

Programming Language Descriptions and Implementations

BNF and context-free grammar

- Backus-Naur Form (1959)
 - way of specifying programming languages using formal grammars and production rules with a particular form of notation.
 - Invented by John Backus to describe Algol 58
 - BNF is nearly equivalent to context-free grammars

A Grammar for a small language

$\langle \text{program} \rangle \rightarrow \text{begin } \langle \text{stmt_list} \rangle \text{ end}$

$\langle \text{stmt_list} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmt_list} \rangle$

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expression} \rangle$

$\langle \text{var} \rangle \rightarrow A \mid B \mid C$

$\langle \text{expression} \rangle \rightarrow \langle \text{integer} \rangle \mid \langle \text{var} \rangle + \langle \text{integer} \rangle \mid \langle \text{var} \rangle - \langle \text{integer} \rangle$

$\langle \text{integer} \rangle \rightarrow [-] \langle \text{unsigned_integer} \rangle$

$\langle \text{unsigned_integer} \rangle \rightarrow \langle \text{any_digit} \rangle \mid \langle \text{non_zero_digit} \rangle \{ \langle \text{any_digit} \rangle \}$

$\langle \text{any_digit} \rangle \rightarrow 0 \mid \langle \text{non_zero_digit} \rangle$

$\langle \text{non_zero_digit} \rangle \rightarrow 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

| : alternative
[]: optional
{ }: zero or more times

? What does this grammar describe?

Derivation Example

- $\langle \text{Program} \rangle \Rightarrow \text{begin } \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } \langle \text{stmt} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } \langle \text{var} \rangle = \langle \text{expression} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = \langle \text{expression} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = \langle \text{integer} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = \langle \text{unsigned_integer} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = \langle \text{non_zero_digit} \rangle \langle \text{any_digit} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 8 \langle \text{any_digit} \rangle; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; \langle \text{stmt_list} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; \langle \text{stmt} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; \langle \text{var} \rangle = \langle \text{expression} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; B = \langle \text{expression} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; B = \langle \text{integer} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; B = \langle \text{unsigned_integer} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; B = \langle \text{non_zero_digit} \rangle \text{ end}$
 - $\Rightarrow \text{begin } A = 80; B = 2 \text{ end}$

Parse Tree

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C \mid D$

$\langle \text{expr} \rangle \rightarrow \langle \text{id} \rangle + \langle \text{expr} \rangle \mid \langle \text{id} \rangle * \langle \text{expr} \rangle \mid (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

$A = B * (A + C)$

$\langle \text{assign} \rangle$

$\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\Rightarrow A = \langle \text{expr} \rangle$

