

Principle of Compiler Construction

compiler → software system that translates one source language into another language

error identification easier in ~~compiler~~ interpreter

Compiler → 6 phases

1. ↗ lexical analysis phase (valid identifier, valid data type, valid exp)
→ finite automata (reg ex pattern match)
2. → Syntax analyzing (whether statement is unambiguous / follows a production rule)
3. → Semantic analyzer (checks for semantic consistency)
4. → Intermediate code generator (role of compiler organization)
5. → Code optimization
6. → Code generator (generates target code)

* Preprocessor, compilers, assemblers, linkers
skeletal source prog

preprocessor

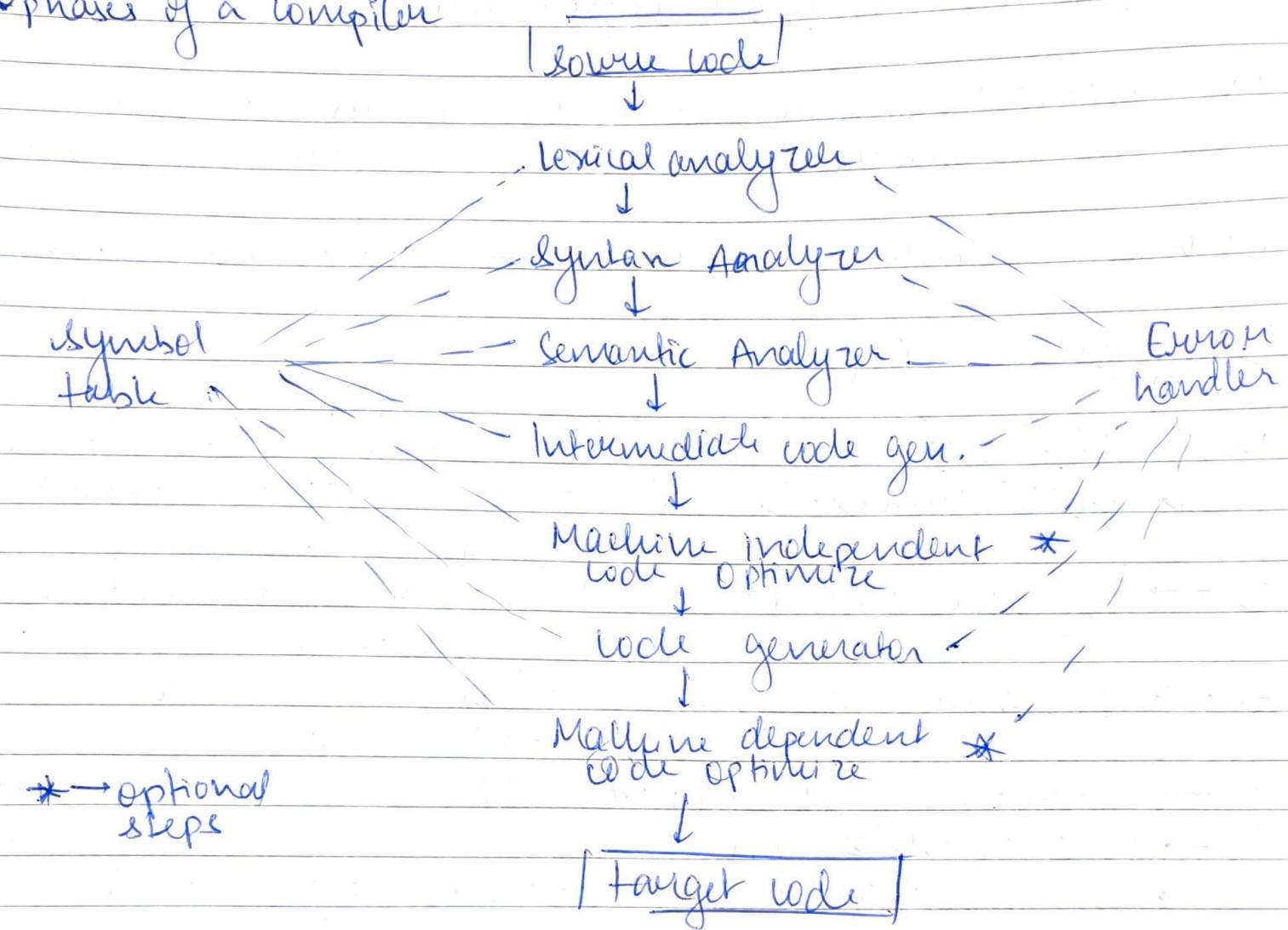
↓
compiler

↓
assembler

↓
linker

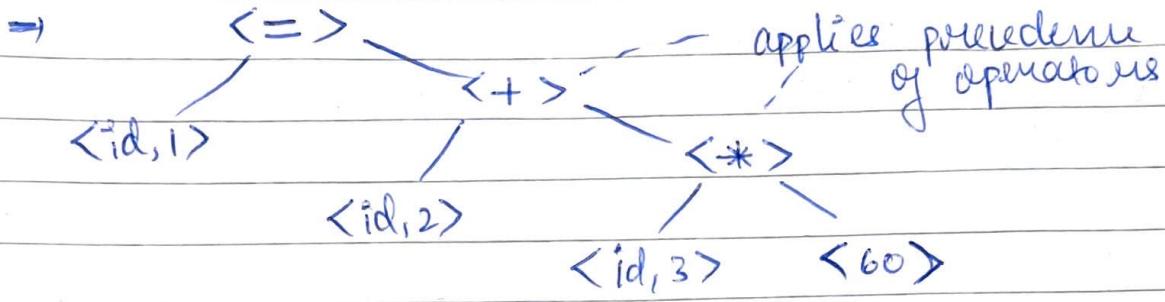
↓
absolute
machine
code

phases of a compiler



e.g. position = initial + ratio * 60
 ↓
lexeme → token
 <token-name, attribute>
 index value
