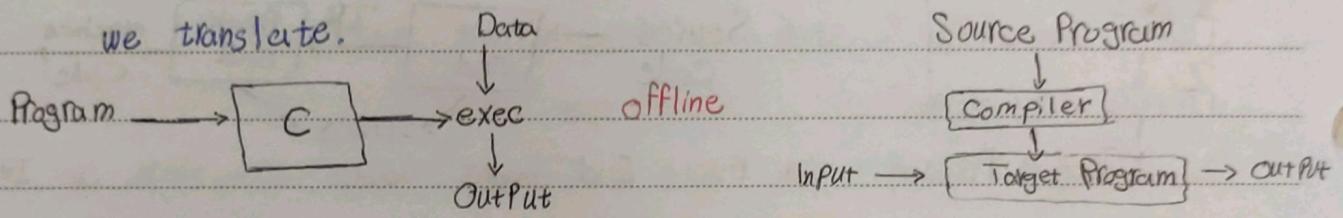


What is a compiler? Takes as input a program written in one language and translates it into a functionally equivalent program in another language.

* Source is usually high-level, target is usually low level.

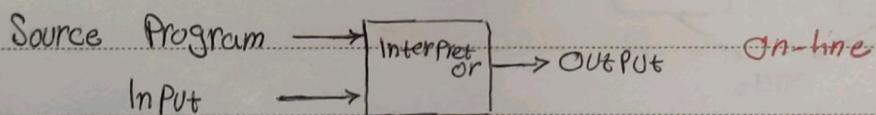
Translation Approaches

newer \leftarrow 1) Compiler $\&$ By following physical structure of program we translate.

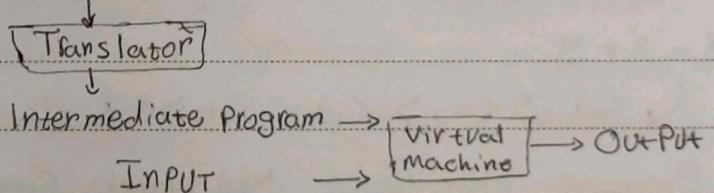


Older \leftarrow 2) Interpreter $\&$ translate the source code to machine code ^{online}

* The code is translated while it is running.



3) hybrid Approach $\&$ Source Program



How Does Compiler work?

tokenize 1) Lexical analysis (Scanning) $\&$ Identifying logical piece of the description

Diagramming 2) Syntax analysis (Parsing) $\&$ Identifying how does those pieces relate to each other sentence

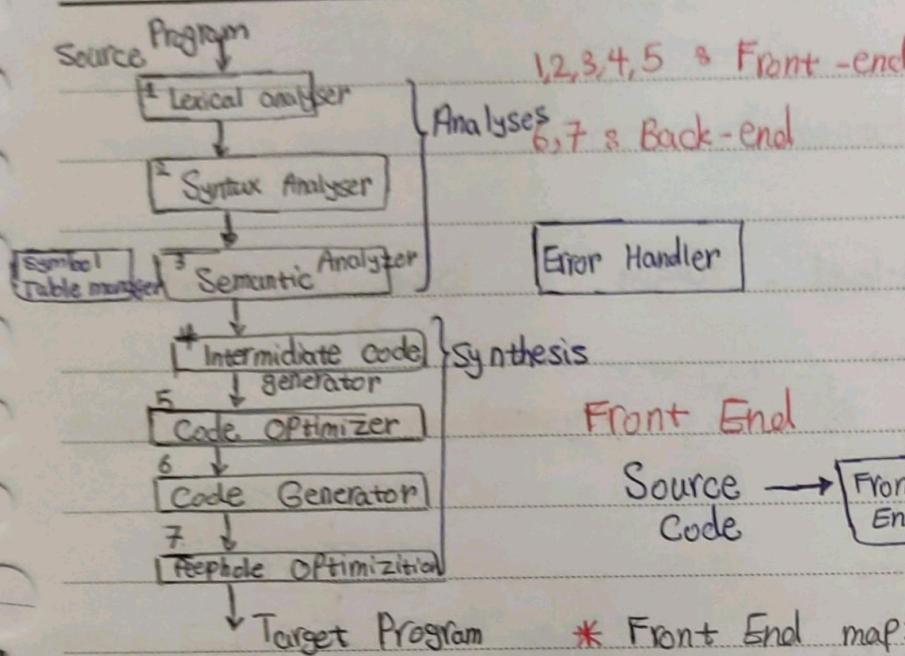
3) Semantic analysis $\&$ Identifying the meaning of the overall structure.

