✓

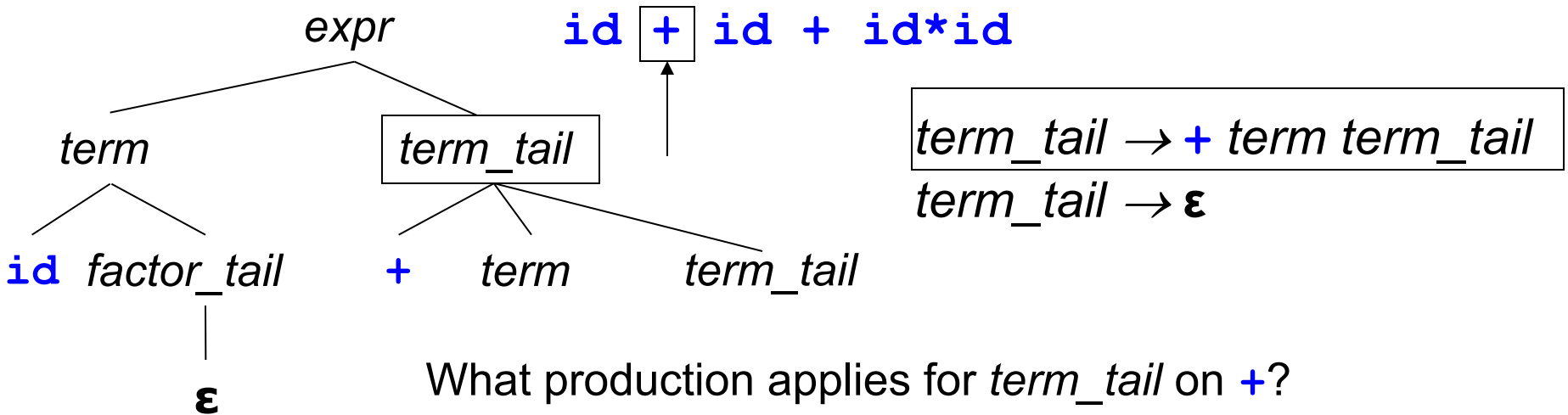
What production applies for $factor_tail$ on $+$?
 $+$ does not belong to an expansion of $factor_tail$.
However, $factor_tail$ has an epsilon production and $+$ belongs to an expansion of $term_tail$ which follows $factor_tail$. Thus, predict the epsilon production.

Intuition

$expr \rightarrow term\ term_tail$
 $term_tail \rightarrow +\ term\ term_tail \mid \epsilon$
 $term \rightarrow id\ factor_tail$
 $factor_tail \rightarrow * id\ factor_tail \mid \epsilon$

■ Top-down parsing

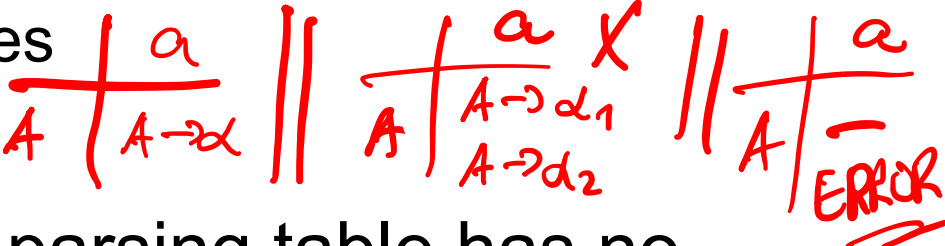
- Parse tree is built from the top to the leaves
- Always expand the leftmost nonterminal



What production applies for $term_tail$ on $+$?
 $+$ is the **first** symbol in expansions of $+\ term\ term_tail$.

Thus, predict production $term_tail \rightarrow +\ term\ term_tail$

LL(1) Tables and LL(1) Grammars

- We can construct an LL(1) parsing table for any context-free grammar
 - In general, the table will **contain multiply-defined entries**. That is, for some nonterminal and lookahead token, more than one production applies
- 
- A grammar whose LL(1) parsing table has no multiply-defined entries is said to be **LL(1) grammar**
 - LL(1) grammars are a very special subclass of context-free grammars. Why?

