

Lecture 2

Compiler: is a software translates programs in a source language into target language.

Interpreter: does not produce a target program, It reads an executable program and produces the results.

Hybrid Compilers: Source program is compiled into an intermediate program, which is later interpreted by an interpreter.

Lexical analyzer (scanner): Reads source program as a stream of characters (Char by Char from left-to-right) and classify it into tokens.

Token Types: Key words, Identifiers, Relop, Num, Op, White spaces.

Regular expression: expression that matches sets of strings.



Lecture 3, 4

Finite Automata consists of:

1- $S \rightarrow$ Set of states $S = \{q_0, q_1, \dots\}$

2- $\Sigma \rightarrow$ Set of symbols $\Sigma =$

State	Input 1	Input 2
-------	---------	---------

3- $\delta \rightarrow$ Transition function

4- $s_0 \rightarrow$ Start state $s_0 = \{q_0\}$

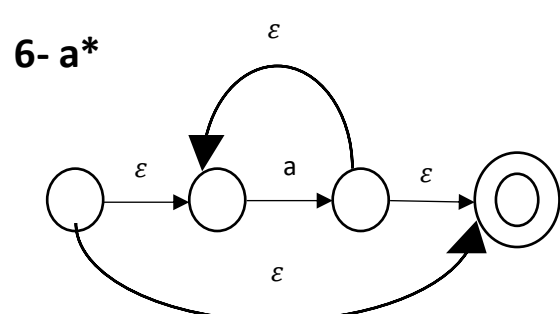
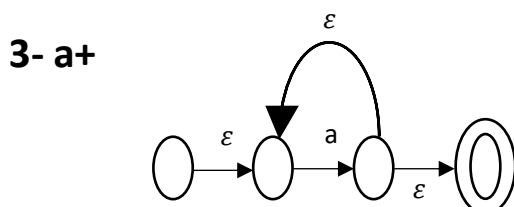
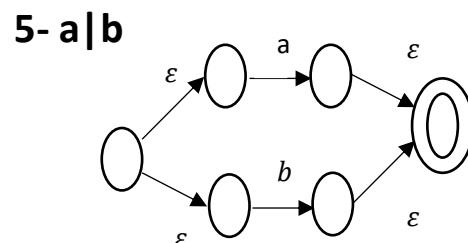
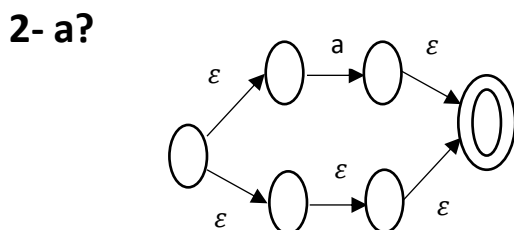
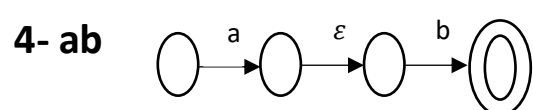
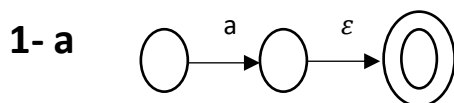
5- $F \rightarrow$ Accepting state $F = \{q_5, q_6, \dots\}$

Scanner phases:

- 1- Pattern specification using regular expression.
- 2- Pattern recognition using finite automata.

1- Thompson's construction (RE \rightarrow NFA)

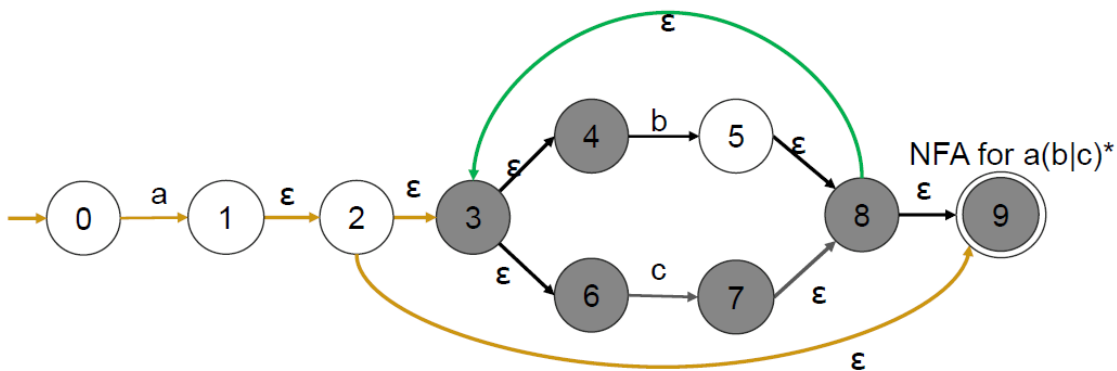
Convert Regular expression to non-deterministic finite automata



Lecture 5

2- Subset construction (NFA -> DFA)

Convert non-deterministic finite automata to deterministic finite automata.