⇒ <id, 1> <= > <id, 2> <+> <id, 3> <*> <60>

Syntax analyzer → converts tokens to syntactic tree
condensed form
of parse tree



- grammar used by syntax analyzer should be unambiguous in nature
- semantic analyzer ~~also~~ checks for semantic consistency.

* Intermediate code generator

→ max 2 operator and 1 output

$$t_1 = \text{inttofloat}(60)$$

$$t_2 = id_3 * t_1$$

$$t_3 = id_2 + t_2$$

$$id_1 = t_3$$

* machine ~~independent~~ independent code optimizer optimizes intermediate code and is responsible for tasks like removing redundant variables

eg $t_1 = id_3 * 60.0$ ↗ 2 line
 $id_1 = id_2 + t_1$ ↗ optimized code

* Code generator

→ generates machine code

on LD RL, id3

MUL R1, R1, #60.0

LD R2, ~~R2~~ id2

SUM R1, R1, R2

ST ~~R1~~ ID31, R1

* Machine dependent code optimizer
→ Optimizes the machine code generated by
code generator

II Passes

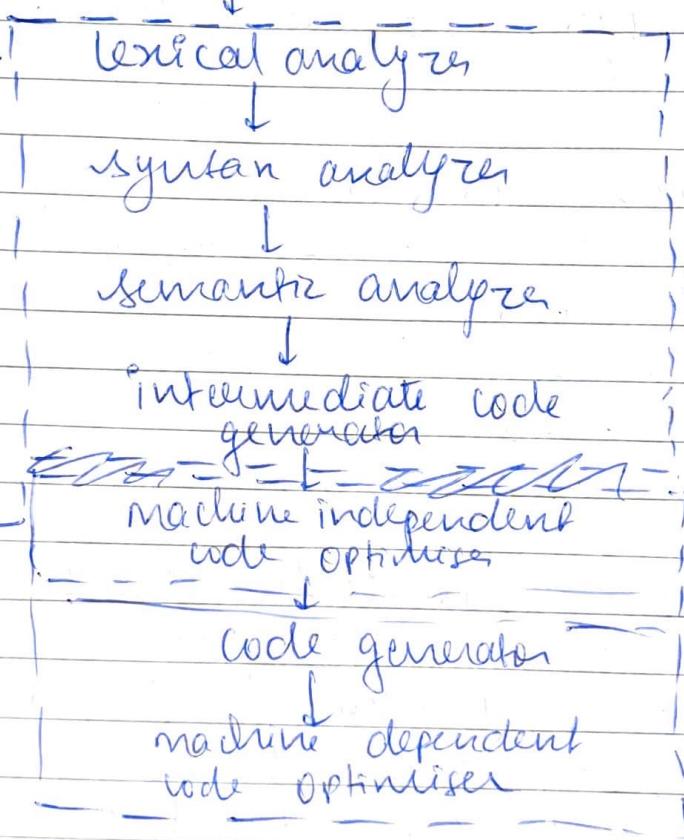
single pass and multi-pass compiler
(low level code)