FIRST and FOLLOW sets

- Let α be any sequence of nonterminals and terminals
 - **FIRST**(α) is the set of terminals **a** that begin the strings derived from α . E.g., $expr \ \$\$ \Rightarrow^* \underline{id}...$, thus **id** in **FIRST**($expr \ \$\$$)
 - If there is a derivation $\alpha \Rightarrow^* \epsilon$, then ϵ is in **FIRST**(α)
- Let A be a nonterminal
 - **FOLLOW**(A) is the set of terminals **b** (including special end-of-input marker $\$$) that can appear immediately to the right of A in some sentential form:

$start \Rightarrow^* \dots A b \dots \Rightarrow^* \dots$

\downarrow
 ϵ

Computing FIRST

Notation:

α is an arbitrary sequence
of terminals and nonterminals

- Apply these rules until no more terminals or ϵ can be added to any $\text{FIRST}(\alpha)$ set

(1) If α starts with a terminal a , then $\text{FIRST}(\alpha) = \{ a \}$

(2) If α is a nonterminal X , where $X \rightarrow \epsilon$, then add ϵ to $\text{FIRST}(\alpha)$

(3) If α is a nonterminal $X \rightarrow \overbrace{Y_1 Y_2 \dots Y_k}^{\epsilon \quad b \dots}$ then add a to $\text{FIRST}(X)$ if for some i , a is in $\text{FIRST}(Y_i)$ and ϵ is in all of $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$. If ϵ is in all of $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_k)$, add ϵ to $\text{FIRST}(X)$.

- Everything in $\text{FIRST}(Y_1) - \{ \epsilon \}$ is surely in $\text{FIRST}(X)$
- If Y_1 does not derive ϵ , then we add nothing more; Otherwise, we add $\text{FIRST}(Y_2) - \{ \epsilon \}$, and so on

Similarly, if α is $Y_1 Y_2 \dots Y_k$, we'll repeat the above

Warm-up Exercise

$start \rightarrow expr \$\$$

$expr \rightarrow term \ term_tail$

$term \rightarrow id \ factor_tail$

$term_tail \rightarrow + \ term \ term_tail \mid \epsilon$

$factor_tail \rightarrow * \ id \ factor_tail \mid \epsilon$

$FIRST(term) = \{ id \}$

$FIRST(expr) = \{ id \}$

$FIRST(start) = \{ id \}$

$FIRST(term_tail) = \{ +, \epsilon \}$

$FIRST(+ \ term \ term_tail) = \{ + \}$

$FIRST(factor_tail) =$

$FIRST(A \ B)$

term-tail \Rightarrow

$+ \ term \ term_tail \Rightarrow$

$+ \ id \ factor_tail \ term_tail$

$\Rightarrow + \ id \ term_tail$

$\Rightarrow + \ id$

$+ \ term \ term_tail \Rightarrow$

$+ \ id \ factor_tail \ term_tail \Rightarrow$

$+ \ id \ term_tail \Rightarrow + \ id$

Exercise

$BCD \Rightarrow CD \Rightarrow D$
 $\Rightarrow WS$

$start \rightarrow S \$\$$	$B \rightarrow \underline{z} S \mid \epsilon$
$S \rightarrow \underline{x} S \mid A \underline{y}$	$C \rightarrow \underline{v} S \mid \epsilon$
$A \rightarrow \underline{BCD} \mid \epsilon$	$D \rightarrow \underline{w} S$

Compute FIRST sets:

$$FIRST(\underline{x} S) = \{x\}$$

$$FIRST(A \underline{y}) = \{z, v, w, y\}$$

$$FIRST(\underline{BCD}) = \{z, v, w\}$$

$$FIRST(\underline{z} S) = \{z\}$$

$$FIRST(\underline{v} S) = \{v\}$$

$$FIRST(\underline{w} S) = \{w\}$$

$$FIRST(S) = \{x, z, v, w, y\}$$

$$FIRST(A) = \{z, v, w, \epsilon\}$$

$$FIRST(B) = \{z, \epsilon\}$$

$$FIRST(C) = \{v, \epsilon\}$$

$$FIRST(D) = \{w\}$$

Computing FOLLOW

Notation:

A, B, S are nonterminals.

α, β are arbitrary sequences of terminals and nonterminals.

- Apply these rules until nothing can be added to any FOLLOW(A) set

(1) If there is a production $A \rightarrow \alpha B \beta$, then everything in FIRST(β) except for ϵ should be added to FOLLOW(B)

start \Rightarrow^ ... A ... \Rightarrow ... $\alpha B \beta$... \Rightarrow^* ... $\alpha B b$...*