4) IR Generation $\&$ Design one possible structure.

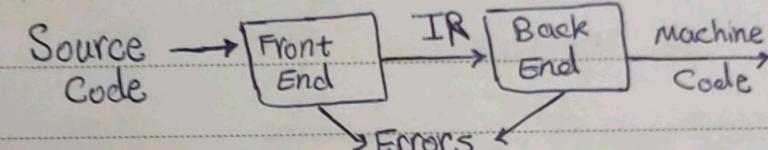
5) IR Optimization $\&$ Simplify the intended structure.

6) Generation $\&$ Fabricate the structure.

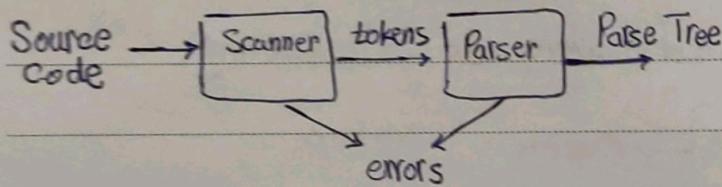
7) Optimization $\&$ Improving the resulting structure.



Front End



- * Front End maps Source code into an IR representation
- * Back End maps IR onto machine code
- * Simplifies retargeting



Scanner & Maps characters into tokens

$x = x + y$

$\hookleftarrow \langle id, x \rangle \langle =, \rangle \langle id, x \rangle \langle +, \rangle \langle id, y \rangle$

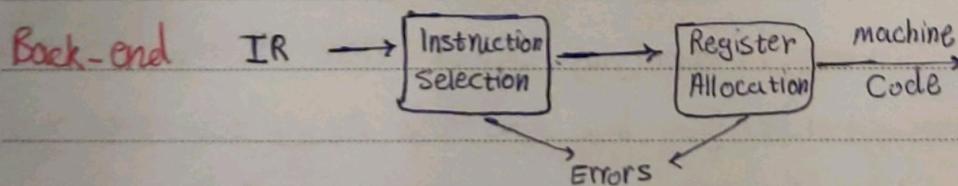
Parser &

Recognize context-free syntax

Guide context-sensitive analysis

Eliminate white space (tabs, blanks, comments)

Produce meaningful error messages

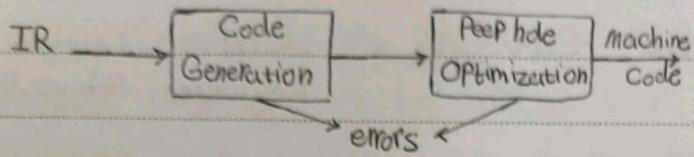


* Translate IR into machine code

* Choose instructions for each IR operation

* Decide what to keep in registers at each point

Two Main Components of Back-end



Code Generator

Produce compact fast code

Use available addressing modes

Peephole Optimization

Limited resources

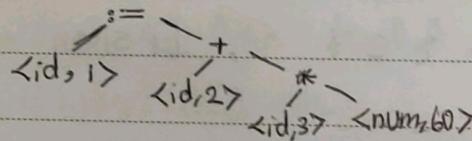
Optimal allocation is difficult

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

Lexical analyzer (Scanner)

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

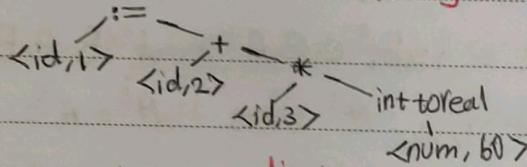
Syntax analyzer (Parser)



