

# Agenda

## Parsers

- Terminology
- LL(1) Parsers
- Overview of LR Parsing

Wednesday, January 19, 2011

Wednesday, January 19, 2011

## Terminology

- Grammar  $G = (V_t, V_n, S, P)$
- $V_t$  is the set of *terminals*
- $V_n$  is the set of *non-terminals*
- $S$  is the *start symbol*
- $P$  is the set of *productions*
  - Each production takes the form:  $V_n \rightarrow \lambda \mid (V_n \mid V_t)^+$
  - Grammar is *context-free* (why?)
- A simple grammar:  
 $G = (\{a, b\}, \{S, A, B\}, \{S \rightarrow A B \$, A \rightarrow A a, A \rightarrow a, B \rightarrow B b, B \rightarrow b\}, S)$

Wednesday, January 19, 2011

## Terminology

- $V$  is the *vocabulary* of a grammar, consisting of terminal ( $V_t$ ) and non-terminal ( $V_n$ ) symbols
- For our sample grammar
  - $V_n = \{S, A, B\}$ 
    - Non-terminals are symbols on the LHS of a production
    - Non-terminals are constructs in the language that are recognized during parsing
  - $V_t = \{a, b\}$ 
    - Terminals are the tokens recognized by the scanner
    - They correspond to symbols in the text of the program

Wednesday, January 19, 2011

## Terminology

- Productions (rewrite rules) tell us how to derive strings in the language
  - Apply productions to rewrite strings into other strings
- We will use the standard BNF form
- $P = \{$   
 $S \rightarrow A B \$$   
 $A \rightarrow A a$   
 $A \rightarrow a$   
 $B \rightarrow B b$   
 $B \rightarrow b$   
 $\}$

Wednesday, January 19, 2011

## Generating strings

$S \rightarrow A B \$$

$A \rightarrow A a$

$A \rightarrow a$

$B \rightarrow B b$

$B \rightarrow b$

- Given a start rule, productions tell us how to rewrite a non-terminal into a different set of symbols
- By convention, first production applied has the start symbol on the left, and there is only one such production

To derive the string "a a b b b" we can do the following rewrites:

$S \Rightarrow A B \$ \Rightarrow A a B \$ \Rightarrow a a B \$ \Rightarrow a a B b \$ \Rightarrow a a B b b \$ \Rightarrow a a b b b \$$

Wednesday, January 19, 2011

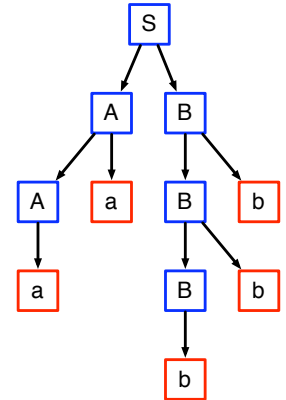
## Terminology

- *Strings* are composed of symbols
  - A A a a B b b A a is a string
  - We will use Greek letters to represent strings composed of both terminals and non-terminals
- $L(G)$  is the language produced by the grammar  $G$ 
  - All strings consisting of only terminals that can be produced by  $G$
  - In our example,  $L(G) = a+b+\$$
  - All regular expressions can be expressed as grammars for context-free languages, but not vice-versa
  - Consider:  $a^i b^j \$$  (what is the grammar for this?)

Wednesday, January 19, 2011

## Parse trees

- Tree which shows how a string was produced by a language
- Interior nodes of tree: non-terminals
  - Children: the terminals and non-terminals generated by applying a production rule
- Leaf nodes: terminals



Wednesday, January 19, 2011

## Leftmost derivation

- Rewriting of a given string starts with the leftmost symbol
- Exercise: do a leftmost derivation of the input program

$F(V + V)$

using the following grammar:

E	→	Prefix (E)
E	→	V Tail
Prefix	→	F
Prefix	→	$\lambda$
Tail	→	+ E
Tail	→	$\lambda$

- What does the parse tree look like?

