

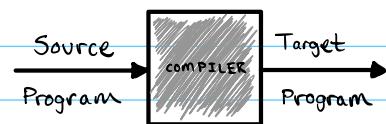
# I. Introduction

a **computer program** is expressed in a **programming language** and transforms an input to an output via the **semantics** of the **symbols** of the **programming language**

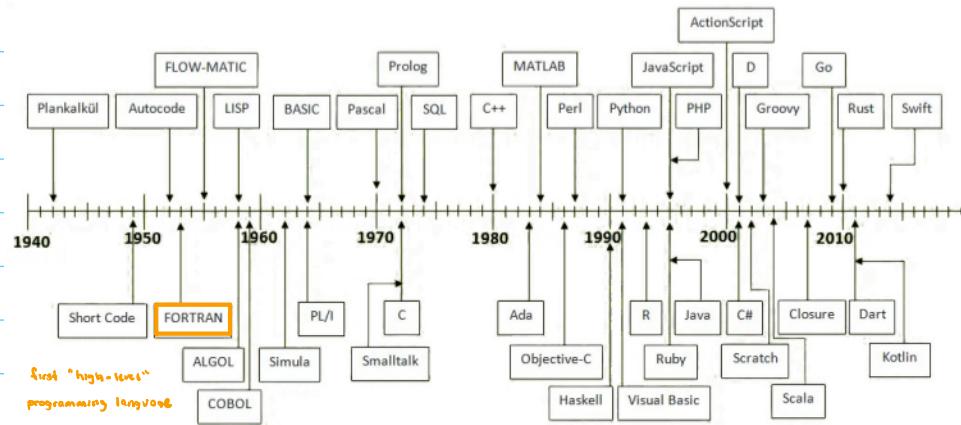
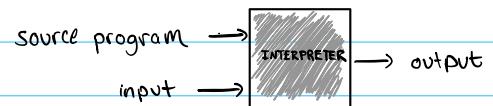
**syntax**: there is a formal set of rules that describes the set of all valid programs that can be written in terms of the **symbols** of the language

**semantics**: gives formal meaning to the **symbols** in a language

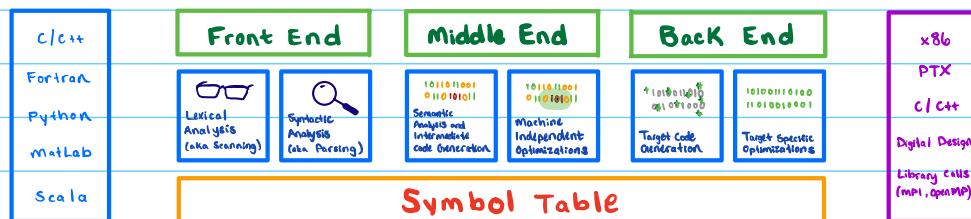
**compiler** translates programs written in one programming language into a program written in another programming language



**interpreter** directly executes the semantics of the symbols in a source program in response to input as its being supplied



## compilation overview



## 2. Lexical Analysis

The process of breaking down a large text into smaller parts

(such as words, phrases, or symbols)

"source tokenization"

deals with "lower level" syntactic structure

strings with same structure get aggregated into the same class

implemented as the scanner

scanner invoked by parser

can be written

- by hand OR - using a scanner generator with reg-ex

SCANNING      input: source program      output: stream of tokens

- decompose the input file (stream of chars) into a stream of strings  
(Lexemes)
- ignore white space (tab, return, etc...)
- assign a token (category) to each string
- job performed by the scanner
- automate tools take regular expressions and produce compilable code  
(scanner generators)  
(e.g. lex, Flex, JFlex, etc...)      (a DFA)

symbol table keeps track of variable names + values

for (int i=0; i < 10; i++)

① ↗ "for", "(", "...", "+", ")"

② ↗ KEYWORD-FOR, LEFT-PAR, ..., PLUS-PLUS, RIGHT-PAR

numbers?:    "(+|-)? [0-9] + \.? [0-9] + " --> NUMBER

removes comments and reports compiler errors

## RE Regular Expressions

denoted by:

- a character      • the empty string  $\epsilon$
- concatenation / "l" or operator      • \* Kleene star

number  $\rightarrow$  integer | real  
integer  $\rightarrow$  digit digit \*  
real  $\rightarrow$  integer exponent | decimal (exponent | e)  
decimal  $\rightarrow$  digit \* (. digit | digit .) digit \*  
exponent  $\rightarrow$  (e | E) (+ | - | e) integer  
digit  $\rightarrow$  [0-9]

## SCANNER RULES

recognizable strings:

1 key words	for, while, if, switch, return
identifiers	i, my-sum, -count-, sum
numbers	integers + floating point
operators	=, +=, +, ++, <, <=, !=, ...
other	(, ), {, }, [., ], ;

