


UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE

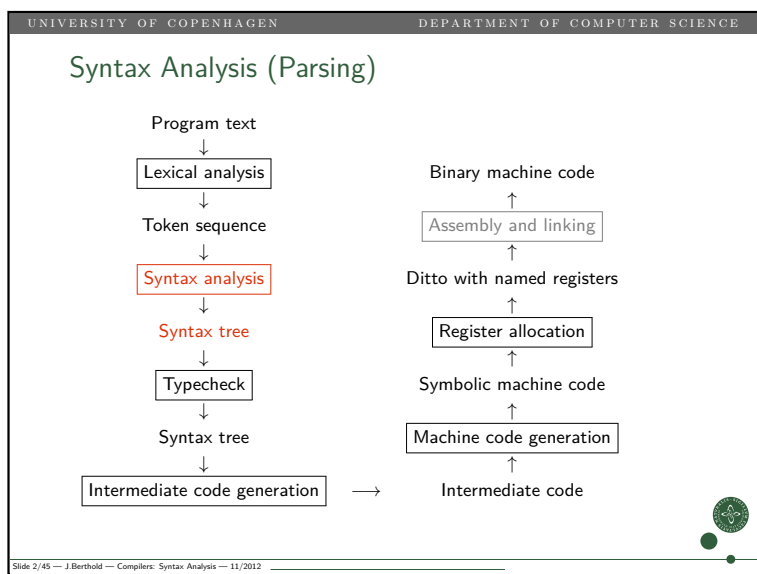
Faculty of Science

# Compilers (Oversættere): Syntax Analysis

Jost Berthold  
berthold@diku.dk  
Department of Computer Science



November 2012  
Slide 1/45



UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE

## Syntax Analysis

syntax analysis covers lecture This.

This analysis covers lecture syntax.

**Syntax error!** (semantic error)

- Words (tokens) need to appear in the right order to form correct sentences (programs) (not necessarily meaningful<sup>1</sup>).
- Syntax analyser, commonly called parser,
- ... analyses token sequence to **build** program **structure**.
- Essential tool and theory used here: Context-free languages.

<sup>1</sup>Traditional example (Noam Chomsky 1957): *Colorless green ideas sleep furiously.*

Slide 3/45 — J.Berthold — Compilers: Syntax Analysis — 11/2012

UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE

## Contents and Goals of this Part

- Context-Free Grammars and Languages
- Top-Down Parsing, LL(1)
  - Recursive Parsing Functions (Recursive-descent)
  - First- and Follow-Sets
  - Look-Ahead Sets and LL(1) Parsing
- Bottom-Up Parsing, SLR
  - Parser Generator Yacc
  - Shift-Reduce Parsing
- Precedence and Associativity

**Goals:**

- Use suitable context-free grammars to describe syntactic structure (especially for programming languages);
- Use parser generators and explain their inner workings;
- Know and use recursive-descent (top-down) parsing;
- Understand concepts and limitations of context-free parsing.

Slide 4/45 — J.Berthold — Compilers: Syntax Analysis — 11/2012

UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE

## Context-Free Grammars

### Definition (Context-Free Grammar)

A context-free grammar is given by

- a set of terminals  $\Sigma$  (the alphabet of the resulting language),
- a set of nonterminals  $N$ ,
- a start symbol  $S \in N$
- a set  $P$  of productions  $X \rightarrow \alpha$  with a single nonterminal  $X \in N$  on the left and a (possibly empty) right-hand side  $\alpha \in (\Sigma \cup N)^*$  of terminals and nonterminals.

$G : S \rightarrow aSB$   
 $S \rightarrow \epsilon$   
 $S \rightarrow B$   
 $B \rightarrow Bb \mid b$   
 (alternative notation)

- Context-free grammars describe (context-free) languages over their terminal alphabet  $\Sigma$ .
- Each nonterminal describes a set of words.
- Nonterminals recursively refer to each other. (cannot do that with regular expressions)

Slide 5/45 — J.Berthold — Compilers: Syntax Analysis — 11/2012

UNIVERSITY OF COPENHAGEN DEPARTMENT OF COMPUTER SCIENCE

## Example, Derivation of Words

$G : S \rightarrow aSB$  (1)  
 $S \rightarrow \epsilon$  (2)  
 $S \rightarrow B$  (3)  
 $B \rightarrow Bb$  (4)  
 $B \rightarrow b$  (5)

$$S = \underbrace{\{a \cdot x \cdot y \mid x \in S, y \in B\}}_{(1)} \cup \underbrace{\{\epsilon\}}_{(2)} \cup \underbrace{B}_{(3)}$$

$$B = \underbrace{\{x \cdot b \mid x \in B\}}_{(4)} \cup \underbrace{\{b\}}_{(5)}$$

- Starting from the start symbol  $S, \dots$
- words of the language can be derived. . .
- by successively replacing nonterminals with right-hand sides.

$$S \xrightarrow{1} aSB \xrightarrow{1} aaSBB \xrightarrow{5} aaSbB \xrightarrow{1} aaaSBbB$$

$$\xrightarrow{2} aaa\_BbB \xrightarrow{4} aaaBbbB \xrightarrow{5} aaaBbbb \xrightarrow{5} aaabbbb$$

Slide 6/45 — J.Berthold — Compilers: Syntax Analysis — 11/2012

### Derivation Relation

#### Definition (Derivation $\Rightarrow$ )

Let  $G = (\Sigma, N, S, P)$  be a grammar.