Wednesday, January 19, 2011

## Rightmost derivation

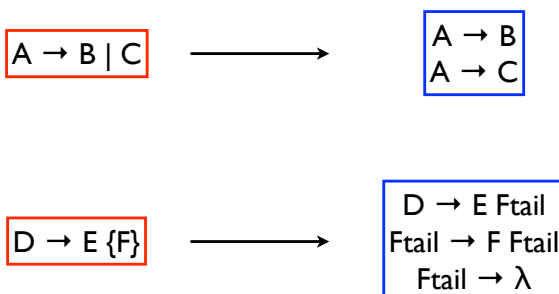
- Rewrite using the rightmost non-terminal, instead of the left
- What is the rightmost derivation of this string?

$F(V + V)$

E	→	Prefix (E)
E	→	V Tail
Prefix	→	F
Prefix	→	$\lambda$
Tail	→	+ E
Tail	→	$\lambda$

Wednesday, January 19, 2011

## Simple conversions



Wednesday, January 19, 2011

## Top-down vs. Bottom-up parsers

- Top-down parsers use left-most derivation
- Bottom-up parsers use right-looking parse
- Notation:
  - LL(1): Leftmost derivation with 1 symbol lookahead
  - LL(k): Leftmost derivation with k symbols lookahead
  - LR(1): Right-looking derivation with 1 symbol lookahead

Wednesday, January 19, 2011

# What is parsing

- Parsing is recognizing members in a language specified/defined/generated by a grammar
- When a construct (corresponding to a production in a grammar) is recognized, a typical parser will take some action
- In a compiler, this action generates an intermediate representation of the program construct
- In an interpreter, this action might be to perform the action specified by the construct. Thus, if  $a+b$  is recognized, the value of  $a$  and  $b$  would be added and placed in a temporary variable

Wednesday, January 19, 2011

# Another simple grammar

PROGRAM  $\rightarrow$  **begin** STMTLIST \$

STMTLIST  $\rightarrow$  STMT ; STMTLIST

STMTLIST  $\rightarrow$  **end**

STMT  $\rightarrow$  **id**

STMT  $\rightarrow$  **if ( id )** STMTLIST

- A sentence in the grammar:  
begin if (id) if (id) id ; end; end; end; \$
- What are the terminals and non-terminals of this grammar?

Wednesday, January 19, 2011

## Parsing this grammar

PROGRAM  $\rightarrow$  **begin** STMTLIST \$

STMTLIST  $\rightarrow$  STMT ; STMTLIST

STMTLIST  $\rightarrow$  **end**

STMT  $\rightarrow$  **id**

STMT  $\rightarrow$  **if ( id )** STMTLIST

- Note
  - To parse STMT in STMTLIST  $\rightarrow$  STMT ; STMTLIST, it is necessary to choose between either STMT  $\rightarrow$  **id** or STMT  $\rightarrow$  **if ...**
  - Choose the production to parse by finding out if next token is **if** or **id**
    - i.e., which production the next input token matches
    - This is the **first** set of the production

Wednesday, January 19, 2011

## Another example

$S \rightarrow A B \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow \lambda$

$B \rightarrow b$

- Consider  $S \Rightarrow A B \$ \Rightarrow x a A B \$ \Rightarrow x a B \$ \Rightarrow x a b \$$
- When parsing  $x a b \$$  we know from the goal production we need to match an A. The next token is  $x$ , so we apply  $A \rightarrow x a A$
- The parser matches  $x$ , matches  $a$  and now needs to parse A again
- How do we know which A to use? We need to use  $A \rightarrow \lambda$ 
  - When matching the right hand side of  $A \rightarrow \lambda$ , the next token comes from a non-terminal that follows A (i.e., it must be  $b$ )
  - Tokens that can follow A are called the **follow** set of A

Wednesday, January 19, 2011

## First and follow sets

- $\text{First}(\alpha)$ : the set of terminals that begin all strings that can be derived from  $\alpha$ 
  - $\text{First}(A) = \{x, y\}$
  - $\text{First}(xaA) = \{x\}$
  - $\text{First}(AB) = \{x, y, b\}$
- $\text{Follow}(A)$ : the set of terminals that can appear immediately after A in some partial derivation
  - $\text{Follow}(A) = \{b\}$

$S \rightarrow A B \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow \lambda$