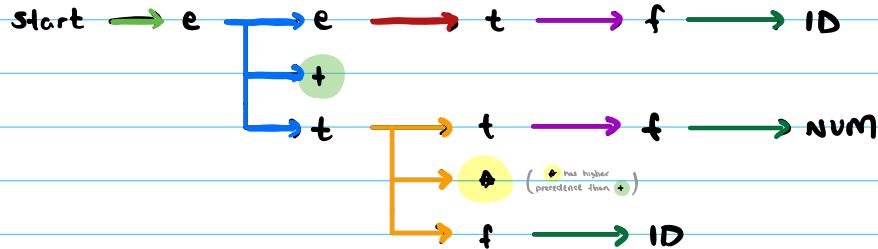
deterministic finite automaton

- A: semi-mechanical, pure DFA (usually realized as nested case statements)
- B: table-driven DFA: write rules, generator produces code
- C: ad-hoc: very case specific

### 3. Syntactic Analysis (parsing)

goal: match stream of tokens to the grammar of the programming language

input: a stream of tokens      output: abstract syntax tree



rules: Start → e      e → t      t → f      f → NUM  
e → e \* t      t → t \* f      f → ID

CONTEXT FREE GRAMMAR: rules applicable independently from "context"

(similar to english tuples: WORD → char | char char\* PUNCTUATION → ., , ; , ? , ! )

four tuples (VTPS):      Variables (non-terminals)      Terminals  
Production Rules      Starting Symbol

expression      expr → id | number | -expr | (expr) | expr op expr      (recursively defined)  
operation      op → + | - | \* | /

- ① begin with starting symbol
- ② chose a production with the starting symbol on the left-hand side
- ③ replace starting symbol with right-hand side of production
- ④ chose a non-terminal A in resulting string
- ⑤ chose a production P with A on the left-hand side,  
and replace A with the right-hand side of P
- ⑥ return to ④ recursively until parsing is complete

**LL parser** (Left to Right, Leftmost derivation)

top-down parse

**LR parser** (Left to Right, Rightmost derivation)

bottom-up parse

Leftmost derivation:

always derive from the first non-terminal

Rightmost derivation:

found in the **sentential form**, moving from:

$$\begin{cases} L \rightarrow R \\ R \rightarrow L \end{cases}$$

(sentential form: only terminal symbols)

"capturing structure"

$A - B - C \Rightarrow$

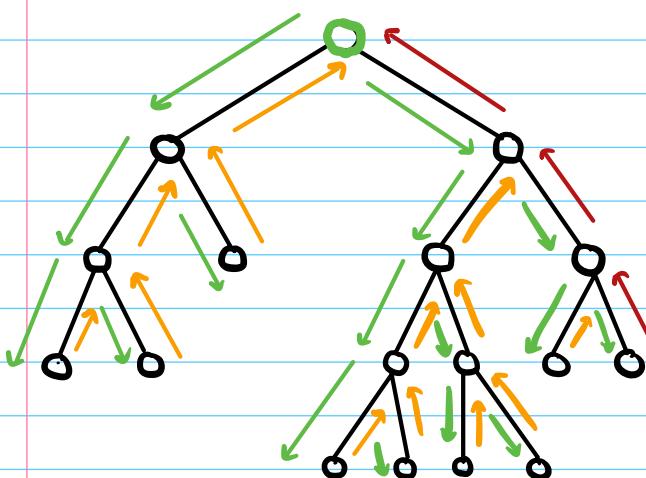
$(A - B) - C$

Left Associativity

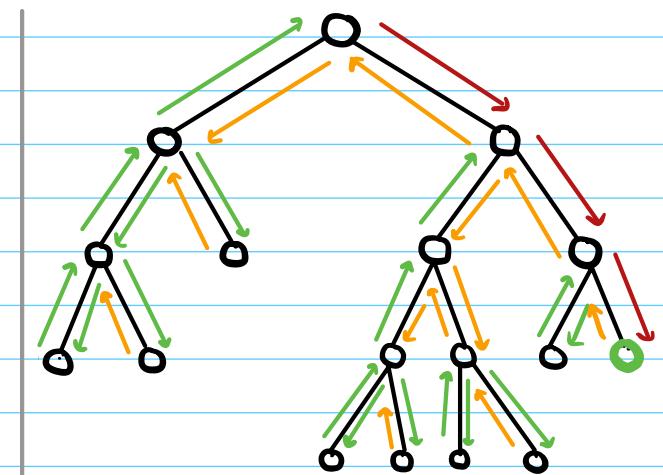
$A - (B - C)$

Right Associativity

**LL parser**



**LR parser**



predictive, recursive-descent

shift-reduce

# LL PARSING TREE EXAMPLE

to parse:  $n^*n$

## production rules

1.  $E \rightarrow TR$
  2.  $R \rightarrow E$
  3.  $R \rightarrow +E$
  4.  $T \rightarrow FS$
  5.  $S \rightarrow E$
  6.  $S \rightarrow ^*T$
  7.  $F \rightarrow n$
  8.  $F \rightarrow (E)$