front end

analysis part

back end

optimiser part



multi-pass / 2 phase compilers generate
machine code for a specific machine

No. of phases

eg

$$\text{position} = \text{initial} + \text{rate} * 60$$



Lexical analysis.

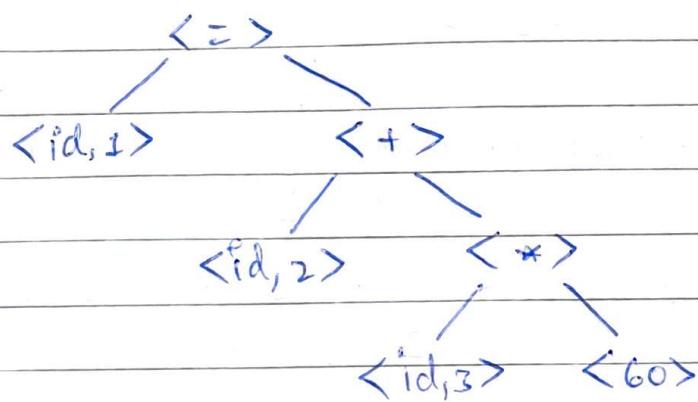
Token stream

$\langle \text{id}, 1 \rangle \langle = \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{60} \rangle$



Syntax Analysis

Syntax tree



* Ambiguous and Unambiguous grammar

eg $E \rightarrow E + E \mid E * E \mid \text{num} \mid E = E$

To remove ambiguity, assign precedence to diff operators

lowest precedence

$$\Rightarrow E \rightarrow E = T \mid T$$

$$T \rightarrow T + F \mid F$$

$$F \rightarrow F * G \mid G$$

$$G \rightarrow \text{num} \mid \langle \text{id} \rangle$$

* To remove ambiguity:

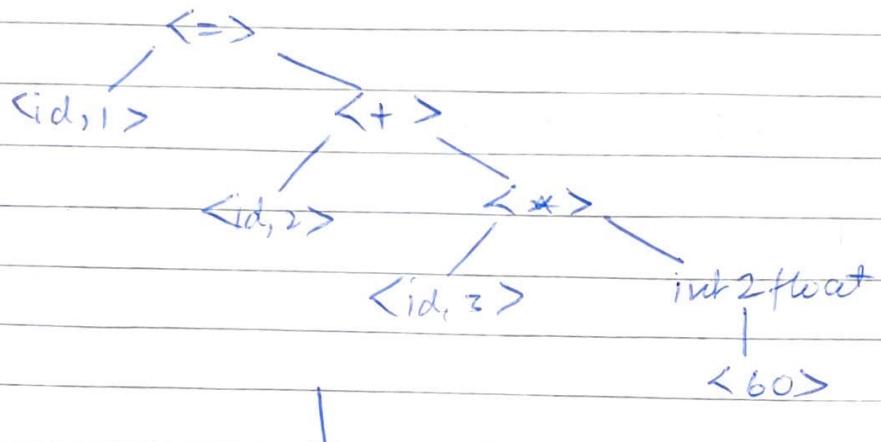
① Precedence

② Left to right (Associativity)

* Syntax Tree

|
Semantic Analysis

|
↓ syntax tree (semantically correct)



|
Intermediate code generation (machine independent component)

$$t_1 = \text{int2float}(60)$$

$$t_2 = t_1 * \text{id}_3$$

$$t_3 = t_2 + \text{id}_2$$

$$\text{id}_1 = t_3$$

right side has max
2 operands and
1 operator

|
Code optimization

using DAG

$$t_1 = \text{id}_3 * 60.0$$

$$\text{id}_1 = \text{id}_2 + t_1$$

|
Code generator (machine dependent component)