The derivation relation  $\Rightarrow$  on  $(\Sigma \cup N)^*$  is defined as follows:

- For an  $X \in N$  and a production  $(X \rightarrow \beta) \in P$  of the grammar,  $\alpha_1 X \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2$  for all  $\alpha_1, \alpha_2 \in (\Sigma \cup N)^*$ .
- Describes one derivation step using one of the productions.
- Can indicate used production by a number ( $\xRightarrow{k}$ ).
- Can indicate left-most (or right-most) derivation ( $\xRightarrow{k}_l, \xRightarrow{k}_r$ ).

```
G : S → aSB (1)
    S → ε (2)
    S → B (3)
    B → Bb (4)
    B → b (5)

S  $\xRightarrow{1}$  aSB  $\xRightarrow{2}$  aaSBB  $\xRightarrow{3}$  aa_BB
    $\xRightarrow{4}$  aaBbB  $\xRightarrow{5}$  aabbB  $\xRightarrow{5}$  aabb
```

### Extended Derivation Relation (Transitive Closure)

#### Definition (Transitive Derivation Relation $\Rightarrow^*$ )

Let  $G = (\Sigma, N, S, P)$  be a grammar and  $\Rightarrow$  its derivation relation.

The transitive derivation relation of  $G$  is defined as:

- $\alpha \Rightarrow^* \alpha$  for all  $\alpha \in (\Sigma \cup N)^*$  (derived in 0 steps).
- For  $\alpha, \beta \in (\Sigma \cup N)^*$ ,  $\alpha \Rightarrow^* \beta$  if there exists a  $\gamma \in (\Sigma \cup N)^*$  such that  $\alpha \Rightarrow \gamma$  and  $\gamma \Rightarrow^* \beta$  (derived in at least one step).

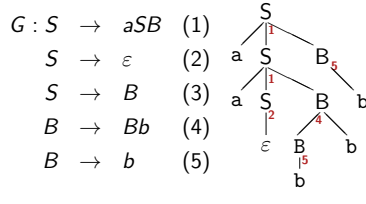
More generally, this is known as the transitive closure of a relation. In our previous examples, we saw  $S \Rightarrow^* aabbbb$  and  $S \Rightarrow^* aabbb$ . That means, both words are in the language of  $G$ .

#### Definition (Language of a Grammar)

Let  $G = (\Sigma, N, S, P)$  be a grammar and  $\Rightarrow$  its derivation relation.

The language of the grammar is  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$ .

### Syntax Tree and Directed Derivation



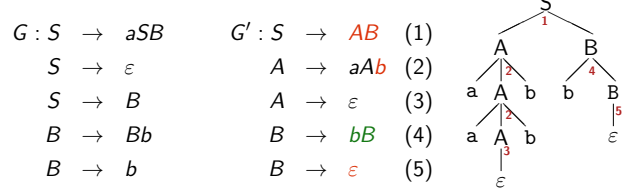
- Syntax trees describe the derivation independent of the direction.
- Left-most derivation: depth-first left-to-right tree traversal.
- $S \xRightarrow{1}_l aSB \xRightarrow{2}_l aaSBB \xRightarrow{3}_l aa\_BB \xRightarrow{4}_l aaBbB \xRightarrow{5}_l aabb$

Nevertheless:  $S \Rightarrow^* aabbb$  can be derived in two ways.

- $S \xRightarrow{1}_l aSB \xRightarrow{2}_l aaSBB \xRightarrow{3}_l aa\_BBB \xRightarrow{4}_l aabBB \xRightarrow{5}_l aabbB \xRightarrow{5}_l aabb$

The grammar  $G$  is said to be ambiguous.

### Avoiding Ambiguity (Changed Grammar)



Modify the grammar to make it non-ambiguous. (describing the same language), give a syntax tree for  $aabbb$ .

- Idea: generate extra bs separately by **new start production**
- **Avoiding left-recursion** (explained later)
- Left-most derivation: (1 2 2 3 4 5)

```
S  $\xRightarrow{1}_l AB \xRightarrow{2}_l aAb \xRightarrow{2}_l aaAb \xRightarrow{3}_l aa\_bB \xRightarrow{4}_l aabbb \xRightarrow{5}_l aabbb$ 
```

### Parsing

Token sequence

Syntax analysis

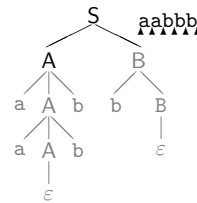
Syntax tree

- Producing a syntax tree from a token sequence.
- Representation of the tree: left-most or right-most derivation

#### Two approaches

- **Top-Down Parsing:** Builds syntax tree from the root. Builds a left-most derivation sequence
- **Bottom-Up Parsing:** Builds syntax tree from the leaves. Builds a reversed right-most derivation sequence
- Both: use stack to keep track of derivation.