$B \rightarrow b$

## First and follow sets

- $\text{First}(\alpha) = \{a \in V_t \mid \alpha \Rightarrow^* a\beta\} \cup \{\lambda \mid \text{if } \alpha \Rightarrow^* \lambda\}$
- $\text{Follow}(A) = \{a \in V_t \mid S \Rightarrow^* \dots Aa \dots\} \cup \{\$ \mid \text{if } S \Rightarrow^* \dots A \$\}$

S: start symbol  
a: a terminal symbol  
A: a non-terminal symbol  
 $\alpha, \beta$ : a string composed of terminals and non-terminals (typically,  $\alpha$  is the RHS of a production)

$\Rightarrow$ : derived in 1 step  
 $\Rightarrow^*$ : derived in 0 or more steps  
 $\Rightarrow^+$ : derived in 1 or more steps

Wednesday, January 19, 2011

Wednesday, January 19, 2011

## Computing first sets

- Terminal:  $\text{First}(a) = \{a\}$
- Non-terminal:  $\text{First}(A)$ 
  - Look at all productions for A
$$A \rightarrow X_1 X_2 \dots X_k$$
  - $\text{First}(A) \supseteq (\text{First}(X_1) - \lambda)$
  - If  $\lambda \in \text{First}(X_1)$ ,  $\text{First}(A) \supseteq (\text{First}(X_2) - \lambda)$
  - If  $\lambda$  is in  $\text{First}(X_i)$  for all  $i$ , then  $\lambda \in \text{First}(A)$
- Computing  $\text{First}(\alpha)$ : similar procedure to computing  $\text{First}(A)$

Wednesday, January 19, 2011

## Exercise

- What are the first sets for all the non-terminals in following grammar:

$$S \rightarrow A B \$$$
$$A \rightarrow x a A$$
$$A \rightarrow y a A$$
$$A \rightarrow \lambda$$
$$B \rightarrow b$$
$$B \rightarrow A$$

Wednesday, January 19, 2011

## Computing follow sets

- $\text{Follow}(S) = \{\$ \}$
- To compute  $\text{Follow}(A)$ :
  - Find productions which have A on rhs. Three rules:
    1.  $X \rightarrow \alpha A \beta$ :  $\text{Follow}(A) \supseteq (\text{First}(\beta) - \lambda)$
    2.  $X \rightarrow \alpha A \beta$ : If  $\lambda \in \text{First}(\beta)$ ,  $\text{Follow}(A) \supseteq \text{Follow}(X)$
    3.  $X \rightarrow \alpha A$ :  $\text{Follow}(A) \supseteq \text{Follow}(X)$
- Note:  $\text{Follow}(X)$  never has  $\lambda$  in it.

Wednesday, January 19, 2011

## Exercise

- What are the follow sets for

$$S \rightarrow A B \$$$
$$A \rightarrow x a A$$
$$A \rightarrow y a A$$
$$A \rightarrow \lambda$$
$$B \rightarrow b$$
$$B \rightarrow A$$

Wednesday, January 19, 2011

## Towards parser generators

- Key problem: as we read the source program, we need to decide what productions to use
- Step 1: find the tokens that can tell which production P (of the form  $A \rightarrow X_1 X_2 \dots X_m$ ) applies

$\text{Predict}(P) =$

$$\begin{cases} \text{First}(X_1 \dots X_m) & \text{if } \lambda \notin \text{First}(X_1 \dots X_m) \\ (\text{First}(X_1 \dots X_m) - \lambda) \cup \text{Follow}(A) & \text{otherwise} \end{cases}$$

- If next token is in  $\text{Predict}(P)$ , then we should choose this production

Wednesday, January 19, 2011

## Parse tables

- Step 2: build a parse table
  - Given some non-terminal  $V_n$  (the non-terminal we are currently processing) and a terminal  $V_t$  (the lookahead symbol), the parse table tells us which production P to use (or that we have an error)
- More formally:
$$T: V_n \times V_t \rightarrow P \cup \{\text{Error}\}$$

Wednesday, January 19, 2011

## Building the parse table

- Start:  $T[A][t] = \text{//initialize all fields to "error"}$

foreach A:

foreach P with A on its lhs:

foreach t in Predict(P):

$T[A][t] = P$

- Exercise: build parse table for our toy grammar

1.  $S \rightarrow A B \$$   
 2.  $A \rightarrow x a A$   
 3.  $A \rightarrow y a A$   
 4.  $A \rightarrow \lambda$   
 5.  $B \rightarrow b$

Wednesday, January 19, 2011

## Recursive-descent parsers

- Given the parse table, we can create a program which generates recursive descent parsers
- Remember the recursive descent parser we saw for MICRO
- If the choice of production is not unique, the parse table tells us which one to take
- However, there is an easier method!