NFA State	a	b	c
$d0 = \{0\}$	$d1 = \{1,2,3,4,6,9\}$	none	none
$d1 = \{1,2,3,4,6,9\}$	none	$d2 = \{3,4,5,6,8,9\}$	$d3 = \{3,4,6,7,8,9\}$
$d2 = \{3,4,5,6,8,9\}$	none	$d2 = \{3,4,5,6,8,9\}$	$d3 = \{3,4,6,7,8,9\}$
$d3 = \{3,4,6,7,8,9\}$	none	$d2 = \{3,4,5,6,8,9\}$	$d3 = \{3,4,6,7,8,9\}$

3- Hopcroft's algorithm (Mini DFA)

Minimizing deterministic finite automata

a- Remove unreachable state -> cannot reach it using any input.

b- Merge equivalent state.

Partition	Set	Input	Action
$\{\text{final group}\} \{\text{non final group}\}$	$\{\text{non final group}\}$	Input test	Split $\{\}$
$\{\text{final group}\} \{\text{non final group}\}$	$\{\text{non final group}\}$	All inputs	None

c- Remove dead state -> not final state that take input and it not have output.

4- DFA -> Code (code, FLex)

1- a. using doubly nested case analysis

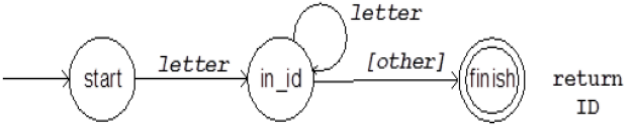
```
state = start;
getchar(input);

while(state != finish || state != error)
{
    switch (state)
    {
        case start:
            if( isalpha(input) )
            {
                advance(input);
                state = in_id;
            }
            else state = error;
            break;

        case in_id:
            if( isalpha(input) ) advance(input);
            else state =
            break;

        default: break;
    }
}

if( state == finish ) return ID;
else return error;
```



The diagram illustrates a Deterministic Finite Automaton (DFA) for recognizing identifiers. It consists of three states: 'start', 'in_id', and 'finish'. The 'start' state is the initial state, indicated by an incoming arrow. A transition labeled 'letter' leads from 'start' to 'in_id'. The 'in_id' state has a self-loop labeled 'letter', indicating that any letter keeps the automaton in the 'in_id' state. A transition labeled '[other]' leads from 'in_id' to 'finish'. The 'finish' state is the final state, represented by a double circle. An arrow points from 'finish' to the text 'return ID', indicating the output of the DFA when it reaches the final state.

1- b. using transition table

```
i = 0;
state = 0;
while ( input[i] )
{
    state = DFA[ state , input[i++] ];
}
```

2- Fast lex (Flex)

RegExp (lex code) -> flex -> Lex.yy.c -> c compiler -> Lex.yy.exe

lex code:

1- global variables
%{

int index = 0;
%}

2- define variables
id [a-zA-Z0-9_]+
digit [0-9]

3- run commands
%%

```
{id}      { print(" <%s , id> " , yylex); }  
{digit} { print(" <%s , digit> " , yylex); }  
"=="     { print(" <== , equal> "); }
```

%%

4- functions
int yywrap()
{
 return 1;
}

```
int main()  
{  
    yylex();  
    return 0;  
}
```

Lecture 6,7

- Any language that requires unbounded counting cannot be represented by regular language. (token -> parser -> syntax tree)

- Parser (syntax analyzer) is responsible for syntax errors.

- Grammar called:

BNF (Backus-Naur Form) or EBNF (Extended Backus-Naur Form).