### Idea of Top-Down Parsing



- Recursive functions modelling the productions ("recursive-descent")

```
fun parseS () = print "parsing_⊔S: prod_⊔1";
                (* one production S -> A B *)
                parseA (); parseB (); match EOF
```

```
and parseA () =
  (* choose A -> a A b or A -> <epsilon> *)
  if should_use_production_2
  then print "parsing_⊔A: prod_⊔2";
        match #"a"; parseA (); match #"b"
  else print "parsing_⊔A: prod_⊔3";()
```

How can we decide which production to use?

```
and parseB () =
  (* choose B -> b B or B -> <epsilon> *)
  if should_use_production_4
  then print "parsing_⊔B: prod_⊔4";
        match #"b"; parseB ()
  else print "parsing_⊔B: prod_⊔5";()
```

## Top-Down Parsing (LL(1) Parsing)

Token sequence

Syntax analysis

Syntax tree

- Producing a left-most derivation from a token sequence.
- Uses a stack (maybe the function call stack) to keep track of derivation.

- Called predictive parsing: needs to "guess" used productions.
- Information to choose the right production (look-ahead):
  - For each right-hand side: What input token can come first?
  - Special attention to empty right-hand sides. What can follow?
- A production  $A \rightarrow \alpha$  is chosen
  - if look-ahead  $c$  and  $\alpha \Rightarrow^* c\beta$  (can start with  $c$ ).
  - or if look-ahead  $c$ ,  $\alpha \Rightarrow^* \epsilon$ , and  $c$  can follow  $A$ .



## FIRST Sets and Property NULLABLE

### Definition (FIRST set and NULLABLE)

Let  $G = (\Sigma, N, S, P)$  a grammar and  $\Rightarrow$  its derivation relation. For all sequences of grammar symbols  $\alpha \in (\Sigma \cup N)^*$ , define

- $FIRST(\alpha) = \{c \in \Sigma \mid \exists \beta \in (\Sigma \cup N)^* : \alpha \Rightarrow^* c\beta\}$   
(all terminals at the start of what can be derived from  $\alpha$ )
- $NULLABLE(\alpha) = \begin{cases} true & , \text{ if } \alpha \Rightarrow^* \epsilon \\ false & , \text{ otherwise} \end{cases}$

Computing NULLABLE and FIRST for right-hand sides:

- Set equations recursively use results for nonterminals.
- Smallest solution found by computing a smallest fixed-point.
- Solved simultaneously for all right-hand sides of the productions.



## Computing NULLABLE by Set Equations

$$\begin{aligned} NULLABLE(\epsilon) &= true \\ NULLABLE(a) &= false \text{ for } a \in \Sigma \\ NULLABLE(\alpha\beta) &= NULLABLE(\alpha) \wedge NULLABLE(\beta) \text{ for } \alpha, \beta \in (\Sigma \cup N)^* \\ NULLABLE(A) &= NULLABLE(\alpha_1) \vee \dots \vee NULLABLE(\alpha_n), \\ &\text{using all productions for } A, A \rightarrow \alpha_i \text{ (} i \in \{1..n\}\text{)} \end{aligned}$$

- Equations for nonterminals of the grammar:

$$\begin{aligned} G' : S &\rightarrow AB & NULLABLE(S) &= NULLABLE(AB) = true \\ A &\rightarrow aAb \mid \epsilon & NULLABLE(A) &= NULLABLE(aAb) \vee NULLABLE(\epsilon) = true \\ B &\rightarrow bB \mid \epsilon & NULLABLE(B) &= NULLABLE(bB) \vee NULLABLE(\epsilon) = true \end{aligned}$$

- Equations for the right-hand side

$$\begin{aligned} NULLABLE(AB) &= NULLABLE(A) \wedge NULLABLE(B) \\ NULLABLE(aAb) &= NULLABLE(a) \wedge NULLABLE(A) \wedge NULLABLE(b) = false \\ NULLABLE(bB) &= NULLABLE(b) \wedge NULLABLE(B) = false \\ NULLABLE(\epsilon) &= true \end{aligned}$$

Compute smallest solution of system, starting by false for all.



## Computing FIRST by Set Equations

$$\begin{aligned} FIRST(\epsilon) &= \emptyset \\ FIRST(a) &= a \text{ for } a \in \Sigma \\ FIRST(\alpha\beta) &= \begin{cases} FIRST(\alpha) \cup FIRST(\beta) & , \text{ if } NULLABLE(\alpha) \\ FIRST(\alpha) & , \text{ otherwise} \end{cases} \\ FIRST(A) &= FIRST(\alpha_1) \cup \dots \cup FIRST(\alpha_n), \\ &\text{using all productions for } A, A \rightarrow \alpha_i \text{ (} i \in \{1..n\}\text{)} \end{aligned}$$

- Equations for nonterminals of the grammar:

$$\begin{aligned} G' : S &\rightarrow AB & FIRST(S) &= FIRST(AB) = FIRST(A) \cup FIRST(B) = \{a, b\} \\ A &\rightarrow aAb \mid \epsilon & FIRST(A) &= FIRST(aAb) \cup FIRST(\epsilon) = \{a\} \\ B &\rightarrow bB \mid \epsilon & FIRST(B) &= FIRST(bB) \cup FIRST(\epsilon) = \{b\} \end{aligned}$$

- Equations for the right-hand side

$$\begin{aligned} FIRST(aAb) &= FIRST(a) = \{a\} \\ FIRST(bB) &= FIRST(b) = \{b\} \\ FIRST(\epsilon) &= \emptyset \end{aligned}$$

Compute smallest solution of system, starting by  $\emptyset$  for all sets.



