# UNIT-1 Finite Automata & Regular Expression

**1. Regular Expression for valid identifies**

Let $\Sigma_A$ be the set of alphabets & $\Sigma_A$ be the set of digits. The regular expression $R$ for all valid identifies (alphabet followed by any sequence of alphabets or digits) is

$$R = (\Sigma_A)(\Sigma_A \cup \Sigma_D)^*$$

The keywords (for, while, if) are excluded in the eemical analysis phase following taken recognition

**2. Design a DFA equivalent to R**

The DFA $M = \{Q, \Sigma, \delta, q_0, F\}$

$Q = (q_0, q_m, q_F)$

$q_0 = $ Start state

$q_1 = $ accepting state (valid identifier started)
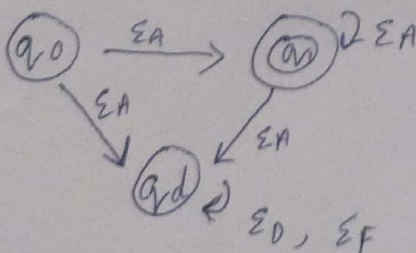
$q_F = $ dead state (invalid start)

$\Sigma = \Sigma_A \cup \Sigma_D$

$A = \{q\}$ (final state)

**Transition Table**

| State | Input $\Sigma_A$ | Input $\Sigma_D$ | Input (other) |
|---|---|---|---|
| $q_0$ | $q_1$ | $q_d$ | $q_d$ |
| $q_1$ | $q_1$ | $q_d$ | $q_1$ (loop) |
| $q_d$ | $q_1$ | $q_d$ | $q_d$ |

**DFA (diagram)**

3> Embedding the DFA is a lexical Analysis

The DFA acts as the state machine for recognizing the part of an identifier

1) DFA recognition : The lexer consumes input characters, tracking DFA's state. when the input stream forces the DFA out of a (encountering a space on operator) the operator reader to that point is identified as a potential token.

2) keyword check : The recognized string is then checked against a small, finite list of reserved keywords (for, while, if).
This is typically done via a fast hash table lookup

3) Token Generation
• If the string is found in the keyword list, a KEYWORD token is generated
• Otherwise, an IDENTIFIER token is generated & its entry (lexeme & type) is stored in the symbol table.

Q2 UNIT2 PDA & content free language.

1) Formulate a CFG for well formed queries
let $O =$ '<open>' & $c =$ '</close>'. The grammar $G$ models balanced nesting

$$S \rightarrow OSC \mid SS \mid \epsilon$$

$S \rightarrow OSC$ handles nested structures (eg <open> ... </close>)

$S \rightarrow SS$ handles concatenated structures (eg <open> </close> <open>

$S \rightarrow \epsilon$ handles empty query

2) Construct a PDA that accept such queries
The PDA accepts the languages by empty stack. It uses the stack to track unmatched <open> tags
$$M = (\{q_0\}, \{O, C\}, \{z_0, x\}, \delta, q_0, z_0)$$

| | Input | Top of stack | New state | stack operation | Rationale |
|---|---|---|---|---|---|
| $q_0$ | O | $Z_0$ | $q_0$ | $x = 0$ | Push x for first O |
| $q_0$ | O | x | $q_0$ | $xx$ | push x for nested O |
| $q_0$ | C | x | $q_0$ | ε | Pop x for matching C |
| $q_0$ | ε | $Z_0$ | $q_0$ | ε | Accept by empty stack |

5> Demonstrate the parse tree
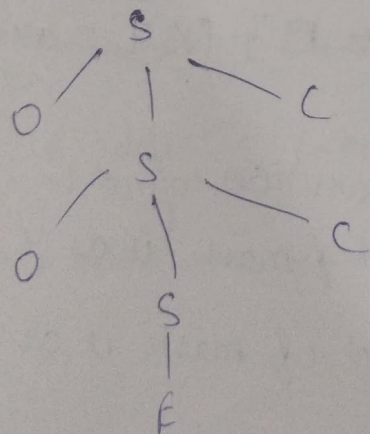
Query : <open> <open> </ulore > </close>

Derivation

$S \Rightarrow OSC$

$S \Rightarrow O(OSC) C \quad (S \rightarrow OSC)$

$S \rightarrow O(OEC) C (S \rightarrow E)$

  parse tree :

Q3 UNIT - 3 Turing Machine & Chomsky hierarchy

1) Justify why $l = \{a^n b^n c^n / n \geq 1\}$ is not content free
   we use pumping lemma for CFLs.

2) Choose string s: let $p$ be the pumping length choose $s = a^p b^p c^p$

3) Decompose S: $S = uvwxy$ where $|vwx| \leq p$ & $|vx| \geq 1$

4) Pumping argument since $|vwx| \leq p$ the pumpable segments
   $vwx$ can only contain symbols from almost two blocks
   (only a's & b's, or only b's & c's)

   • case (vwx is in a's & b's) pumping up (setting $i = 2$) ↑ the
   no of a's and/or b's, but leaves the number of c's flexed
   atp.

   • Resulting string $s' = u^2 w x^2 y$ has unequal counts of a's b's
   & c's (specifically, count(a) + count(b) > count(c))

5) Conclusion $s' \notin l$, since the condition of the pumping lemma
   are violated, $l$ can't be a CFL

2) Design a Turing Machine (TM) that accepts $l$
   The TM mark one a, one b & one c in the cycle all symbols are.
   marked.