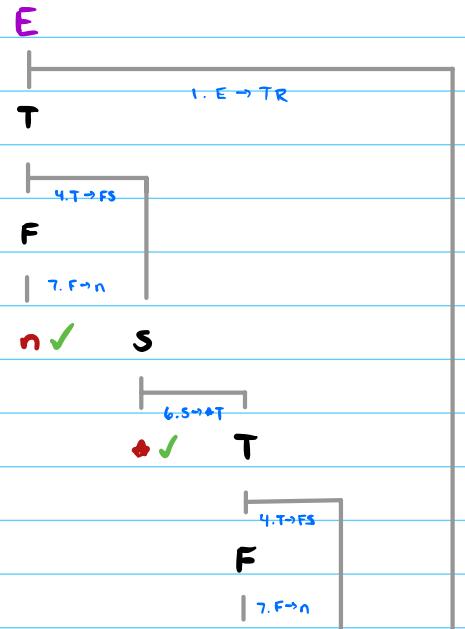
$$D(E, \cap) = 1$$

$$D(T, \textcolor{red}{n}) = 4$$

$$D(F, n) = 7$$

$$D(S, \star) = 6$$

$$D(T, n) = 4$$



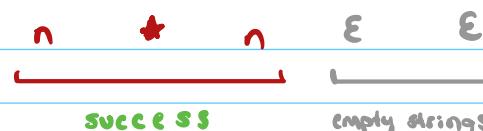
## Table D

	?	+	*	(	)	\$
E	1			1		
R		3	2		2	2
T	4			4		
S		5	6		5	5
F	7			8		

$$D(F, \textcolor{red}{n}) = 7$$

1

## 5. $s \rightarrow e$



## LR PARSING EXAMPLE

productions marked with a (.) become **LR items**

the **LR item** tells us in which production we can be, and in what position

1. $e \rightarrow e \cdot t$	1. $e \rightarrow e + t$	1. $e \rightarrow e \cdot t$	1. $e \rightarrow e + t$	1. $e \rightarrow e + t$
2. $e \rightarrow t$	2. $e \rightarrow \cdot t$	2. $e \rightarrow t$	2. $e \rightarrow t$	2. $e \rightarrow t$
3. $t \rightarrow t^* f$	3. $t \rightarrow t^* \cdot f$	3. $t \rightarrow \cdot t^* f$	3. $t \rightarrow t^* f$	3. $t \rightarrow t^* f$
4. $t \rightarrow f$	4. $t \rightarrow f$	4. $t \rightarrow f$	4. $t \rightarrow \cdot f$	4. $t \rightarrow f$
5. $f \rightarrow (e)$	5. $f \rightarrow (e)$	5. $f \rightarrow (e)$	5. $f \rightarrow (e)$	5. $f \rightarrow (e)$
6. $f \rightarrow ID$	6. $f \rightarrow ID$	6. $f \rightarrow ID$	6. $f \rightarrow ID$	6. $f \rightarrow \cdot ID$
7. $f \rightarrow NUM$	7. $f \rightarrow NUM$	7. $f \rightarrow NUM$	7. $f \rightarrow NUM$	7. $f \rightarrow \cdot NUM$

**basis**      **closure**

\* I DONT UNDERSTAND THE REST OF LR PARSING \*

# Quiz 1

1. what compilation phase is responsible for detecting and reporting syntax errors

a. Syntactic Analysis

2. what compiler component permits the use of specific identifiers

a. Lexical Analyzer (aka scanner)

3. the decision to replace (or not) expressions occurs in what compilation phase

a. Target Specific Optimizations

4. sort the compilation phases

a. 1. Lexical Analysis 2. Syntactic Analysis 3. Syntactic Analysis + Intermediate Code Generation  
4. Machine Independent Optimizations 5. Target Code Generation 6. Target Specific Optimization

5. what is considered the first "high level" programming language

a. FORTRAN

b. what is the output of lexical analyzer

a. set of tokens

7. a VPN —

1a. stands for Virtual Private Network

2a. is used to appear as though you're physically on a network you're not

8. on Linux, return to home directory from anywhere with —

a. cd /~

9. command to compile hello.java using CLI

a. javac hello.java

10. command to get from directory home/A/B/C to home/A/E/F

1a. cd ../../E/F

2a. cd ~/A/E/F

11. command to create new text file

- a. touch book.txt

12. which are lexemes

- a. Identifiers, constants, keywords

13. token category of 3 in sum=3+2

- a. Integer Literal

14. lexical analyzer takes \_\_\_\_\_ as input, returns \_\_\_\_\_

1a. Source Program

2a. Stream of Tokens

15. what goes over the characters of the lexeme to produce a value

- a. Scanner

# Quiz 2

for 1-6 consider the following grammar:

$$1) v \rightarrow v[\text{expr}]$$

$$2) v \rightarrow \text{ID}$$

$$3) \text{expr} \rightarrow \text{NUM}$$

$$4) \text{expr} \rightarrow v$$

① top-down, right-most derivation of the string "ID [NUM] [ID]"