Wednesday, January 19, 2011

## Stack-based parser for LL(1)

- Given the parse table, a stack-based algorithm is much simpler to generate than a recursive descent parser

- Basic algorithm:

1. Push the RHS of a production onto the stack
2. Pop a symbol, if it is a terminal, match it
3. If it is a non-terminal, take its production according to the parse table and go to 1

- Algorithm on page 121
- Note: always start with start state

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

1.  $S \rightarrow A B \$$   
 2.  $A \rightarrow x a A$   
 3.  $A \rightarrow y a A$   
 4.  $A \rightarrow \lambda$   
 5.  $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

1.  $S \rightarrow A B \$$   
 2.  $A \rightarrow x a A$   
 3.  $A \rightarrow y a A$   
 4.  $A \rightarrow \lambda$   
 5.  $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

1.  $S \rightarrow A B \$$   
 2.  $A \rightarrow x a A$   
 3.  $A \rightarrow y a A$   
 4.  $A \rightarrow \lambda$   
 5.  $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)
A B \$	$y a b \$$	predict 3

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)
A B \$	$y a b \$$	predict 3
$y a A B \$$	$y a b \$$	match(y)

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)
A B \$	$y a b \$$	predict 3
$y a A B \$$	$y a b \$$	match(y)
a A B \$	$a b \$$	match(a)

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)
A B \$	$y a b \$$	predict 3
$y a A B \$$	$y a b \$$	match(y)
a A B \$	$a b \$$	match(a)
A B \$	$b \$$	predict 4

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	$x a y a b \$$	predict 1
A B \$	$x a y a b \$$	predict 2
$x a A B \$$	$x a y a b \$$	match(x)
a A B \$	$a y a b \$$	match(a)
A B \$	$y a b \$$	predict 3
$y a A B \$$	$y a b \$$	match(y)
a A B \$	$a b \$$	match(a)
A B \$	$b \$$	predict 4
B \$	$b \$$	predict 5

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	x a y a b \$	predict 1
A B \$	x a y a b \$	predict 2
x a A B \$	x a y a b \$	match(x)
a A B \$	a y a b \$	match(a)
A B \$	y a b \$	predict 3
y a A B \$	y a b \$	match(y)
a A B \$	a b \$	match(a)
A B \$	b \$	predict 4
B \$	b \$	predict 5
b \$	b \$	match(b)

Wednesday, January 19, 2011

## An example

- How would a stack-based parser parse:

$x a y a b$

- $S \rightarrow A B \$$
- $A \rightarrow x a A$
- $A \rightarrow y a A$
- $A \rightarrow \lambda$
- $B \rightarrow b$

Parse stack	Remaining input	Parser action
S	x a y a b \$	predict 1
A B \$	x a y a b \$	predict 2
x a A B \$	x a y a b \$	match(x)
a A B \$	a y a b \$	match(a)
A B \$	y a b \$	predict 3
y a A B \$	y a b \$	match(y)
a A B \$	a b \$	match(a)
A B \$	b \$	predict 4
B \$	b \$	predict 5
b \$	b \$	match(b)
\$	\$	Done!

Wednesday, January 19, 2011

## LL(k) parsers

- Can use similar techniques for LL(k) parsers
- Use more than one symbol of look-ahead to distinguish productions
- Why might this be bad?