- Context Free Grammars (CFG) using in parser define as:

$$G = (N, T, S, P)$$

N-> non terminal (symbol have rule)

T-> terminal (symbol have no rule)

S-> start symbol

P-> set of Production rules $\{ \alpha \rightarrow \beta \mid \alpha \in N \wedge \beta \in (N \cup T)^* \}$

- Derivation have two types (Left Most Derivation, Right Most Derivation)

LMD -> expand the leftmost nonterminal in the production, it's using with (Top-Down parser)

RMD -> expand the rightmost nonterminal in the production, it's using with (Bottom-Up parser)

- Trees have two types (parse tree (All details), syntax tree (less details))

ex:

E-> E+E | E*E | num check exp-> 1+2*3

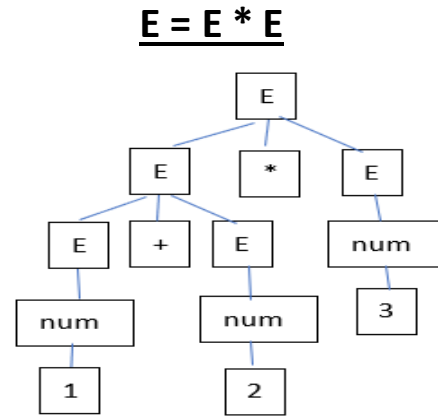
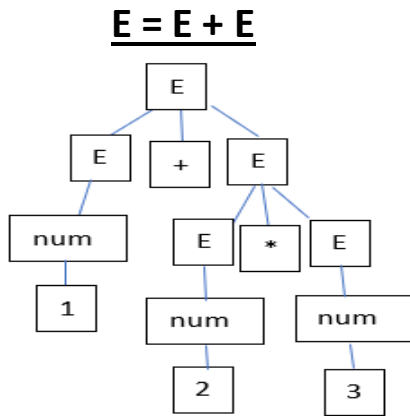
Solve:

LMD

RMD

E = E + E	E = E * E	E = E + E	E = E * E
E = num + E	E = E + E * E	E = E + E * E	E = E * num
E = 1 + E	E = num + E * E	E = E + E * num	E = E * 3
E = 1 + E * E	E = 1 + E * E	E = E + E * 3	E = E + E * 3
E = 1 + num * E	E = 1 + num * E	E = E + num * 3	E = E + num * 3
E = 1 + 2 * E	E = 1 + 2 * E	E = E + 2 * 3	E = E + 2 * 3
E = 1 + 2 * num	E = 1 + 2 * num	E = num + 2 * 3	E = num + 2 * 3
E = 1 + 2 * 3	E = 1 + 2 * 3	E = 1 + 2 * 3	E = 1 + 2 * 3

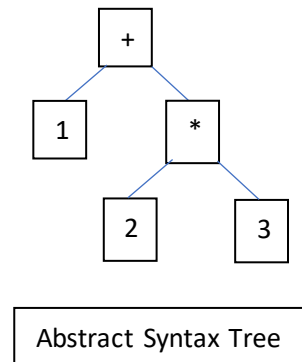
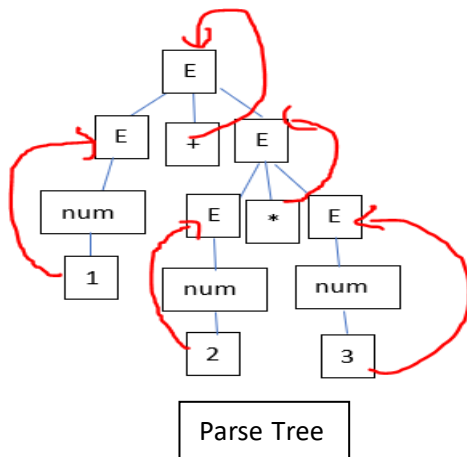
- NOTE: Trees different by starting rule not by LMD and RMD



Convert parse tree to syntax tree

- Every node that have single child (junk nodes), its child rewrite instead.

Ex:



Resolve ambiguous rules:

- Every operation has same precedence must have non terminal.
- Concatenate each rule with next.
- Low precedence come first.
- Remove left recursion or right recursion according to operation associativity (+, -, /, * from left to right) (power from right to left).

Ex: (Correct Grammar)

E -> **E** + **T** | **E** - **T** | **T**

T -> **T** * **P** | **T** / **P** | **P**

P -> **N** ^ **P** | **N**

N -> num

Note: to show if this grammar is ambiguous, it must draw two different parse trees for the same string that user input using the same grammar.