   • Tape alphabets $f: \{a, b, c, x, y, z, \square\}$ (x, y, z are markers)

     core logic
     $q_0$: Mark the leftmost a as x & transition to find b
     $q_a$: find the leftmost unmarked b & mask it as y, then transition
     to find c
     $q_b$: find the leftmost unmarked c & mask it as z, then transition
     to return
     $q_{ret}$: soon tell to the starting point (x)
     $q_{check}$: After all a's are marked, scan right to ensure the
   rest of the tape is first 4b, 2's & finally $\square$ 'B'

Step by step configuration for 'aaabbbccc'

The TM cycles thrice times to mark the three pairs

•) Cycle 1 (mark $a_1 b_1 c_1$): $q_0, aaabbbccc \rightarrow q_a, Xaabbccc$ (Mark $a$)
$\rightarrow q_b, Xaabbbccc$ (Mark $b$) $\rightarrow q_{ner} \rightarrow Xaa Ybb Zcc$ (Mark $c$,
return left) $\rightarrow q_0, XaaYbb Zcc$ (Restart)

Cycle 2 (Mark $a_2 b_2 c_2$): $q_0, Xaabccc \rightarrow q_{ner}, XXaYYY Yzzz$

Final check ($q_{check}$): $q_0$ reads the marked $a$'s ($x$'s), transition
to $q_{check}$. $q_{check}$ scans over $y$'s & $z$'s until it hits $\square$ $q_{check}$;
$XX xyYYYZZZ\square \rightarrow q_{accept}$ (Accept)

UNIT-4  Code generation & optimization

Expression $(A+B) * (C-D) + E$

Syntax - Directed translation scheme (s attributed)
using a simple procedure - based grammar

| Production | Semantic Rules |
|---|---|
| $\rightarrow E_1 + T$ | $E_0.addr = new\_temp()$ ; $E_{code} = E_1.code \| T.code \|$ |
| | $E_0.addr = E_1.addr + T.addr$ |
| $\rightarrow T_1 * F$ | $T.addr = new\_temp()$ ; $T.code = T_1.code \| F.code \|$ |
| | $T.addr = T_1.addr * F.addr$ |
| $\rightarrow (E_1)$ | $F.addr = E_1.addr$ ; $F.code = E_1.code$ |
| $\rightarrow id$ | $F.addr = id.lexeme$ ; $F.code = E$ |

2) Generate three address code (TAC)

The TAC is generated based on expression's evaluation order procedure: parenthesis → multiplication → addition

1) $t_1 = A + B$

2) $t_2 = C - D$

3) $t_3 = t_1 * t_2$

4) $t_4 = t_3 + E$

3) Optimize the generated TAC

There is no duplicate expression (common subexpressions) in lines 1 & 2. The code is already optimal wrt to CSE.

Dead Code Removal

Assume the final result $t_4$ is used, all intermediate variables $(t_1, t_2, t_3)$ are nescessary inputs for subsequent lines. No dead code can be removed.

Optimized TAC (unchanged)

1) $t_1 = A + B$

2) $t_2 = C - D$

3) $t_3 = t_1 * t_2$

4) $t_4 = t_3 + E$

Q3 (Commutative

Language L = Equal no of 0's & 1's & no prefix has more 1's than 0's (Dyck Paths)

i) Prove that L is content free but not regular

• Not Regular : Use the pumping lemma for regular languages. Choose $s = 0^p 1^p$ Pumping down (i=0) gives $0^{p-v} 1^p$, $P(v \geq 1)$ which has unequal counts voilating L. Thus L is not Regular

• Content free : The language L is accepted by a pushdown Automaton (PDA) shown below) which demonstrates its content free nature. The PDA's stack is essential for counting and comparing the non local sequ dependencies ($0's \ vs \ 1's$)

Provide a CFG for his language (1)

The Grammar G must enforce that every 1 is matched by a preceding 0

$$S \rightarrow 0S1 \mid S \mid \epsilon$$

This CFG generate all valid cycle paths

3) Design a PDA & trace '0011'

A) PDA Design (M)

M accepts by empty stack, using x to count the excess no. of 0's

start $z_0$ : $\delta(q_0, 0, z_0) = \{(q_0, x z_0)\}$

Push 0 : $\delta(q_0, 0, x) = \{(q_0, x x)\}$

Pop 1 (prefix check) : $\delta(q_0, 1, x) = \{(q_0, \epsilon)\}$ (pops only if 0's are in excess.

B) Trace the acceptance of '0011'

| Input | State | Stack (1→R) | Transition | Condition check (0's ≥ 1's) |
|-------|-------|-------------|------------|------------------------------|
| 0011 | $q_0$ | $z_0$ | push 0 | |
| 0011 | $q_0$ | $x z_0$ | push 0 | 2 ≥ 0 |
| 0011 | $q_0$ | $x x z_0$ | pop 1 | 2 ≥ 1 |
| 0011 | $q_0$ | $x z_0$ | pop 1 | 2 ≥ 2 |
| 0011 $\epsilon$ | $q_0$ | $z_0$ | empty stack | 2 = 2 |
| | | $\epsilon$ | accept | |