Symbol Table

1	Position	real
2	initial	real
3	rate	real

Semantic analyzer



Intermediate Code Generator

$t_1 := \text{int to real} (60)$

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

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

$\text{id}_1 := t_3$

Code optimizer

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

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

Final Code Generator

LD R1, id3

MUL R1, R1, # 60.0

LD R2, id2

ADD R1, R1, R2

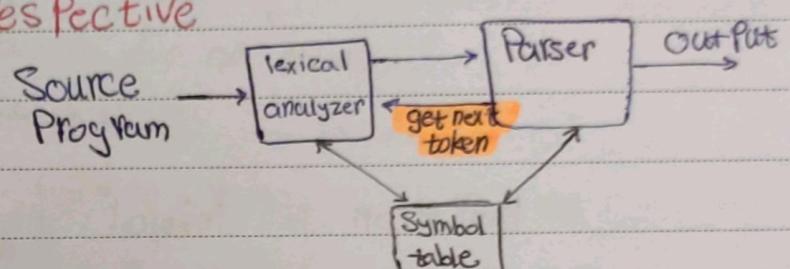
ST id1, R1

Lexical Analysis Partition input string into tokens

* Token is A syntactic Category : Identifier, Integer keyword, whitespace

Lexical Analysis in Perspective

Symbol key	lexeme	Table type
1	Position	real
2	initial	real
3	rate	real
	<u>token</u>	<type, attribute>



Atomic Regular Expressions

Single character ' c ' = $\{ "c" \}$, Epsilon $\epsilon = \{ "" \}$

Compound Regular Expressions

Union $A + B = \{ s | s \in A \text{ or } s \in B \}$

Concatenation $AB = \{ ab | a \in A \text{ and } b \in B \}$

Iteration $A^* = \bigcup_{i \geq 0} A^i$ where $A^i = A - i \text{ times} - A$

Regular Expressions The regular expressions over Σ are the smallest set of expressions including $\rightarrow L(\cdot)$

ϵ , ' c ' where $c \in \Sigma$, AB where A, B are rexp over Σ
 A^* where A is a rexp over Σ , $AB \cup \cap \cap \cap \cap \cap$

Syntax vs. Semantic

$L(\epsilon) = \{ "" \}$, $L('c') = \{ "c" \}$, $L(A+B) = L(A) \cup L(B)$

$L(AB) = \{ ab | a \in L(A) \text{ and } b \in L(B) \}$, $L(A^*) = \bigcup_{i \geq 0} L(A^i)$

* integer = digit digit*, $A^+ = AA^*$

identifier = letter (letter + digit)*

option $A + \epsilon$ $A ?$ $\frac{\epsilon}{\epsilon}$ بصيغة

Regular Expressions \Rightarrow Lexical Spec.

1) Write a regexp for the lexemes of each token

* Number: digit + * keyword: 'if' + 'else' + ... * Identifier: letter (letter + digit)*

2) Construct R, matching all lexemes for all tokens.

$R = \text{keyword} + \text{Identifier} + \text{Number} + \dots = R_1 + R_2 + \dots$

3) Let input be $x_1 \dots x_n$. For $1 \leq i \leq n$, check $x_1 \dots x_i \in L(R)$

4) If success, then we know that $x_1 \dots x_i \in L(R_j)$ for some j

5) Remove $x_1 \dots x_i$ from input and go to (3)

What if $x_1 \dots x_i \in L(R)$ and also $x_1 \dots x_k \in L(R)$, $i \neq k$

→ Rule "maximal munch": Pick longest possible string in $L(R)$

What if $x_1 \dots x_i \in L(R_j)$ and also $x_1 \dots x_i \in L(R_k)$, $j \neq k$

→ Rule: use rule listed first (j if $j < k$)

Dr. حقيقة العيت ديم: identifier \rightarrow this keyword

What if $x_1 \dots x_i \notin L(R_j)$. No rule matches a prefix of input?

هذا خطأ ديني

ما لا يحوله إلى زبائن عرضت خارج عن ولئن وقتي لست بحاجة محسنة تكون \Rightarrow accepting state; if it is

حالاً ما لا يحول زبائن عرضت start state \rightarrow start state

اول كونه ترين سويها لغير بعد حصلت زبون

* IF (end of input) and (in accepting state) \Rightarrow accept

Epsilon moves ~~امثلة~~ ε-moves $(A) \xrightarrow{\epsilon} (B)$ \rightarrow $\text{لعن دليل جعلية ريجار}$

Machine can move from state A to state B without reading input

Deterministic Finite Automata (DFA)

Non-deterministic Finite Automata (NFA)

One transition per input per state

Can have multiple transitions for one input given state.