## FOLLOW Sets for Nonterminals

FIRST Sets are often not enough.

In production  $X \rightarrow \alpha$ , if  $NULLABLE(\alpha)$ , we need to know what can follow  $X$  (FIRST set of  $\alpha$  cannot provide this information).

### Definition (FOLLOW Set of a Nonterminal)

Let  $G = (\Sigma, N, S, P)$  a grammar and  $\Rightarrow$  its derivation relation. For each nonterminal  $X \in N$ , define

- $FOLLOW(X) = \{c \in \Sigma \mid \exists \alpha, \beta \in (\Sigma \cup N)^* : S \Rightarrow^* \alpha X c \beta\}$   
(all input tokens that follow  $X$  in sequences derivable from  $S$ )

To recognise the end of the input

- add a new character  $\$$  to the alphabet
- add start production  $S' \rightarrow S\$$  to the grammar.

Thereby, we always have  $\$ \in FOLLOW(S)$ .



## Set Equations for FOLLOW Sets

FOLLOW sets solve a collection of set constraints.

Constraints derived from right-hand sides of grammar productions

For  $X \in N$ , consider all productions  $Y \rightarrow \alpha X \beta$  where  $X$  occurs on the right.

- $FIRST(\beta) \subseteq FOLLOW(X)$
- If  $NULLABLE(\beta)$  or  $\beta = \epsilon$ :  $FOLLOW(Y) \subseteq FOLLOW(X)$

If  $X$  occurs several times, each occurrence contributes separate equations.

$$\begin{aligned} S' &\rightarrow S\$ & \dots & FIRST(\$) = \{\$ \} \subseteq FOLLOW(S) \\ S &\rightarrow AB & \dots & FIRST(B) = \{b\} \subseteq FOLLOW(A) \\ & & & FOLLOW(S) \subseteq FOLLOW(A) \text{ (B nullable)} \\ & & & FOLLOW(S) \subseteq FOLLOW(B) \\ A &\rightarrow aAb & \dots & FIRST(b) = \{b\} \subseteq FOLLOW(A) \\ B &\rightarrow bB & \dots & FOLLOW(B) \subseteq FOLLOW(B) \end{aligned}$$

Example:

$A \rightarrow \epsilon$  and  $B \rightarrow \epsilon$  do not contribute.

Solve iteratively, starting by  $\emptyset$  for all nonterminals.

$$\begin{aligned} FOLLOW(S) &= FOLLOW(B) = \{\$ \} \\ FOLLOW(A) &= \{\$, b\} \end{aligned}$$



### Putting it Together: Look-ahead Sets and LL(1)

After computing NULLABLE and FIRST for all right-hand sides and FOLLOW for all nonterminals, a parser can be constructed.

#### Definition (Look-ahead Sets of a Grammar)

For every production  $X \rightarrow \alpha$  of a context-free grammar G, we define the Look-ahead set of the production as:

$$la(X \rightarrow \alpha) = \begin{cases} \text{FIRST}(\alpha) \cup \text{FOLLOW}(X) & , \text{ if } \text{NULLABLE}(\alpha) \\ \text{FIRST}(\alpha) & , \text{ otherwise} \end{cases}$$

#### LL(1) Grammars

If for each nonterminal  $X \in N$  in grammar G, all productions of X have disjoint look-ahead sets, the grammar G is LL(1) (left-to-right, left-most, look-ahead 1).

For an LL(1) grammar, a parser can be constructed which constructs a left-most derivation for valid input with one token look-ahead (predicting the next production from look-ahead).

### Recursive Descent with Look-Ahead

The grammar in our example is LL(1):

$$\begin{aligned} G' : S &\rightarrow AB & la(S \rightarrow AB) &= \text{FIRST}(AB) \cup \text{FOLLOW}(S) = \{a, b, \$\} \\ A &\rightarrow aAb & la(A \rightarrow aAb) &= \text{FIRST}(aAb) = \{a\} \\ A &\rightarrow \epsilon & la(A \rightarrow \epsilon) &= \text{FIRST}(\epsilon) \cup \text{FOLLOW}(A) = \{b, \$\} \\ B &\rightarrow bB & la(B \rightarrow bB) &= \text{FIRST}(bB) = \{b\} \\ B &\rightarrow \epsilon & la(B \rightarrow \epsilon) &= \text{FIRST}(\epsilon) \cup \text{FOLLOW}(B) = \{\$ \} \end{aligned}$$

```

fun parseS ()
= if next = #"a" orelse next = #"b" orelse next = EOF
  then parseA (); parseB (); match EOF else error

and parseA () (* choose by look-ahead *)
= if next = #"a" then match #"a"; parseA (); match #"b"
  else if next = #"b" orelse next = EOF then ()
  else error

and parseB () = if next = #"b" then match #"b"; parseB ()
  else if next = EOF then ()
  else error
    
```

### Table-Driven LL(1) Parsing

- Stack, contains unprocessed part of production, initially S\$.
- Parser Table: action to take, depends on stack and next input
- Actions (*pop* consumes input, derivation only reads it)
- Pop:** remove terminal from stack (on matching input).
- Derive:** pop nonterminal from stack, push right-hand side (in table).
- Accept input when stack empty at end of input.