Wednesday, January 19, 2011

## Dealing with semantic actions

- Recall: we can annotate a grammar with *action symbols*
  - Tell the parser to invoke a semantic action routine
- Can simply push action symbols onto stack as well
- When popped, the semantic action routine is called

Wednesday, January 19, 2011

## Non-LL(1) grammars

- Not all grammars are LL(1)!
- Consider
  - $\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt list} \rangle \text{ endif}$
  - $\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt list} \rangle \text{ else } \langle \text{stmt list} \rangle \text{ endif}$
- This is not LL(1) (why?)
- We can turn this in to
  - $\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt list} \rangle \langle \text{if suffix} \rangle$
  - $\langle \text{if suffix} \rangle \rightarrow \text{endif}$
  - $\langle \text{if suffix} \rangle \rightarrow \text{else } \langle \text{stmt list} \rangle \text{ endif}$

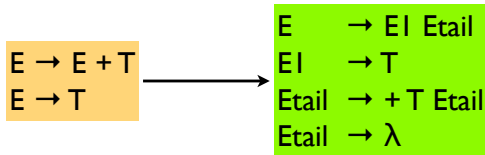
Wednesday, January 19, 2011

## Left recursion

- Left recursion* is a problem for LL(1) parsers
  - LHS is also the first symbol of the RHS
- Consider:
  - $E \rightarrow E + T$
- What would happen with the stack-based algorithm?

Wednesday, January 19, 2011

## Removing left recursion



Algorithm on page 125

Wednesday, January 19, 2011

## Are all grammars LL(k)?

- No! Consider the following grammar:

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow (E + E) \\ E &\rightarrow (E - E) \\ E &\rightarrow x \end{aligned}$$

- When parsing E, how do we know whether to use rule 2 or 3?
- Potentially unbounded number of characters before the distinguishing '+' or '-' is found
- No amount of lookahead will help!

Wednesday, January 19, 2011

## In real languages?

- Consider the if-then-else problem
- if x then y else z
- Problem: else is optional
- if a then if b then c else d
  - Which if does the else belong to?
- This is analogous to a "bracket language":  $[ \ ]^i$  ( $i \geq j$ )

$$\begin{aligned} S &\rightarrow [ S C \\ S &\rightarrow \lambda \\ C &\rightarrow ] \\ C &\rightarrow \lambda \end{aligned}$$

[ ] can be parsed:  $SS\lambda C$  or  $SSC\lambda$   
(it's ambiguous!)

Wednesday, January 19, 2011

## Solving the if-then-else problem

- The ambiguity exists at the language level. To fix, we need to define the semantics properly
- " ] " matches nearest unmatched "[ "
- This is the rule C uses for if-then-else
- What if we try this?

$$\begin{aligned} S &\rightarrow [ S \\ S &\rightarrow S I \\ S I &\rightarrow [ S I ] \\ S I &\rightarrow \lambda \end{aligned}$$

This grammar is still not LL(1)  
(or LL(k) for any k!)

Wednesday, January 19, 2011

## Two possible fixes

- If there is an ambiguity, prioritize one production over another
- e.g., if C is on the stack, always match " ] " before matching "  $\lambda$  "

$$\begin{aligned} S &\rightarrow [ S C \\ S &\rightarrow \lambda \\ C &\rightarrow ] \\ C &\rightarrow \lambda \end{aligned}$$

- Another option: change the language!
- e.g., all if-statements need to be closed with an endif

$$\begin{aligned} S &\rightarrow \text{if } S \text{ E} \\ S &\rightarrow \text{other} \\ E &\rightarrow \text{else } S \text{ endif} \\ E &\rightarrow \text{endif} \end{aligned}$$

Wednesday, January 19, 2011

## Parsing if-then-else

- What if we don't want to change the language?
- C does not require { } to delimit single-statement blocks
- To parse if-then-else, we need to be able to look ahead at the entire rhs of a production before deciding which production to use
- In other words, we need to determine how many " ] " to match before we start matching "[ "s
- LR parsers can do this!

Wednesday, January 19, 2011



## LR Parsers

- Parser which does a **Left-to-right, Right-most** derivation
  - Rather than parse top-down, like LL parsers do, parse bottom-up, starting from leaves
- Basic idea: put tokens on a stack until an entire production is found
- Issues:
  - Recognizing the endpoint of a production
  - Finding the length of a production (RHS)
  - Finding the corresponding nonterminal (the LHS of the production)

Wednesday, January 19, 2011

## Data structures

- At each state, given the next token,
  - A **goto table** defines the successor state
  - An **action table** defines whether to
    - shift** – put the next state and token on the stack
    - reduce** – an RHS is found; process the production
    - terminate** – parsing is complete