No ε-move

Can have ε-move

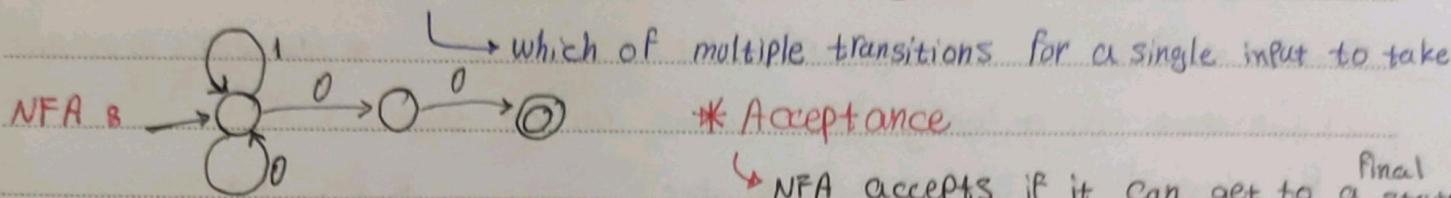
دون عروض داريم لكن

PAPCO

- comes \Rightarrow NFA \rightarrow if DFA \Rightarrow *

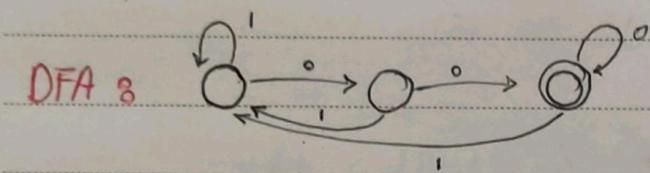
* A DFA can take only one path through the state graph. (completely determined by input)

* NFAs can choose → whether to make ϵ -moves



* DFAs are faster to execute (there are no choices to consider)

* NFAs are, in general, smaller (Sometime exponentially smaller)



Regular expression to Finite Automata (High-Level Sketch)

Lexical Specification → Regular expressions → NFA → DFA → Table-Driven Implementation of DFA

* For each kind of rexp, define an NFA

↳ NFA for rexp M → $M \circlearrowright$

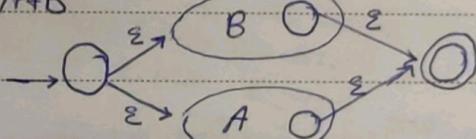
For ϵ → $\circlearrowright \xrightarrow{\epsilon} \circlearrowright$

for input a → $\circlearrowright \xrightarrow{a} \circlearrowright$

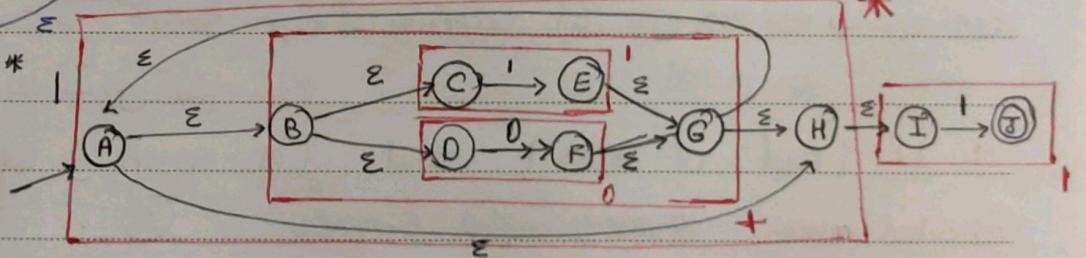
For AB → $A \xrightarrow{\epsilon} B \circlearrowright$

For $A+B$

For A^* → $\circlearrowright \xrightarrow{\epsilon} A \xrightarrow{\epsilon} \circlearrowright$



clue 8 $(1+0)^*$



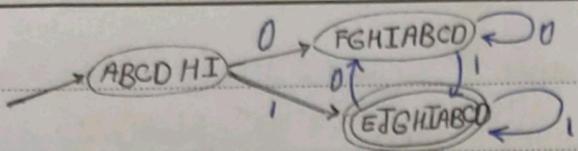
NFA to DFA : The Trick

States of the NFA

- 1) Simulate the NFA
- 2) Each state of DFA = a non-empty subset of states of the NFA
- 3) Start State = ϵ -closure of the start state of NFA
- 4) Add a transition $S \xrightarrow{a} S'$ to DFA iff

PAPCO S' is the set of NFA states reachable from any state in S after seeing the input a , considering ϵ -moves as well.