$\Rightarrow A = \langle \text{id} \rangle * \langle \text{expr} \rangle$

$\Rightarrow A = B * \langle \text{expr} \rangle$

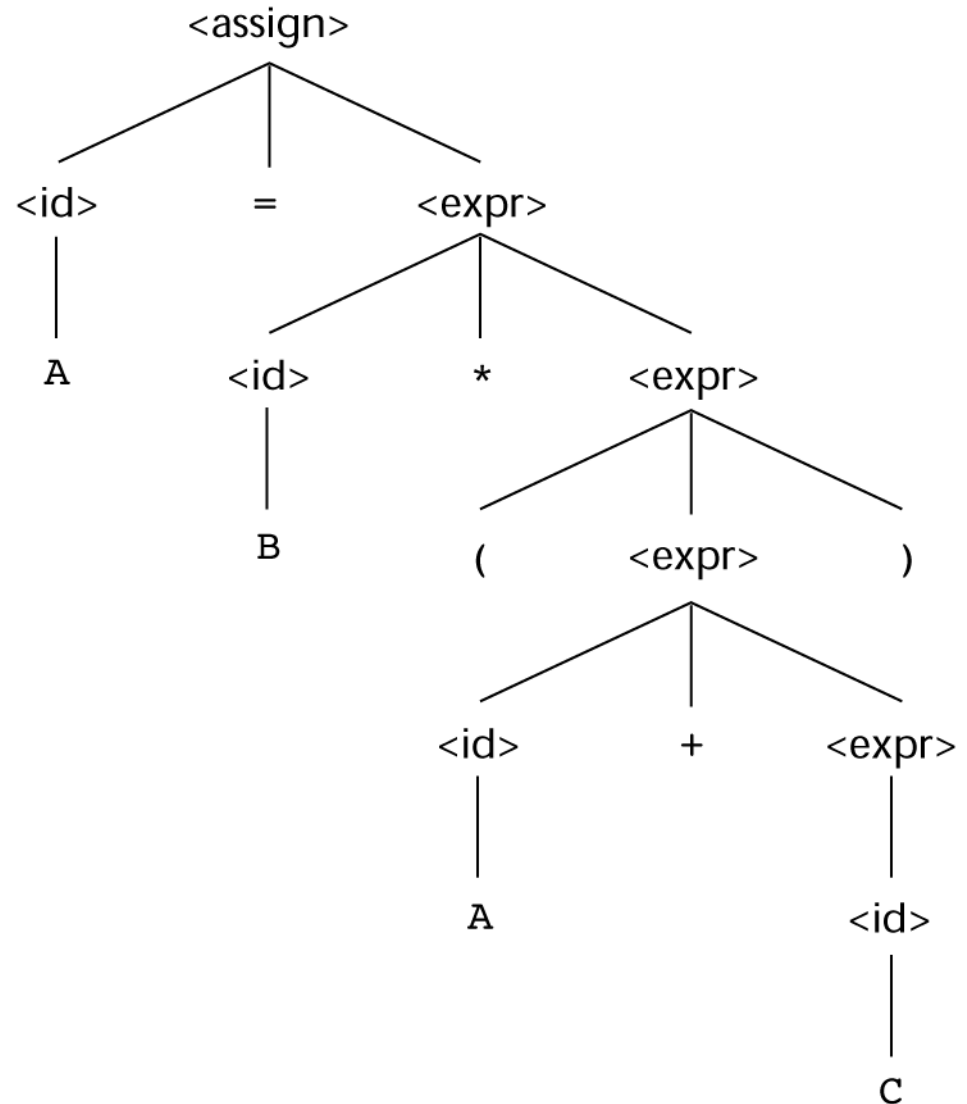
$\Rightarrow A = B * (\langle \text{expr} \rangle)$

$\Rightarrow A = B * (\langle \text{id} \rangle + \langle \text{expr} \rangle)$

$\Rightarrow A = B * (A + \langle \text{expr} \rangle)$

$\Rightarrow A = B * (A + \langle \text{id} \rangle)$

$\Rightarrow A = B * (A + C)$



Introduction

Two classes: static and dynamic semantics

- **Static semantics** defines restrictions on the structure of valid texts that are hard or impossible to express in standard syntactic formalisms
 - Typed variables can only accept values of that type
 - A variable must be declared before it can be used
- **Dynamic semantics** express the meaning of the expressions, statements, and program.
 - After statement `int x = 44 - y; x == 42` is true

Attribute Grammars (Example)

- ❖ Syntax rule: $\langle \text{assign} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
Semantic rule: $\langle \text{expr} \rangle.\text{expected_type} \leftarrow \langle \text{var} \rangle.\text{actual_type}$
- ❖ Syntax rule: $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle[2] + \langle \text{var} \rangle[3]$
Semantic rule:
$$\langle \text{expr} \rangle.\text{actual_type} \leftarrow \text{if (} \langle \text{var} \rangle[2].\text{actual_type} = \text{int) and}$$
$$\langle \text{var} \rangle[3].\text{actual_type} = \text{int)}$$
$$\text{then int}$$
$$\text{else real}$$
$$\text{end if}$$

Predicate: $\langle \text{expr} \rangle.\text{actual_type} = \langle \text{expr} \rangle.\text{expected_type}$
- ❖ Syntax rule: $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle$
Semantic rule:
$$\langle \text{expr} \rangle.\text{actual_type} \leftarrow \langle \text{var} \rangle.\text{actual_type}$$