Example run (input aabbb):

Stack:	Look-ahead/Input:			Input	Stack	Action	Output
	a	b	\$				
S	AB, 1	AB, 1	AB, 1	aabbbb\$	S\$	derive	$\epsilon$
A	aAb, 2	$\epsilon$ , 3	$\epsilon$ , 3	aabbbb\$	AB\$	derive	1
B	error	bB, 4	$\epsilon$ , 5	aabbbb\$	aAbB\$	pop	12
a	pop	error	error	abbb\$	AbB\$	derive	12
b	error	pop	error	abbb\$	aAbbB\$	pop	122
\$	error	error	accept	bbb\$	AbbB\$	derive	122
				bbb\$	bbB\$	pop	1223
				bb\$	bB\$	pop	1223
				b\$	B\$	derive	1223
				b\$	bB\$	pop	12234
				\$	B\$	derive	12234
				\$	\$	accept	122345

### Eliminating Left-Recursion and Left-Factorisation

Problems that often occur when constructing LL(1) parsers:

- Identical prefixes:** Productions  $X \rightarrow \alpha\beta \mid \alpha\gamma$ .  
Requires look-ahead longer than the common prefix  $\alpha$ .  
Solution: Left-Factorisation, introducing new productions  $X \rightarrow \alpha Y$  and  $Y \rightarrow \beta \mid \gamma$ .
- Left-Recursion:** a nonterminal reproducing itself on the left.  
Direct: production  $X \rightarrow X\alpha \mid \beta$ , or indirect:  $X \Rightarrow^* X\alpha$ .  
Cannot be analysed with finite look-ahead!  
 $X \rightarrow X\alpha \mid \beta$ , thus  $\text{FIRST}(X) \subset \text{FIRST}(X\alpha) \cup \text{FIRST}(\beta)$   
Solution: new (nullable) nonterminals and swapped recursion.  
 $X \rightarrow \beta X'$  and  $X' \rightarrow \alpha X' \mid \epsilon$   
Also works in case of multiple left-recursive productions.  
For indirect recursion: first transform into direct recursion.

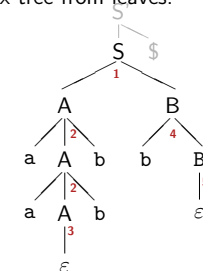
### Contents

- Context-Free Grammars and Languages
- Top-Down Parsing, LL(1)
  - Recursive Parsing Functions (Recursive-descent)
  - First- and Follow-Sets
  - Look-Ahead Sets and LL(1) Parsing
- Bottom-Up Parsing, SLR
  - Parser Generator Yacc
  - Shift-Reduce Parsing
- Precedence and Associativity

### Bottom-Up Parsing

LL(1) Parser works top-down. Needs to guess used productions.  
Bottom-Up approach: build syntax tree from leaves.

$$\begin{aligned} G'' : S' &\rightarrow S\$ \quad (0) \\ S &\rightarrow AB \quad (1) \\ A &\rightarrow aAb \quad (2) \\ A &\rightarrow \epsilon \quad (3) \\ B &\rightarrow bB \quad (4) \\ B &\rightarrow \epsilon \quad (5) \end{aligned}$$



$$S' \xrightarrow{0}_r S\$ \xrightarrow{1}_r AB \xrightarrow{4}_r AbB \xrightarrow{5}_r Ab_\epsilon \xrightarrow{2}_r aAbb \xrightarrow{2}_r aaAbbb \xrightarrow{3}_r aa\_bbb$$

Right-most derivation: 1 4 5 2 2 3

### Bottom-Up Parsing: Idea for a Machine

Stack	Input	Action
$\epsilon$	aabbb\$	shift
a	abbb\$	shift
aa	bbb\$	reduce 3
aaA	bbb\$	shift
aaAb	bb\$	reduce 2
aA	bb\$	shift
aAb	b\$	reduce 2
A	b\$	shift
Ab	\$	reduce 5
AbB	\$	reduce 4
AB	\$	reduce 1
S	\$	accept

$G'' : S' \rightarrow S\$ (0)$   
 $S \rightarrow AB (1)$   
 $A \rightarrow aAb (2)$   
 $A \rightarrow \epsilon (3)$   
 $B \rightarrow bB (4)$   
 $B \rightarrow \epsilon (5)$

Questions:

- When to accept (solved: separate start production)
- When to shift, when to reduce? Especially  $R \rightarrow \epsilon$ .

### mosmlyac: Yet Another Compiler Compiler in MosML

- Generates bottom-up parser from a grammar specification
- Grammar specification also includes token datatype declaration and other declarations.