Lecture 8

Types of parse

- Top – down parser:
 - Predictive parsing (LL1 -> Left lookahead 1 token)

LL1 problems

left factoring: $S \rightarrow SAB|z$ (temp solution: ignore)

left recursion: $S \rightarrow +A|+B|-C$ (Not allowed)

- Bottom – up parser:

First calculations: (calculate first from down to top)

- First of any terminal is itself.
- All first terminal from each exp in each |.
- If there is non-terminal take its first.
- If any first of non-terminal have epsilon then take first of next one.
- If reached to last non-terminal and its first has epsilon then epsilon will be in result.
- Duplication not allowed.

Follow calculations: (calculation of terminal and non-terminal is the same)

- Follow of start symbol must start with \$.
- Find follow for each non-terminal by finding it in all expressions and take first of next symbol after it.
- If first of next symbol has ϵ then take first of next symbol.
- If there is no next take follow of left non-terminal.
- Duplication and epsilon are not allowed.

Ex:

$S \rightarrow ABC|z|a|b$

$A \rightarrow a|\epsilon$

$B \rightarrow b|\epsilon$

$C \rightarrow cBx|\epsilon$

Solu:

Follow $S = \{\$ \}$

Follow $A = \{\text{first } B\} = \{b, \text{first } C\} = \{b, c, \text{follow } S\} = \{b, c, \$ \}$

Follow $B = \{\text{first } C, x\} = \{c, \text{follow } S, x\} = \{c, \$, x\}$

Follow $C = \{\text{follow } S\} = \{\$ \}$

Parsing table

- For each first of non-terminal intersection between this non-terminal and every terminal in first, write rule in table related to terminal until |
- If terminal have ϵ write ϵ rule at intersect of follow terminals and non-terminal

○ Ex:

Rules: $E \rightarrow TX$ $X \rightarrow +E | \epsilon$ $T \rightarrow (E) | int Y$ $Y \rightarrow *T | \epsilon$

- First $E = \{ (, int \}$
- First $X = \{ +, \epsilon \}$
- First $T = \{ (, int \}$
- First $Y = \{ *, \epsilon \}$
- Follow $E = \{ \$,) \}$
- Follow $X = \{ \$,) \}$
- Follow $T = \{ +, \$,) \}$
- Follow $Y = \{ +, \$,) \}$

Non terminal		Terminal					
		Int	*	+	()	\$
	E	$E \rightarrow TX$			$E \rightarrow TX$		
	X			$X \rightarrow +E$		$X \rightarrow \epsilon$	$X \rightarrow \epsilon$
	T	$T \rightarrow int Y$			$T \rightarrow (E)$		
	Y		$Y \rightarrow *T$	$Y \rightarrow \epsilon$		$Y \rightarrow \epsilon$	$Y \rightarrow \epsilon$

Build Stack from parsing table

- Put \$ in stack and start symbol in stack column.
- Put input in input column.
- Take most right symbol in stack column and most left symbol in input column and get the rule in intersection between them from parsing table and make replace in reverse order of symbol in stack column with this rule (write this replace in action column).
- If most right symbol in stack column and most left symbol in input column are the same match them with each other and take pop action.
- Do that again until stack empty (Accept) or get error (Reject).
- If there is a symbol still in stack column this mean that user have not been write this symbol (like ;), and If there is a symbol still in input column this mean that user have been write this symbol without needing it.
- Note:
 - grammar in stack and input in queue.
 - Multiple entry in (LL1) parsing table not allowed (more one rule in the same cell in table)

Ex: match this input $\text{int}^*\text{int}\$$

Stack	Input	Action
$\$E$	$\text{int}^*\text{int}\$$	Replace $E \rightarrow TX$
$\$XT$	$\text{int}^*\text{int}\$$	Replace $T \rightarrow \text{int } Y$
$\$XY \text{ int}$	$\text{int}^*\text{int}\$$	Match (pop int)
$\$XY$	$^*\text{int}\$$	Replace $Y \rightarrow ^*T$
$\$XT^*$	$^*\text{int}\$$	Match (pop *)
$\$XT$	$\text{int}\$$	Replace $T \rightarrow \text{int } Y$
$\$XY \text{ int}$	$\text{int}\$$	Match (pop int)
$\$XY$	$\$$	Replace $Y \rightarrow \epsilon$
$\$X$	$\$$	Replace $X \rightarrow \epsilon$
$\$$	$\$$	Accept