Wednesday, January 19, 2011

## Example

- Consider the simple grammar:
  - `<program>` → `begin <stmts> end $`
  - `<stmts>` → `SimpleStmt ; <stmts>`
  - `<stmts>` → `begin <stmts> end ; <stmts>`
  - `<stmts>` → `λ`
- Shift-reduce driver algorithm on page 142

Wednesday, January 19, 2011

## Action and goto tables

	begin	end	:	SimpleStmt	\$	<program>	<stmts>
0	S / 1						
1	S / 4	R4		S / 5			S / 2
2		S / 3					
3					A		
4	S / 4	R4		S / 5			S / 7
5			S / 6				
6	S / 4	R4		S / 5			S / 10
7		S / 8					
8			S / 9				
9	S / 4	R4		S / 6			S / 11
10		R2					
11		R3					

Wednesday, January 19, 2011

## Example

- Parse: `begin SimpleStmt ; SimpleStmt ; end $`

Step	Parse Stack	Remaining Input	Parser Action
1	0	begin S ; S ; end \$	Shift 1
2	0 1	S ; S ; end \$	Shift 5
3	0 1 5	; S ; end \$	Shift 6
4	0 1 5 6	S ; end \$	Shift 5
5	0 1 5 6 5	; end \$	Shift 6
6	0 1 5 6 5 6	end \$	Reduce 4 (goto 10)
7	0 1 5 6 5 6 10	end \$	Reduce 2 (goto 10)
8	0 1 5 6 10	end \$	Reduce 2 (goto 2)
9	0 1 2	end \$	Shift 3
10	0 1 2 3	\$	Accept

Wednesday, January 19, 2011

## LR Parsers

- Basic idea:
  - shift** tokens onto the stack. At any step, keep the set of productions that could generate the read-in tokens
  - reduce** the RHS of recognized productions to the corresponding non-terminal on the LHS of the production. Replace the RHS tokens on the stack with the LHS non-terminal.

Wednesday, January 19, 2011

# LR(k) parsers

- LR(0) parsers
  - No lookahead
  - Predict which action to take by looking only at the symbols currently on the stack
- LR(k) parsers
  - Can look ahead  $k$  symbols
  - Most powerful class of deterministic bottom-up parsers
  - LR(1) and variants are the most common parsers

Wednesday, January 19, 2011

# Terminology for LR parsers

- Configuration: a production augmented with a “•”
 
$$A \rightarrow X_1 \dots X_i \bullet X_{i+1} \dots X_j$$
- The “•” marks the point to which the production has been recognized. In this case, we have recognized  $X_1 \dots X_i$
- Configuration set: all the configurations that can apply at a given point during the parse:
 
$$A \rightarrow B \bullet CD$$

$$A \rightarrow B \bullet GH$$

$$T \rightarrow B \bullet Z$$
- Idea: every configuration in a configuration set is a production that can possibly be matched

Wednesday, January 19, 2011

## Configuration closure set

- Include all the configurations necessary to recognize the next symbol after the •

$\text{closure0}(\text{configuration\_set})$  defined on page 146

- Example:

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

$\text{closure0}(\{S \rightarrow \bullet E \$\}) = \{$   
 $S \rightarrow \bullet E \$$   
 $E \rightarrow \bullet E + T$   
 $E \rightarrow \bullet T$   
 $T \rightarrow \bullet ID$   
 $T \rightarrow \bullet (E)$   
 $\}$

Wednesday, January 19, 2011

## Successor configuration set

- Starting with the initial configuration set

$s_0 = \text{closure0}(\{S \rightarrow \bullet \alpha \$\})$

an LR(0) parser will find the successor given the next symbol  $X$

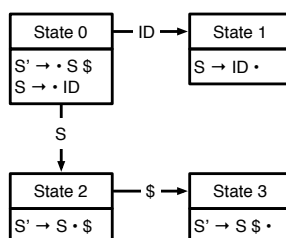
- $X$  can be either a terminal (the next token from the scanner) or a non-terminal (the result of applying a reduction)
- Determining the successor  $s' = \text{go\_to0}(s, X)$ :
  - For each configuration in  $s$  of the form  $A \rightarrow \beta \bullet X \gamma$  add  $A \rightarrow \beta X \bullet \gamma$  to  $t$
  - $s' = \text{closure0}(t)$