(2) If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where FIRST(β) contains ϵ , then everything in FOLLOW(A) should be added to FOLLOW(B)

start \Rightarrow^ ... A b ... \Rightarrow ... $\alpha B b$...*

Warm-up

$start \rightarrow expr \$\$$

$expr \rightarrow term \text{ term_tail}$

$term \rightarrow id \text{ factor_tail}$

$term_tail \rightarrow + \text{ term term_tail} \mid \epsilon$

$factor_tail \rightarrow * \text{ id factor_tail} \mid \epsilon$

$FOLLOW(expr) = \{ \$\$ \}$

$FOLLOW(term_tail) = \{ \$\$ \}$

$FOLLOW(term) = \{ \$\$, + \}$

$FOLLOW(\text{factor_tail}) = \{ \$\$, + \}$

$expr \$\$ \Rightarrow term \text{ term_tail} \$\$ \Rightarrow \dots + \text{ term term_tail} \$\$$

$\Rightarrow \dots + term + term \text{ term_tail} \$\$$

Exercise

$start \rightarrow S \ \$\ \$$
 $S \rightarrow \mathbf{x} S \mid A \mathbf{y}$
 $A \rightarrow BCD \mid \epsilon$

$B \rightarrow \mathbf{z} S \mid \epsilon$
 $C \rightarrow \mathbf{v} S \mid \epsilon$
 $D \rightarrow \mathbf{w} S$

Compute FOLLOW sets:

$FOLLOW(A) = \{y\}$

$FOLLOW(B) = \{v, w\}$

$FOLLOW(C) = \{w\}$

$FOLLOW(D) = \{y\}$

$FOLLOW(S) = \{\$, v, w, y\}$

PREDICT Sets

$$\text{PREDICT}(A \rightarrow \alpha) = \begin{cases} \text{FIRST}(\alpha) & \text{if } \alpha \text{ does not derive } \epsilon \\ (\text{FIRST}(\alpha) - \{\epsilon\}) \cup \text{FOLLOW}(A) & \text{if } \alpha \text{ derives } \epsilon \end{cases}$$

Constructing LL(1) Parsing Table

- Algorithm uses PREDICT sets:

```
foreach production  $A \rightarrow \alpha$  in grammar  $G$ 
  foreach terminal  $a$  in PREDICT( $A \rightarrow \alpha$ )
    add  $A \rightarrow \alpha$  into entry parse_table[A,a]
```

- If each entry in `parse_table` contains at most one production, then G is said to be LL(1)

Exercise

$start \rightarrow S \$\$$	$B \rightarrow z S \mid \epsilon$
$S \rightarrow x S \mid A y$	$C \rightarrow v S \mid \epsilon$
$A \rightarrow BCD \mid \epsilon$	$D \rightarrow w S$

Compute PREDICT sets:

$PREDICT(S \rightarrow x S) =$

$PREDICT(S \rightarrow A y) =$

$PREDICT(A \rightarrow BCD) = \{z, v, w\}$

$PREDICT(A \rightarrow \epsilon) = \{y\}$

... etc...

$PREDICT(A \rightarrow BCD) \cap$

$PREDICT(A \rightarrow \epsilon) = \emptyset$

Writing an LL(1) Grammar

- Most context-free grammars are not LL(1) grammars
- Obstacles to LL(1)-ness

- Left recursion is an obstacle. Why?

$$\begin{array}{l} \text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term} \\ \text{term} \rightarrow \text{term} * \text{id} \mid \text{id} \end{array}$$

Handwritten annotations in red: A bracket above the first rule groups "expr + term" and "term". A bracket above the second rule groups "term * id" and "id".

- Common prefixes are an obstacle. Why?

$$\begin{array}{l} \text{stmt} \rightarrow \text{if } b \text{ then } \text{stmt} \text{ else } \text{stmt} \mid \\ \text{if } b \text{ then } \text{stmt} \mid \\ a \end{array}$$

Removal of Left Recursion

- Left recursion can be removed from a grammar mechanically
- Started from this left-recursive expression grammar:

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{term} \mid \text{term} \\ \text{term} &\rightarrow \text{term} * \text{id} \mid \text{id} \end{aligned}$$

- After removal of left recursion, we obtain this equivalent grammar, which is LL(1):

$$\begin{aligned} \text{expr} &\rightarrow \text{term} \text{term_tail} \\ \text{term_tail} &\rightarrow + \text{term} \text{term_tail} \mid \epsilon \\ \text{term} &\rightarrow \text{id} \text{factor_tail} \\ \text{factor_tail} &\rightarrow * \text{id} \text{factor_tail} \mid \epsilon \end{aligned}$$

Removal of Common Prefixes

- Common prefixes can be removed mechanically as well by using **left-factoring**
- Original if-then-else grammar:

$stmt \rightarrow \underline{if\ b\ then\ stmt}\ \boxed{else\ stmt} \mid$
 $\underline{if\ b\ then\ stmt} \boxed{\mid}$
 a

- After left-factoring:

● $stmt \rightarrow \underline{if\ b\ then\ stmt}\ \boxed{else_part} \mid a$
 $else_part \rightarrow \underline{else\ stmt} \mid \epsilon$

Exercise

$start \rightarrow stmt \$\$$

$stmt \rightarrow \mathbf{if\ b\ then\ stmt\ else_part} \mid \mathbf{a}$

$else_part \rightarrow \mathbf{else\ stmt} \mid \epsilon$

- Compute FIRSTs:

$FIRST(stmt\ \$\$)$, $FIRST(\mathbf{if\ b\ then\ stmt\ else_part})$,
 $FIRST(\mathbf{a})$, $FIRST(\mathbf{else\ stmt})$

- Compute FOLLOW:

$FOLLOW(else_part)$

- Compute PREDICT sets for all 5 productions
- Construct the LL(1) parsing table. Is this grammar
an LL(1) grammar?

Exercise

$start \rightarrow stmt \$\$$

$stmt \rightarrow \underline{if} \ b \ \underline{then} \ stmt \ \underline{else_part} \ | \ \underline{a}$

$\underline{else_part} \rightarrow \underline{else} \ stmt \ | \ \epsilon$

- Compute FIRSTs:

$$FIRST(stmt \$\$) = \{if, a\}$$

$$FIRST(if \ b \ then \ stmt \ else_part) = \{if\}$$

$$FIRST(a) = \{a\}$$

$$FIRST(else \ stmt) = \{else\}$$

Exercise

$start \rightarrow \underline{stmt} \underline{\$ \$}$
 $stmt \rightarrow \text{if } b \text{ then } stmt \text{ else_part } | a$
 $\text{else_part} \rightarrow \text{else } stmt | \epsilon$

- Compute FOLLOW:

$FOLLOW(stmt) =$
 $FOLLOW(\text{else_part}) = \{ \underline{\$ \$}, \text{else} \}$ $FOLLOW(\text{else_part})$

	<u>else</u>
<u>else_part</u>	else part $\rightarrow \epsilon$
	else part $\rightarrow \underline{\text{else } stmt}$

Exercise

$start \rightarrow stmt \$\$$

$stmt \rightarrow \text{if } b \text{ then } stmt \text{ else_part} \mid a$

$\text{else_part} \rightarrow \text{else } stmt \mid \epsilon$

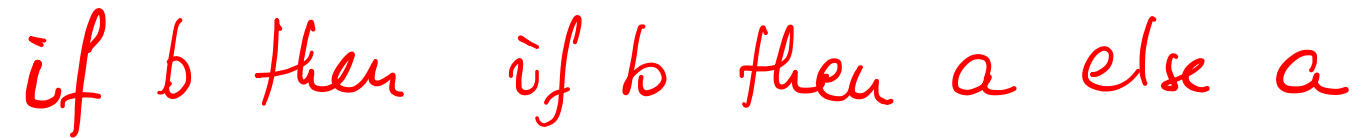
- Construct the LL(1) parsing table

- Is this grammar an LL(1) grammar?

No!

Exercise

if b then if b then a else a



Lecture Outline

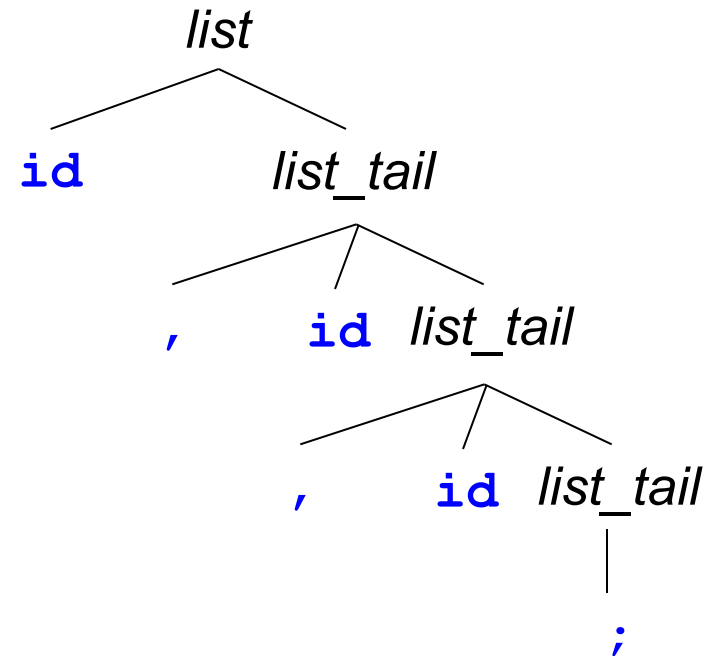
- Top-down parsing (also called LL parsing)
 - LL(1) parsing table
 - FIRST, FOLLOW, and PREDICT sets
 - LL(1) grammars
- Bottom-up parsing (also called LR parsing)
 - A brief overview, no detail