### Demo mosmlyac

Tradition: **Lex** and **Yacc** (GNU: flex and bison)

- Parser generators usually use LALR(1) Parsing<sup>2</sup>.
- We use **SLR parsing** instead:  
Simple Left-to-right Right-most analysis with look-ahead 1.

<sup>2</sup>More information about LALR(1) and LR(1) parsing can be found in the Red-Dragon book.

### Constructing an SLR-Parser: Items

Each production in the grammar leads to a number of items:

#### Shift Items and Reduce Items of a Production

Let  $X \rightarrow \alpha$  be a production in a grammar.

The production implies:

- **Shift items:**  $[X \rightarrow \alpha_1 \bullet \alpha_2]$  for every decomposition  $\alpha = \alpha_1 \alpha_2$  (including  $\alpha_1 = \epsilon$  and  $\alpha_2 = \epsilon$ );
- One **reduce item:**  $[X \rightarrow \alpha \bullet]$  per production.

Items give information about the next action:

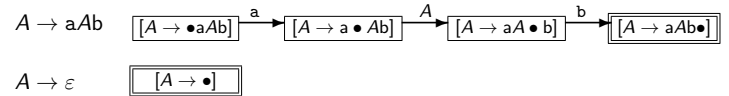
- Either to shift an item to the stack and read input
- or to reduce the top of stack (a production's right-hand side).
- Stack of the parser will contain items, not grammar symbols.
- Therefore, no need to read into the stack for reductions.

### Constructing an SLR Parser: Production DFAs

Each production  $X \rightarrow \alpha$  suggests a DFA with items as states, and doing the following transitions:

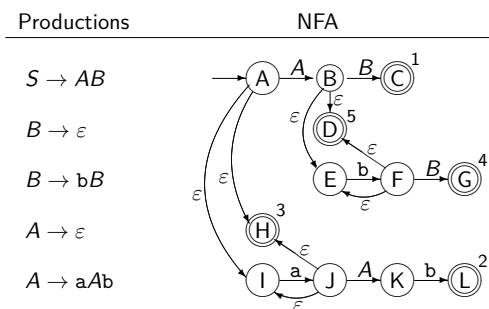
- From  $[X \rightarrow \alpha \bullet a\beta]$  to  $[X \rightarrow \alpha a \bullet \beta]$  for input tokens a. These will be **Shift** action that read input later.
- From  $[X \rightarrow \alpha \bullet Y\beta]$  to  $[X \rightarrow \alpha Y \bullet \beta]$  for nonterminals Y. These will be **Go** actions later, without consuming input.

All items are states, start state is the first item  $[X \rightarrow \bullet \alpha]$ .



While traversing the DFA: items pushed on the stack.  
 When reaching a reduce item: use stack to back-track (later).

### SLR Parser Construction: Example

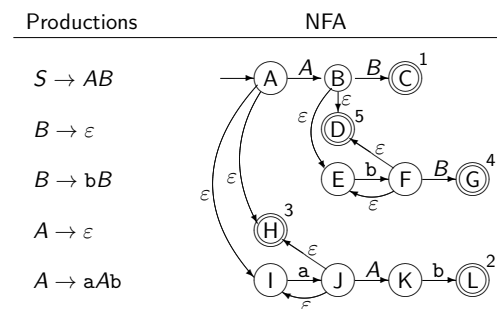


Extra  $\epsilon$ -transitions connect the DFAs for all productions:

- From  $[X \rightarrow \alpha \bullet Y\beta]$  to  $[Y \rightarrow \bullet \gamma]$  for all productions  $Y \rightarrow \gamma$

When in front of a nonterminal Y in a production DFA: try to run the DFA for one of the right-hand sides of Y productions.

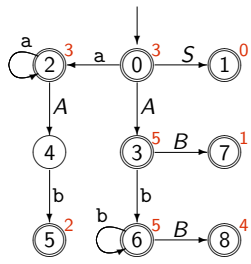
### SLR Parser Construction: Example(2)



Next step: Subset construction of a combined DFA.

**Blackboard...**

### SLR Parsing: Internal DFA and Stack



- Transitions: Shift actions (terminals) and Go actions (nonterminals).
- Final DFA states: contain reduce items. Reduce actions need to be added to the transition table.
- Reduce action: remove items from stack corresponding to right-hand side, then do a Go action with the left-hand side.

- SLR Parse Table: actions indexed by symbols and DFA states
- Shift  $n$  Terminal transition: push state  $n$  on stack, consume input
- Go  $n$  Nonterminal transition: push state  $n$  on stack, (no input read)
- Reduce  $p$  Reduce with production  $p$
- Accept Parsing has succeeded (reduce with production 0).

### SLR Parser Construction: Conflicts

- After constructing a DFA: shift and go actions.
- Next: add reduce actions for states containing reduce items

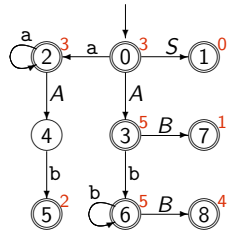
#### SLR Parser Conflicts

Subset construction of the DFA might join conflicting items in one DFA state. We call these conflicts

- **Shift-Reduce conflict**, if a DFA state contains both shift and reduce items. Typically, productions to  $\epsilon$  generate these conflicts.
- **Reduce-Reduce conflict**, if a DFA state contains reduce items for two different productions.