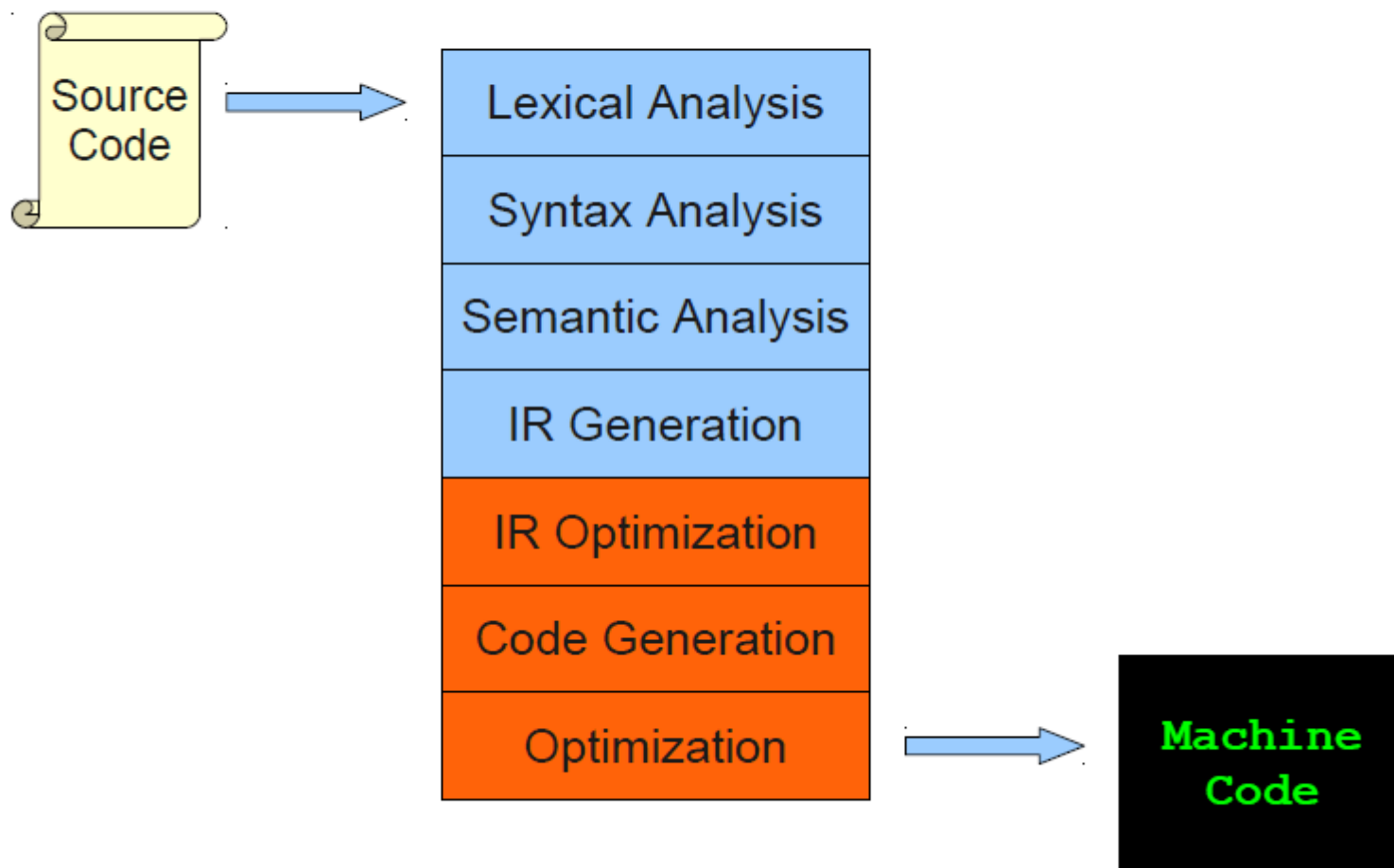
Predicate: $\langle \text{expr} \rangle.\text{actual_type} = \langle \text{expr} \rangle.\text{expected_type}$
- ❖ Syntax rule: $\langle \text{var} \rangle \rightarrow A \mid B \mid C$
Semantic rule:
$$\langle \text{var} \rangle.\text{actual_type} \leftarrow \text{lookup}(\langle \text{var} \rangle.\text{string})$$

Blue: Basic BNF
Red: Semantic Func.
Green: Predicate Func.

Describing (Dynamic) Semantics

- There is no single widely acceptable notation or formalism for describing dynamic semantics
- Three formal methods:
 - ❖ Operational Semantics
 - ❖ Axiomatic Semantics
 - ❖ Denotational Semantics

The Structure of a Modern Compiler



IR: Intermediate Representation

Introduction

For all practical purposes, we can think of the source code as a very long string of characters from some alphabet.