LD R1, id3

MUL R1, R1, #60.0

LD R2, id2

ADD R1, R1, R2

ST id1, R1

Compiler Construction Tools

① Scanner generator

↳ lexical analysis

② Parser generator

↳ syntax and semantic analysis

③ Syntax directed translation tool

↳ provides multiple routines to go through
syntax trees and generate intermediate code

④ Data flow analysis generator (gives flow of program)

⑤ Code optimizer generator

⑥ Compiler construction toolkit

Evolution of Languages

1940 → machine language → 1st generation

1950 → Assembly → 2nd generation

1955 → HLL (C, Pascal, etc) → 3rd gen

→ Problem specific lang → 4th gen
(SQL, postscript)

→ Logic and constraint based → 5th gen
(PROLOG, OPS 5)

based on nature / working culture

↓
Imperative

→ flow is defined
you how the
work is to be
done

e.g. C, Pascal

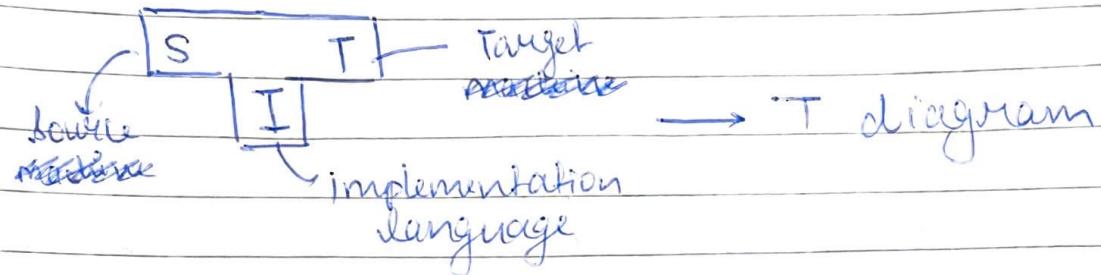
↓
Declarative

→ flow is not
defined. Only
work is defined

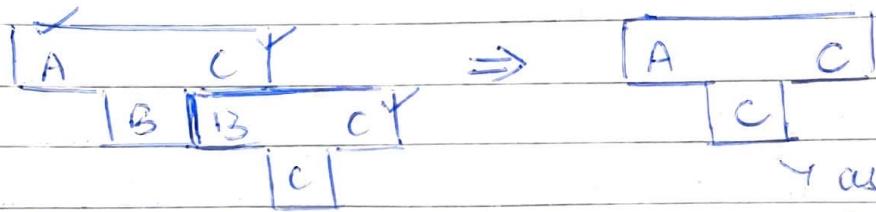
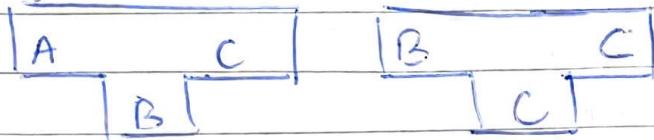
e.g. Relational algebra, PROLOG

* Cross compiler

→ a compiler which runs on one machine and generates target for another machine



case -1



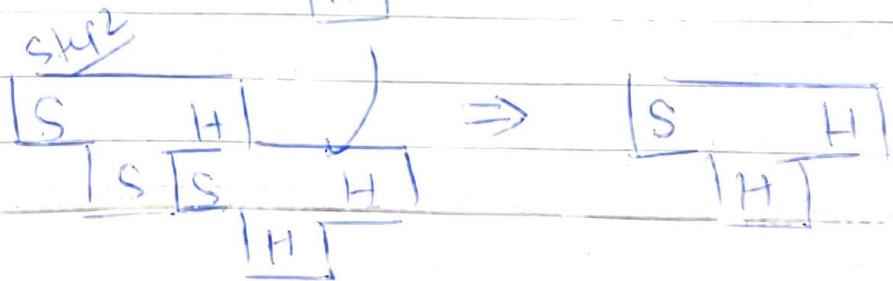
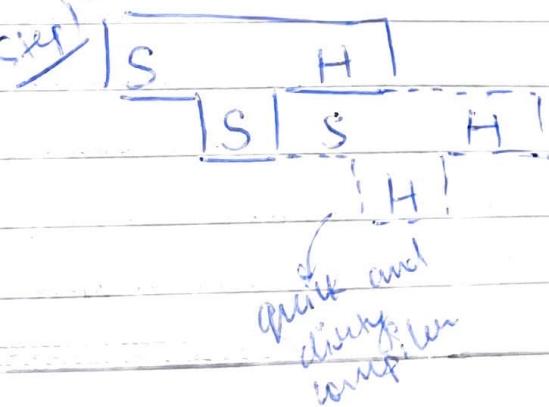
as B ultimately compiles to C, C is the implementation language