In SLR parsing: FOLLOW sets of nonterminals are compared to the look-ahead to resolve conflicts.

### SLR Parser Construction: The Parser Table



• Parser Table:

	a	b	\$	S	A	B
0	2	red.3	red.3	Go 1	Go 3	
1			acc.			
2	2	red.3	red.3		Go 4	
3		6	red.5			Go 7
4		5				
5		red.2	red.2			
6		6	red.5			Go 8
7			red.1			
8			red.4			

- $G'' : S' \rightarrow SS$  (0)  
 $S \rightarrow AB$  (1)  
 $A \rightarrow aAb$  (2)  
 $A \rightarrow \epsilon$  (3)  
 $B \rightarrow bB$  (4)  
 $B \rightarrow \epsilon$  (5)

- FOLLOW Sets of Nonterminals:  
 FOLLOW( $S$ ) =  $\{\$ \}$   
 FOLLOW( $A$ ) =  $\{b, \$ \}$   
 FOLLOW( $B$ ) =  $\{\$ \}$

### Table-Driven SLR Parsing

- Stack contains DFA states, initially start state 0.
- SLR Parse Table: actions and transitions
- Shift**: do a transition consuming input, push new state on stack
- Reduce**: pop length of right-hand-side from stack, then go to a new state with left-hand side non-terminal, push new state on stack
- Accept input when accept state reached at end of input.

	a	b	\$	S	A	B	Example run (aabb):
0	2	red.3	red.3	Go 1	Go 3		Stack: 0, Input: aabbb\$, Action: shift
1			acc.				Stack: 02, Input: abbb\$, Action: shift
2	2	red.3	red.3		Go 4		Stack: 022, Input: bbb\$, Action: reduce 3
3		6	red.5			Go 7	Stack: 0224, Input: bbb\$, Action: shift
4		5					Stack: 02245, Input: bb\$, Action: reduce 2
5		red.2	red.2				Stack: 024, Input: bb\$, Action: shift
6		6	red.5			Go 8	Stack: 0245, Input: b\$, Action: reduce 2
7			red.1				Stack: 03, Input: b\$, Action: shift
8			red.4				Stack: 036, Input: \$, Action: reduce 5
							Stack: 0368, Input: \$, Action: reduce 4
							Stack: 037, Input: \$, Action: reduce 1
							Stack: 01, Input: \$, Action: accept

### Contents

- 1 Context-Free Grammars and Languages
- 2 Top-Down Parsing, LL(1)
  - Recursive Parsing Functions (Recursive-descent)
  - First- and Follow-Sets
  - Look-Ahead Sets and LL(1) Parsing
- 3 Bottom-Up Parsing, SLR
  - Parser Generator Yacc
  - Shift-Reduce Parsing
- 4 Precedence and Associativity

### Ambiguity, Precedence and Associativity

Arithmetic Expressions:

- $E \rightarrow E + E \mid E - E$
- $E \rightarrow E * E \mid E / E$
- $E \rightarrow a \mid (E)$

- In many cases, grammars are rewritten to remove ambiguity.
- Sometimes, ambiguity is resolved by changes in the parser.

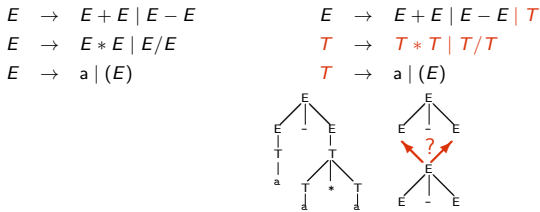
- In both cases: Precedence and associativity guide decisions.

Problems with this grammar:

- 1 Ambiguous derivation of  $a - a * a$ .  
Want precedence of  $*$  over  $+$ ,  $a + (a * a)$ .
- 2 Ambiguous derivation of  $a - a - a$ .  
Want a left-associative interpretation,  $(a - a) - a$ .

### Operator Precedence in the Grammar

- Introduce precedence levels to get operator priorities
- New Grammar: own nonterminal for each level
- Here: 2 levels, mathematical interpretation:  
 $a - a \cdot a = a - (a \cdot a)$  Precedence of \* and / over + and -.  
 More precedence levels could be added (exponentiation).



### About Operator Associativity

#### Definition (Operator Associativity)

A binary operator  $\oplus$  is called

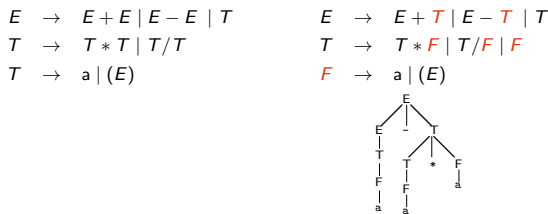
- **left-associative**, if the expression  $a \oplus b \oplus c$  should be evaluated from left to right, as  $(a \oplus b) \oplus c$ .
- **right-associative**, if the expression  $a \oplus b \oplus c$  should be evaluated from right to left, as  $a \oplus (b \oplus c)$ .
- **non-associative**, if expressions  $a \oplus b \oplus c$  are disallowed, (and **associative**, if both directions lead to the same result).