Wednesday, January 19, 2011

## CFSM

- CFSM = Characteristic Finite State Machine
- Nodes are configuration sets (starting from  $s_0$ )
- Arcs are  $\text{go\_to}$  relationships

$S' \rightarrow S \$$   
 $S \rightarrow ID$



Wednesday, January 19, 2011

## Building the goto table

- We can just read this off from the CFSM

		Symbol		
		ID	\$	S
State	0	1		2
	1			
	2		3	
	3			

Wednesday, January 19, 2011

## Building the action table

- Given the configuration set  $s$ :
  - We **shift** if the next token matches a terminal after the  $\cdot$  in some configuration  
 $A \rightarrow \alpha \cdot a \beta \in s$  and  $a \in V_t$ , else error
  - We **reduce** production  $P$  if the  $\cdot$  is at the end of a production  
 $B \rightarrow \alpha \cdot \in s$  where production  $P$  is  $B \rightarrow \alpha$
- Extra actions:
  - shift** if goto table transitions between states on a non-terminal
  - accept** if we are about to shift  $\$$

Wednesday, January 19, 2011

## Action table

		Symbol		
		ID	\$	S
State	0	S		S
	1	R2	R2	R2
	2		A	
	3			

Wednesday, January 19, 2011

## Conflicts in action table

- For LR(0) grammars, the action table entries are unique: from each state, can only shift or reduce
- But other grammars may have conflicts
  - Reduce/reduce conflicts: multiple reductions possible from the given configuration
  - Shift/reduce conflicts: we can either shift or reduce from the given configuration

Wednesday, January 19, 2011

## Shift/reduce example

- Consider the following grammar:  
 $S \rightarrow A y$   
 $A \rightarrow \lambda \mid x$
- This leads to the following initial configuration set:  
 $S \rightarrow \cdot A y$   
 $A \rightarrow \cdot x$   
 $A \rightarrow \lambda \cdot$
- Can shift or reduce here

Wednesday, January 19, 2011

## Lookahead

- Can resolve reduce/reduce conflicts and shift/reduce conflicts by employing *lookahead*
  - Looking ahead one (or more) tokens allows us to determine whether to shift or reduce
  - (cf how we resolved ambiguity in LL(1) parsers by looking ahead one token)

Wednesday, January 19, 2011

## Semantic actions

- Recall: in LL parsers, we could integrate the semantic actions with the parser
  - Why? Because the parser was **predictive**
- Why doesn't that work for LR parsers?
  - Don't know which production is matched until parser reduces
- For LR parsers, we put semantic actions at the end of productions
  - May have to rewrite grammar to support all necessary semantic actions

Wednesday, January 19, 2011

# Parsers with lookahead

- Adding lookahead creates an **LR(I) parser**
  - Built using similar techniques as LR(0) parsers, but uses lookahead to distinguish states
  - LR(I) machines can be much larger than LR(0) machines, but resolve many shift/reduce and reduce/reduce conflicts
  - Other types of LR parsers are SLR(I) and LALR(I)
    - Differ in how they resolve ambiguities
    - yacc and bison produce LALR(I) parsers

Wednesday, January 19, 2011

# LR(I) parsing

- Configurations in LR(I) look similar to LR(0), but they are extended to include a lookahead symbol
$$A \rightarrow X_1 \dots X_i \bullet X_{i+1} \dots X_j, l \text{ (where } l \in V_t \cup \lambda \text{)}$$
- If two configurations differ only in their lookahead component, we combine them
$$A \rightarrow X_1 \dots X_i \bullet X_{i+1} \dots X_j, \{l_1 \dots l_m\}$$

Wednesday, January 19, 2011

## Building configuration sets

- To close a configuration
$$B \rightarrow \alpha \bullet A \beta, l$$
- Add all configurations of the form  $A \rightarrow \bullet \gamma, u$  where  $u \in \text{First}(\beta l)$
- Intuition: the parse could apply the production for A, and the lookahead after we apply the production should match the next token that would be produced by B

Wednesday, January 19, 2011

## Example

$\text{closure}(\{S \rightarrow \bullet E \$, \{\lambda\}\}) =$
---

$S \rightarrow E \$$
$E \rightarrow E + T \mid T$
$T \rightarrow ID \mid (E)$

Wednesday, January 19, 2011

## Example

$\text{closure}(\{S \rightarrow \bullet E \$, \{\lambda\}\}) =$
$S \rightarrow \bullet E \$, \{\lambda\}$