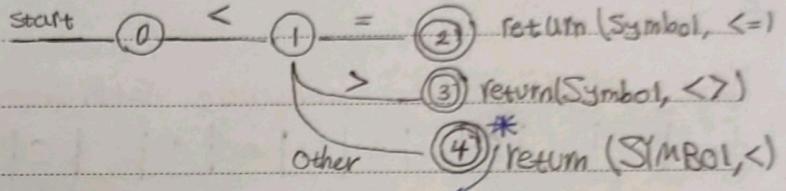
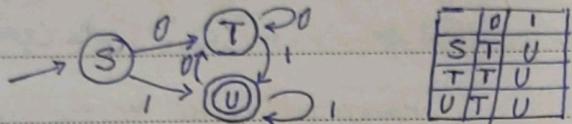
- 5) Final States & Subsets that include at least one final state of NFA.



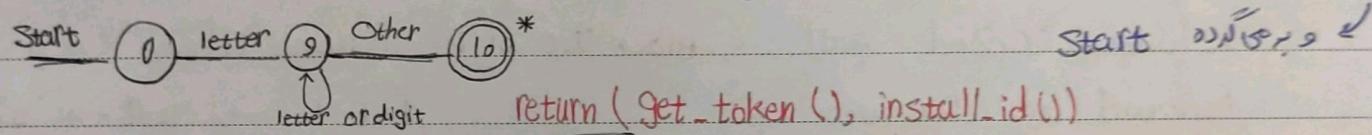
8 قبلي NFA , DFA

* A DFA can be implemented by a 2D table T : States vs. input symbol

$S_i \xrightarrow{a} S_k$ define $T[i, a] = k$



we have accepted ">" and have read "other" character that must be unread. That is moving the input pointer one character back

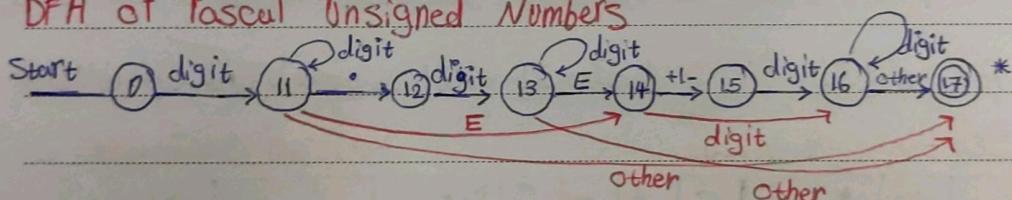


```
return (get_token(), install_id())
```

• returns either a keyword or ID based on the type of the token.

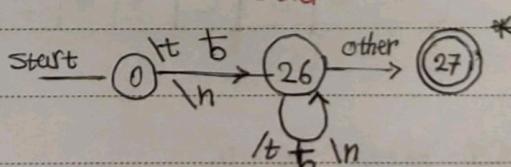
→ If the token is an ID, its lexeme is inserted into the symbol table (only one record for each lexeme) and lexeme of token returned.

DFA of Pascal Unsigned Numbers



return (NUM, lexeme ^{of the} number)

DFA of whitespace



Lexical % تعریضی

Syntax 7.60

Error Recovery abc \$ de \$:= 0

Cementie ½ 30

از پیش رانسون اخوند تا پیش رانسون (Panic de \$) دور

دباره روش دخنی de\$ و \$ Panic می باشد

و بہ نبی مسیح ارسلان

Cascaded errors

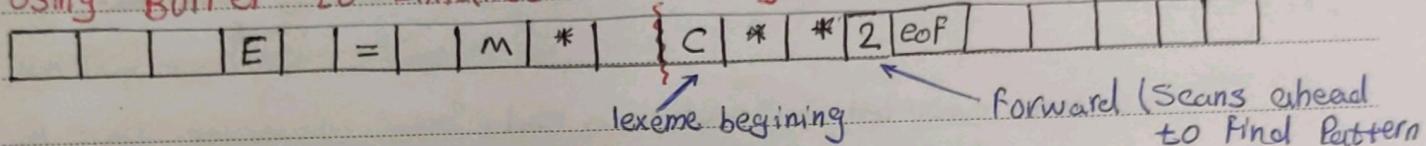
error handling → error recovery b error correction

هین دالن
لقتیم

لِصَاحِبِ الْحُكْمِ

- 1) Delete one character from the remaining input لکی اس این جا کار ام اسٹرنگ
حالتہ درکرد و بحال میکنے
- 2) Insert a missing character into the remaining input حال پر ایدھ سنت
حالتہ درکرد و بحال میکنے
- 3) Replace a character by another character دستیابی درجی همایوں
حالتہ درکرد و بحال میکنے
- 4) Transpose two adjacent character آخرین جا کا دو حروف دستیار سوچ کر تبدیل کر لیں رونگزیری