Examples:

- Arithmetic operators like - and /: left-associative.
- List constructors in functional languages: right-associative.
- Function arrows in types: right-associative.
- 'less-than' (<) in C: **left-associative**  
`if (3 < 2 < 1) { fprintf(stdout, "Awesome!\n"); }`

### Establishing the Intended Associativity

- limit recursion to the intended side
- When operators are indeed associative, use same associativity as comparable operators.
- Cannot mix left- and right-associative operators at same precedence level.



### Precedence and Associativity in SLR Parse Tables

Precedence and ambiguity usually materialise as **shift-reduce conflicts** in SLR parsers.



Instead of rewriting the grammar, resolve conflicts by targeted changes to parser table.

- if operator symbol with higher precedence follows: **Shift**
- if operator should be right-associative: **Shift**
- if symbol of lower precedence or left-associative: **Reduce**

### Example: Resolving Precedence and Ambiguity

Regular expressions:

- $R \rightarrow R' \mid R$
- $R \rightarrow RR$
- $R \rightarrow R^*$
- $R \rightarrow \text{char} \mid (R)$

- 1 **Precedence: star, sequence, alternative.**  
 $a \mid b \mid a^* \text{ is } a \mid (b \mid (a^*))$ .
- 2 **Left-associative derivations:**  
 $\alpha \mid \beta \mid \gamma \text{ is } (\alpha \mid \beta) \mid \gamma$ .

New grammar:

- $R \rightarrow R' \mid R2 \mid R2$
- $R2 \rightarrow R2R3 \mid R3$
- $R3 \rightarrow R4^* \mid R4$
- $R4 \rightarrow \text{char} \mid (R)$

Precedence/Associativity declarations:

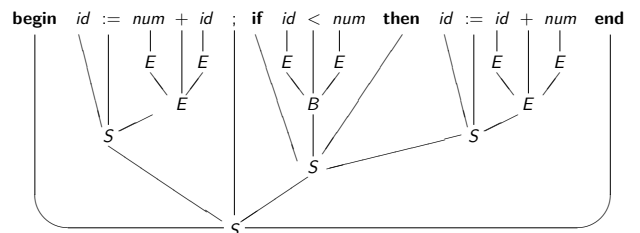
```

mosmlyac file
%token BAR STAR LPAREN RPAREN ...
%left BAR /* lowest precedence */
%nonassoc CHAR LPAREN
%left seq /* pseudo-token for sequence */
%nonassoc STAR /* highest precedence */
...
R : R BAR R { ... }
  R R %prec seq { ... }
  R STAR { ... }
  CHAR { ... }
  LPAREN R RPAREN { ... }
  ...
    
```

Full example: Mosmlyac Demo (regular expressions)

### One word about the Syntax Trees

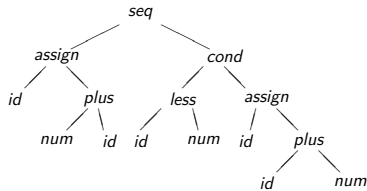
- Concrete Syntax contains many extra tokens for practical reasons:
  - Parentheses, braces, ... for grouping,
  - Semicolons, commas, ... to separate statements or arguments.
  - begin, end ... (also a kind of parentheses).
- Following stage works on abstract syntax tree without those



## One word about the Syntax Trees

- Concrete Syntax contains many extra tokens for practical reasons:
  - Parentheses, braces, ... for grouping,
  - Semicolons, commas, ... to separate statements or arguments.
  - `begin`, `end` ... (also a kind of parentheses).
- Following stage works on abstract syntax tree without those

```
begin id := num + id ; if id < num then id := id + num end
```



## More about Context-Free Languages

- Context-Free languages are commonly processed using a stack machine (Push-Down Automaton, PDA)
- Can count one thing at a time, or remember input.
  - $\{a^n b^n \mid n \in \mathbb{N}\}$  context-free.
  - $\{a^n b^n c^n \mid n \in \mathbb{N}\}$  not context-free!
- Palindromes over  $\Sigma$ : context-free language.
  - However: non-deterministic (need to guess the middle).
  - Non-deterministic stack machines are more powerful than deterministic ones (unlike NFAs and DFAs)!
- Context-free languages are closed under union:
  - $L_1, L_2$  context-free  $\leadsto L_1 \cup L_2$  context-free.
- ... but not closed under intersection (famous counter examples above) and complement (by de Morgan's laws).



## Summary

### Context-free grammars and languages

- Writing and rewriting grammars can be tricky! :-)

### Top-down parsing (recursive-descent)

- FIRST- and FOLLOW-sets;
- Look-ahead sets for decisions in recursive-descent parser.

### Bottom-up parsing (shift-reduce parsing, SLR parsing)

- Items, grammar-implied NFA and subset construction;
- Reduce actions in transition table, stack of visited states.

### Precedence and associativity

- Solved in the grammar or by manipulation of the SLR parser.