Lexical analyzers collect characters into logical groupings (lexemes) and assign internal codes (tokens) to the grouping.

As an analogy, we can think of the lexical analyzer as a machine that chops up source code into the individual “words”.

What lexical analyzers do

- The “front-end” for the parser
- A **pattern matcher** for character strings
- Identifies substrings of the source program that belong together => lexemes

– Example: sum = B -5;

<u>Lexeme</u>	<u>Token</u>
sum	ID (identifier)
=	ASSIGN_OP
B	ID
-	SUBTRACT_OP
5	INT_LIT (integer literal)
;	SEMICOLON

State Transition Diagram: Example

CFG

$\langle id \rangle \rightarrow \langle letter \rangle \{ \langle letter \rangle | \langle int \rangle \}$

$\langle int \rangle \rightarrow \langle digit \rangle \{ \langle digit \rangle \}$

$\langle letter \rangle \rightarrow A | B | C | \dots | Z | a | b | \dots | z$

$\langle digit \rangle \rightarrow 0 | 1 | 2 | \dots | 9$

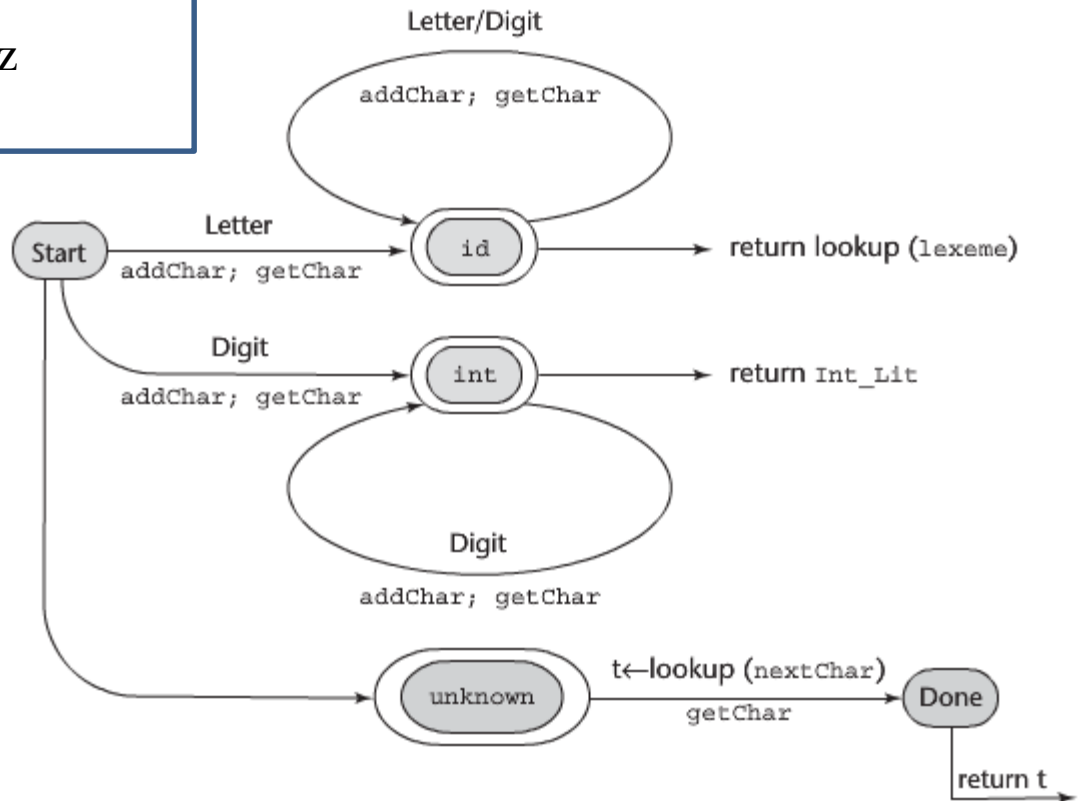
Regular grammar

$Id \rightarrow letter(letter|digit)^*$

$int \rightarrow digit\ digit^*$

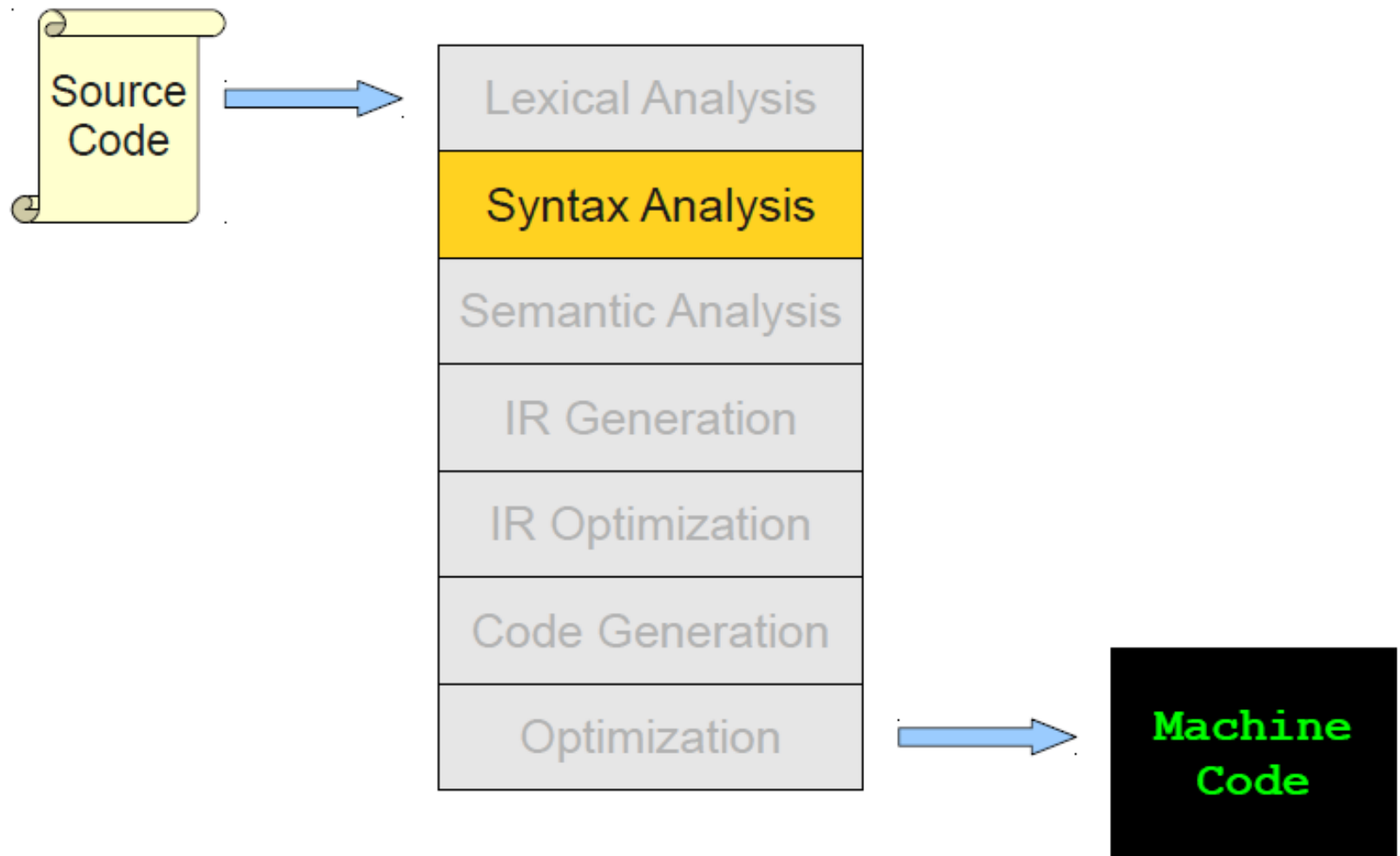
$digit \rightarrow [0-9]$

$letter \rightarrow [a-zA-Z]$



Suppose we need a lexer that recognizes only names, reserved words, integer literals, parentheses, and arithmetic operators.

Where We Are



Introduction

- The syntax analyzer receives the source code in the form of tokens from the lexical analyzer and performs syntax analysis, which determines **the structure** of the program.
- Typically, the result generated by the syntax analyzer is a parse tree.