$S \rightarrow E \$$
$E \rightarrow E + T \mid T$
$T \rightarrow ID \mid (E)$

Wednesday, January 19, 2011

## Example

$\text{closure}(\{S \rightarrow \bullet E \$, \{\lambda\}\}) =$
$S \rightarrow \bullet E \$, \{\lambda\}$
$E \rightarrow \bullet E + T, \{\$ \}$

$S \rightarrow E \$$
$E \rightarrow E + T \mid T$
$T \rightarrow ID \mid (E)$

Wednesday, January 19, 2011

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	
$T \rightarrow \bullet ID, \{\$ \}$	

Wednesday, January 19, 2011

Wednesday, January 19, 2011

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	
$T \rightarrow \bullet ID, \{\$ \}$	
$T \rightarrow \bullet (E), \{\$ \}$	

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	
$T \rightarrow \bullet ID, \{\$ \}$	
$T \rightarrow \bullet (E), \{\$ \}$	
$E \rightarrow \bullet E + T, \{+\}$	

Wednesday, January 19, 2011

Wednesday, January 19, 2011

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	
$T \rightarrow \bullet ID, \{\$ \}$	
$T \rightarrow \bullet (E), \{\$ \}$	
$E \rightarrow \bullet E + T, \{+\}$	
$E \rightarrow \bullet T, \{+\}$	

## Example

$S \rightarrow E \$$   
 $E \rightarrow E + T \mid T$   
 $T \rightarrow ID \mid (E)$

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =	
$S \rightarrow \bullet E \$, \{\lambda\}$	
$E \rightarrow \bullet E + T, \{\$ \}$	
$E \rightarrow \bullet T, \{\$ \}$	
$T \rightarrow \bullet ID, \{\$ \}$	
$T \rightarrow \bullet (E), \{\$ \}$	
$E \rightarrow \bullet E + T, \{+\}$	
$E \rightarrow \bullet T, \{+\}$	
$T \rightarrow \bullet ID, \{+\}$	

Wednesday, January 19, 2011

Wednesday, January 19, 2011

## Example

S → E \$  
E → E + T | T  
T → ID | (E)

closure( $\{S \rightarrow \bullet E \$, \{\lambda\}\}$ ) =
$S \rightarrow \bullet E \$, \{\lambda\}$
$E \rightarrow \bullet E + T, \{\$ \}$
$E \rightarrow \bullet T, \{\$ \}$
$T \rightarrow \bullet ID, \{\$ \}$
$T \rightarrow \bullet (E), \{\$ \}$
$E \rightarrow \bullet E + T, \{+\}$
$E \rightarrow \bullet T, \{+\}$
$T \rightarrow \bullet ID, \{+\}$
$T \rightarrow \bullet (E), \{+\}$

Wednesday, January 19, 2011

## Building goto and action tables

- The function **goto** (configuration-set, symbol) is analogous to **goto** (configuration-set, symbol) for LR(0)
- Build goto table in the same way as for LR(0)
- Key difference: the action table.  
action[s][x] =
  - reduce** when  $\bullet$  is at end of configuration and  $x \in$  lookahead set of configuration  
 $A \rightarrow \alpha \bullet, \{\dots x \dots\} \in s$
  - shift** when  $\bullet$  is before x  
 $A \rightarrow \beta \bullet x \gamma \in s$

Wednesday, January 19, 2011

## Problems with LR(1) parsers

- LR(1) parsers are very powerful ...
  - But the table size is much larger than LR(0) — as much as a factor of  $|V_t|$  (why?)
  - Example: Algol 60 (a simple language) includes several thousand states!
- Storage efficient representations of tables are an important issue

Wednesday, January 19, 2011

## Solutions to the size problem

- Different parser schemes
  - SLR (simple LR): build an CFSM for a language, then add lookahead wherever necessary (i.e., add lookahead to resolve shift/reduce conflicts)
    - What should the lookahead symbol be?
    - To decide whether to reduce using production  $A \rightarrow \alpha$ , use Follow(A)
  - LALR: merge LR states in certain cases (we won't discuss this)

Wednesday, January 19, 2011