

Automata Design Document

Muhammad Shaffan Ahmad 23i-0673

Zaid Bin Umer 23i-0671

1. Overview

The Custom Language Scanner uses a Deterministic Finite Automaton (DFA) generated by JFlex to recognize tokens. This document describes the regular expressions, individual NFAs, and the final optimized DFA.

2. Regular Expressions (Lexical Specification)

The following regular expressions define the token classes for the language:

Token Class	Regular Expression	Description
Single Line Comment	<code>##[^\n]*</code>	Starts with ##, followed by non-newline chars
Boolean Literal	<code>(true false)</code>	
Identifier	<code>[A-Z][a-z0-9]{0,30}</code>	Uppercase start, mixed case tail (Max 31 chars)
Integer Literal	<code>[+-]?[0-9]+</code>	Optional sign, followed by digits
Floating Point	<code>[+-]?[0-9]+\.[0-9]{1,6}([eE][+-]?[0-9]+)?</code>	Decimal required, optional exponent. Max 6 decimal digits.
Operator	<code>[+\-*/%<>=!]</code>	Single character operators
Punctuator	<code>[(){}[\],;:]</code>	Braces, parentheses, and delimiters

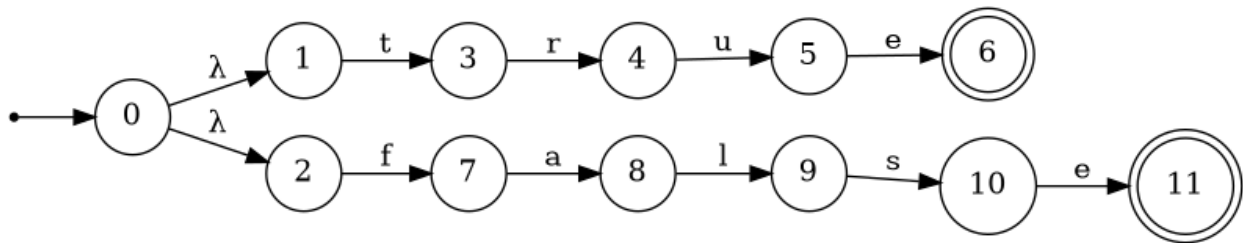
3. Individual NFA Transition Tables

These tables represent the Non-Deterministic Finite Automata for each individual token class before combination.

3.1 Boolean Literal NFA

Regex: true | false

Accept States: [6, 11]



State Input Next State

q0 λ q1, q2

q1 t q3

q2 f q7

q3 r q4

q4 u q5

q5 e q6 (Accept)

q7 a q8

q8 l q9

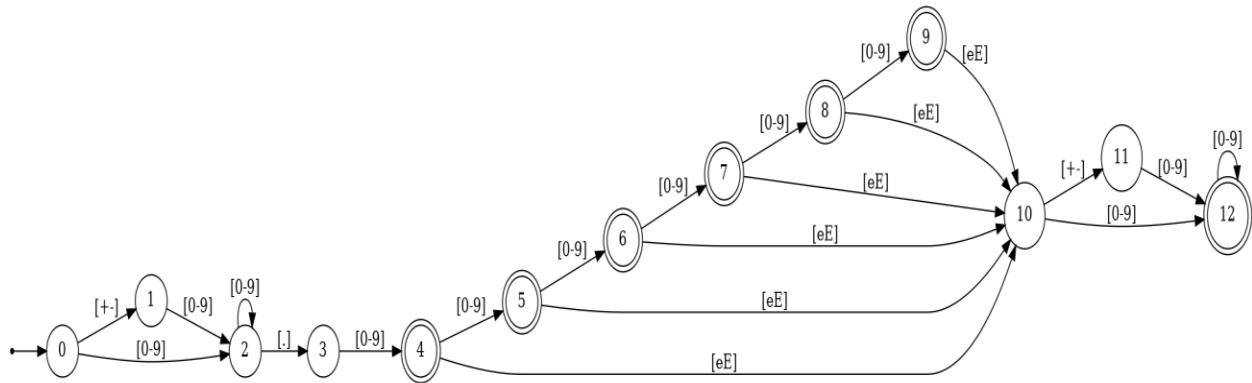
q9 s q10

q10 e q11 (Accept)

3.2 Floating Point Literal NFA

Regex: [+]?[0-9]+\.[0-9]+([eE][+]?[0-9]+)?