Using Buffer to Enhance Efficiency



if forward at end of first half then begin

reload second half; ← Block I/O is 0 = forward lexeme, 1B
beginning

forward := forward + 1 $F = F - 1$ new state

end

else if forward at end of second half then begin

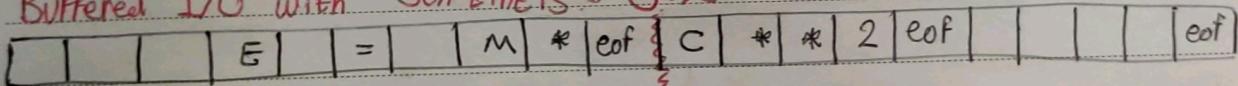
reload first half; ← Block I/O

move forward to beginning of first half

end

else forward := forward + 1 طالی قرین کوں یا منی در اندھہ باخڑوئہ!

Buffered I/O with Sentinels & ijkts



forward := forward + 1

if forward is at eof then begin

if forward at end of first half then begin

reload second half;

forward := forward + 1

end

else if forward at end of second half then begin

reload first half; move forward to beginning of first half

PAPCO

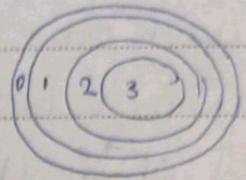
end

else /* eof within buffer Signifying end of Point 7
lexical analyisis

end 2nd eof \Rightarrow no more input!

Chomsky Hierarchy نحویات چومسکی Parser

0 Unrestricted $\alpha / \beta \rightarrow \alpha * \beta$



1 Context-Sensitive $|LHS| \leq |RHS|$

Right linear لکی طایر مکعبات

2 Context-Free $|LHS| = 1$

left linear لکی طایر مکعبات

3 Regular $|RHS| = 1 \text{ or } 2, A \rightarrow a \text{ or } aB, \text{ or } A \rightarrow aB a$

هزیانی یک گرام دارد لکی طایر مکعبات

$$G = \langle N, T, S, P \rangle$$

↓ ↓ ↓ ↑
 non-terminal terminal start symbol Product-rule

a, b, c لکی طایر مکعبات, A, B, C لکی طایر مکعبات α, β, γ لکی طایر مکعبات, x, y, z لکی طایر مکعبات (هرجی) Σ لکی طایر مکعبات

Plus $G: \langle \{E\}, \{+, *\}, \{., (\), id\}, E, P \rangle$ درست آوریدم P که در اینجا درست آوریدم

$P: E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E), E \rightarrow id$

Parse tree درخت نحوی $* \rightarrow 1 + 2$ سبل $A \rightarrow \Sigma$ فرآورده ایمیشن

هستی دیگر $=$ terminals متون

* Non-terminals are written uppercase, Terminals lower case

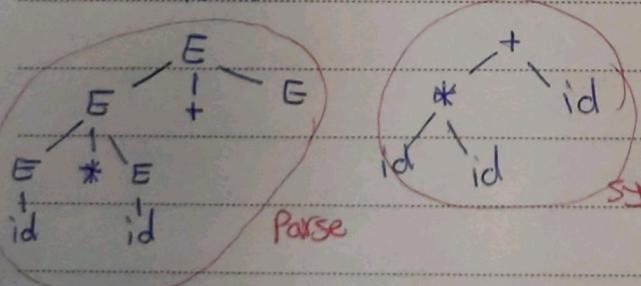
* An in-order traversal of the leaves is the original input

The right-most and left-most derivations have the same Parse tree.

(Parse) tree VS. (syntax tree) درخت نحو VS. درخت ساختار

$id * id + id$

terminal لکی طایر مکعبات



LMD (Left most Derivation) از چپ عذر طایر مکعبات

العیت های بالاتر پاسخ داده اند لکی طایر مکعبات

parse tree لکی طایر مکعبات

برای لکی طایر مکعبات

: $E \rightarrow E+E | E^*E | (E) | id$

$\xrightarrow{\text{but}}$ $E \rightarrow E+T | T$

$T \rightarrow T^*F | F$

$F \rightarrow (E) | id$

enforce Precedence of $*$ over $+$

PAPCO