case -2



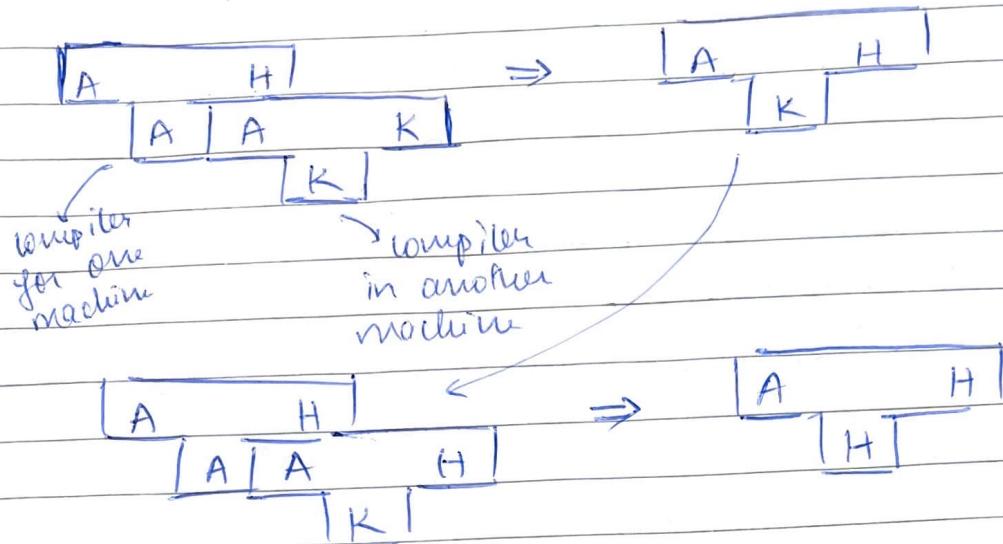
* Bootstrapping (self compiling compiler)

write a compiler in a simple language, which further generates compiler for a more complex language



* Portability

Any generated program should be executable on multiple machines



Just In Time Compilation

Compilation process occurs at the time of execution.

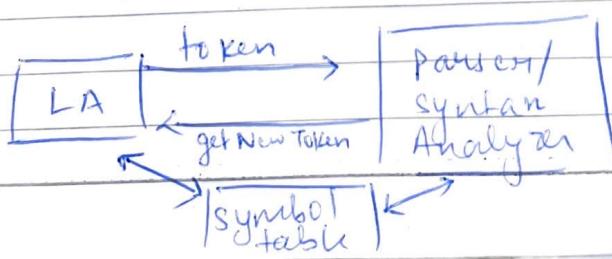
e.g. Java (JVM)

C# (Common language Runtime)

Lexical Analysis (Scanner Generation)

takes source programs, identifies lexemes in it and generates token stream

source prog → [LA] → token stream



① Token generation

eg int a = b + c;
T₁ T₂ T₃ T₄ T₅ T₆ T₇

② Error handling

→ tries to implicitly handle errors, if not
then send information to error handler

③ Remove / Handles whitespace

(space, tabs, newline, comments etc.)