Bottom-up Parsing

- Terminals are seen in the order of appearance in the token stream

id , **id** , **id** ;
↑ ↑ ↑ ↑ ↑ ↑

- Parse tree is constructed
 - From the leaves to the top
 - A rightmost derivation in reverse



$list \rightarrow id\ list_tail$
 $list_tail \rightarrow ,\ id\ list_tail\ | ;$

Bottom-up Parsing

$list \rightarrow id\ list_tail$
 $list_tail \rightarrow ,\ id\ list_tail \mid ;$

Stack	Input	Action
	<code>id, id, id;</code>	shift
<code>id</code>	<code>, id, id;</code>	shift
<code>id,</code>	<code>id, id;</code>	shift
<code>id, id</code>	<code>, id;</code>	shift
<code>id, id,</code>	<code>id;</code>	shift
<code>id, id, id</code>	<code>;</code>	shift
<code>id, id, id,</code>		reduce by $list_tail \rightarrow ;$

Bottom-up Parsing

$list \rightarrow id\ list_tail$
 $list_tail \rightarrow ,\ id\ list_tail \mid ;$

Stack

Input

Action

$id, id, id\ list_tail$

reduce by

$list_tail \rightarrow ,\ id\ list_tail$

$id, id\ list_tail$

reduce by

$list_tail \rightarrow ,\ id\ list_tail$

$id\ list_tail$

reduce by

$list \rightarrow id\ list_tail$

$list$

ACCEPT

Bottom-up Parsing

- Also called LR parsing
- LR parsers work with LR(k) grammars
 - L stands for “left-to-right” scan of input
 - R stands for “rightmost” derivation
 - k stands for “need k tokens of lookahead”
- We are interested in LR(0) and LR(1) and variants in between
- LR parsing is better than LL parsing!
 - Accepts larger class of languages
 - Just as efficient!

LR Parsing

- The parsing method used in practice
 - LR parsers recognize virtually all PL constructs
 - LR parsers recognize a much larger set of grammars than predictive parsers
 - LR parsing is efficient
- LR parsing variants
 - SLR (or Simple LR)
 - LALR (or Lookahead LR) – `yacc/bison` generate LALR parsers
 - LR (Canonical LR)
 - $SLR < LALR < LR$

Main Idea

- Stack \leftarrow Input
- Stack: holds the part of the input seen so far
 - A string of both terminals and nonterminals
- Input: holds the remaining part of the input
 - A string of terminals
- Parser performs two actions
 - **Reduce**: parser pops a “suitable” production right-hand-side off top of stack, and pushes production’s left-hand-side on the stack
 - **Shift**: parser pushes next terminal from the input on top of the stack

Example

- Recall the grammar

$expr \rightarrow expr + term \mid term$
 $term \rightarrow term * id \mid id$

- This is not LL(1) because it is left recursive
- LR parsers can handle left recursion!

- Consider string

$id + id * id$

id + id*id

Stack	Input	Action
	id+id*id	shift id
<u>id</u>	+id*id	reduce by <i>term</i> → id
<u>term</u>	+id*id	reduce by <i>expr</i> → <i>term</i>
<u>expr</u>	+id*id	shift +
expr +	id*id	shift id
expr + <u>id</u>	* id	reduce by <i>term</i> → id

$expr \rightarrow expr + term \mid term$
 $term \rightarrow term * id \mid id$

id + id*id

Stack	Input	Action
<i>expr+term</i>	<i>*id</i>	shift <i>*</i>
<i>expr+term*</i>	<i>id</i>	shift <i>id</i>
<i>expr+<u>term*id</u></i>		reduce by <i>term</i> → <i>term *id</i>
<i><u>expr+term</u></i>		reduce by <i>expr</i> → <i>expr+term</i>
<i>expr</i>		ACCEPT, SUCCESS

expr → *expr + term* | *term*
term → *term * id* | *id*

*id + id*id*

Sequence of reductions performed by parser

↑
*id+id*id*
*term+id*id*
*expr+id*id*
*expr+term*id*
expr+term
expr

- A rightmost derivation in reverse
- The stack (e.g., *expr*) concatenated with remaining input (e.g., *+id*id*) gives a sentential form (*expr+id*id*) in the rightmost derivation

expr → *expr + term* | *term*
term → *term * id* | *id*

The End