Top-down parsers

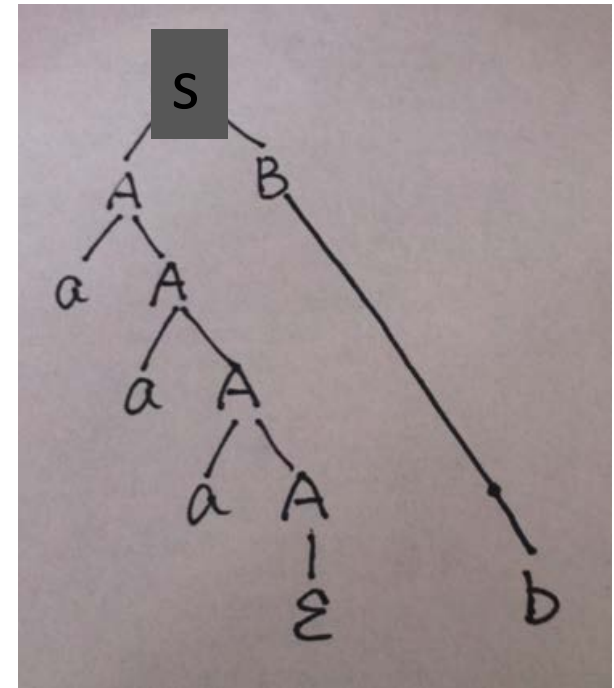
- Start with the start symbol and apply the productions until you arrive at the desired string.
 - Builds a parse tree in **preorder**
 - Corresponds **leftmost derivation**
 - LL parser**

$$\begin{array}{lcl} S & \rightarrow & AB \\ A & \rightarrow & aA \mid \epsilon \\ B & \rightarrow & b \mid bB \end{array}$$

With this grammar

Here is a top-down parse of `aaab`

S	$S \rightarrow AB$
AB	$A \rightarrow aA$
aAB	$A \rightarrow aA$
aaAB	$A \rightarrow aA$
aaaAB	$A \rightarrow \epsilon$
aaa ϵ B	$B \rightarrow b$
aaab	



(Predicative) Top-down Parsers


From deviation aspect

- A nonterminal has more than one RHS
- The correct RHS is chosen on the basis of the **next token** of input (the lookahead)
 - The next token is compared with the first token that can be generated by each RHS until a match is found
 - If no match is found, it is a syntax error

$A \rightarrow bB \mid cBb \mid a$

Example

```
<expr> → <term> {('+' | '-' ) <term>}  
<term> → <factor> {('*' | '/')<factor>}  
<factor> → id | '('<expr>')'
```



```
void expr()  
{...}
```

```
void term()  
{...}
```

```
void factor()  
{...}
```

```
void error()  
{...}
```

- “id” here is not a terminal, but it stands for the identifier token
- language literals are enclosed with ‘ ‘

Program 6:
read x
y = x + 5
write y
halt

```
else if (nextToken == READ) {  
    lex();  
    if (nextToken == IDENT) {  
        printf("INP\nSTA %s\n", lexeme);  
        lex();  
    } else {
```

Homework 2

LMC assembly:

INP		; 1. Read a value from the input into the accumulator (A)
STA x		; 2. Store the value from the accumulator into memory location 'x'
LDA x		; 3. Load the value from memory location 'x' back into the accumulator
STA TMP0		; 4. Store the accumulator value into a temporary memory location 'TMP0'
LDA NUM0		; 5. Load the constant value '5' (NUM0) into the accumulator
ADD TMP0		; 6. Add the value stored in 'TMP0' (original x) to the accumulator
STA y		; 7. Store the result (x + 5) into memory location 'y'
LDA y		; 8. Load the value of 'y' into the accumulator
OUT		; 9. Output the value of the accumulator (x + 5)
HLT		; 10. Halt the program
x	DAT 0	; Memory location for variable x, initialized to 0
TMP0	DAT 0	; Temporary memory to hold x, initialized to 0
NUM0	DAT 5	; Memory location holding the constant 5
y	DAT 0	; Memory location for variable y, initialized to 0

<https://www.peterhigginson.co.uk/LMC/>

Homework 2

Finding more than one error at time:

```
} else {  
    error("Invalid assignment statement.");  
    while (nextToken != IDENT && nextToken != READ &&  
           nextToken != WRITE && nextToken != HALT && nextToken != EOF) {  
        lex();  
    }  
    return;  
}
```