④ Macro Expansion

* Tokens

Identifiers
eg <id, pointer>

Keywords
tokens are same as

Keywords

Operators
eg <op>

numbers

eg <number, pointer>

TBD

eg <int>

punctuations
eg <;>

Brackets

eg <(>, <{>}

Token

↳ is an interpretation of valid lexeme

↳ atomic element which cannot
be further divided and is
useful to the program

to identify valid lexeme, patterns are used

↳ are basically regular expressions

Input Buffering

$E = m * c ** 2$

—
↳ to resolve this token, LA needs to
look ahead

⇒ operators as they come are put into an
input buffer up until a valid lexeme is
obtained / error generated / invalid lexeme
encountered.

* One buffer technique, Two Buffer technique
if buffer ~~exist~~ is filled i.e lexeme → A second buffer is initialized
every time current buffer

* to check if buffer is filled, sentinels are used

↳ a pointer at end of lexeme
in a buffer is a sentinel

acaction:

* Regular expressions

S. All strings of lowercase letters that contains
5 vowels in-order

$$\rightarrow C, \text{no-vowel-set} = (b+c+d+f+g+h+j+k+l+m+n+p+q+r+s+t+v+w+x+y+z)$$

$$\Rightarrow \cancel{\text{no-vowel-set} * a} \cancel{\text{no-vowel-set} * e} \cancel{\text{no-vowel-set} * i}$$

$$\cancel{\text{no-vowel-set} * o} \cancel{\text{no-vowel-set} * u} \cancel{\text{no-vowel-set}}$$

$$\Rightarrow (\text{non-vowel})^* aa^* (\text{non-vowel})^* ee^* (\text{non-vowel})^* ii^* (\text{non-vowel})^* oo^* (\text{non-vowel})^* uu^* (\text{non-vowel})^*$$

ab a ababa

~~$$(\text{non-vowel} + a)^* a a^*$$~~
~~$$(b+o+o) * a (c+e)^* e (c+i)^* i (c+o)^* o (c+u)^* u c^*$$~~

S. All strings of as and bs, that don't contain subsequence ab

~~$$(a^* b^*)$$~~
~~$$(aa^* ba a^*)^*$$~~
~~$$(bb^* b^*)$$~~
~~$$(a^* ba a^*)^*$$~~

$$b^* (a^* b a^* (e+b))^*$$

$$b^* (a + ab)^*$$

$$b^* (a + (ab))^*$$

B. No subsequence ab

~~$$b^* (a^* a^* + ab) a^*$$~~

$$b^* a^* (ab + e) a^*$$

$$b^* (a^* b) a^*$$

b a a a b

1. RE \rightarrow minimized DFA
2. NFA \rightarrow DFA
3. RE \rightarrow NFA

/* / / / */

& ~~lex~~

Q Regen for comments consisting of string surrounded by /* */ without an intervening */ until it is inside " "

$C \rightarrow \text{non } \{, *, /, "$

~~/*~~ /*(C*)(ee(CC+*+/)*)) + e)C*/*/ /

/* (C*(C*+/)*+e)C* (C*+/*+e))

(/* (C(C*+/)*+e)*) */ (/) * (*) *

111*

~~/* C*(C*+/)*+e)C*~~

* Lex form of Regular Expression

Expression

matches

Examples

C

any non operator character

a, b, c, ...

\c

escape characters

*

"S"

any character except newline

"**"

*

any character except newline

a.*b

^

beginning of a line

^abc.

\$

ending of a line

abc\$

[s]

any one character from list

[abc]

[^s]

no character from set

[^abc]

g1*

any 0 or more occurrence

a*

g1+

only 1 or more occurrence

a+

g1?

Zero or one

a?

$$\pm (x), (-\delta)^{\wedge})$$

$\{M_1, M_2\}$	M_1 to M_2 occurrences	$a \{1, 2\}$
M_1, M_2	M_1 , followed by M_2	
$M_1 \mid M_2$	M_1 , or M_2	
(M_1)		(a)
M_1 / M_2	M_1 when followed by M_2	
wordspace	(blank / newline). tab) +	

* Recognition of Tokens

Stmt → if expr then Stmt |
 if expr then Stmt else Stmt | e
 Expr → term relational operator term | term

term → id / num

Since exp^n is a loop
eg if ($a \equiv 2$) then
 $a++;$ } exp^n
start

RegEx for if statement:
~~if \([a-zA-Z]+\)~~