Accept States: [4, 5, 6, 7, 8, 9, 12] (Enforcing precision limits)



State Input Next State

q0 [+ -] q1

q0 [0-9] q2

q1 [0-9] q2

q2 [0-9] q2

q2 . q3

q3 [0-9] q4

q4 [0-9] q5

q4 [eE] q10

q5 [0-9] q6

q5 [eE] q10

... (states q6-q9 enforce max 6 decimal digits) ...

q10 [+ -] q11

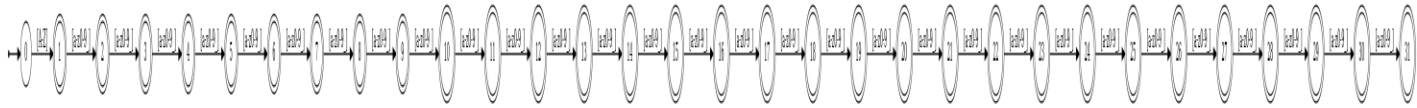
q10 [0-9] q12

q12 [0-9] q12 (Accept)

3.3 Identifier NFA

Regex: [A-Z][a-z0-9_]* (Max 31 chars)

Accept States: [1..31]



State Input Next State

q0 [A-Z] q1

q1 [a-z0-9_] q2

q2 [a-z0-9_] q3

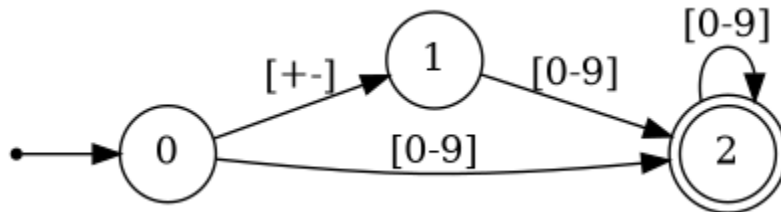
...

q30 [a-z0-9_] q31 (Accept, max length reached)

3.4 Integer Literal NFA

Regex: `[+-]?[0-9]+`

Accept States: [2]



State Input Next State

q0 [+] q1

q0 [0-9] q2

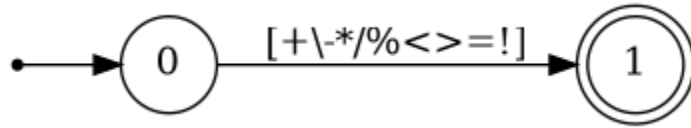
q1 [0-9] q2

q2 [0-9] q2 (Accept)

3.5 Single Char Operator NFA

Regex: `[+\-*/%<>=!]`

Accept States: [1]



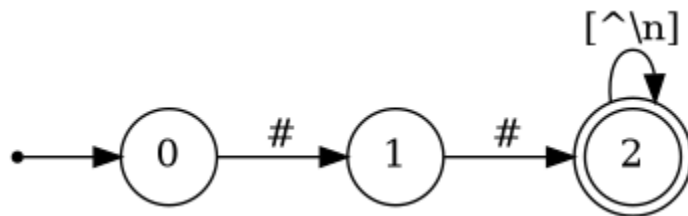
State	Input	Next State
-------	-------	------------

q0	[+\\-*/%<>=!]	q1 (Accept)
----	---------------	-------------

3.6 Single Line Comment NFA

Regex: ##[^\\n]*

Accept States: [2]



State	Input	Next State
-------	-------	------------

q0	#	q1
----	---	----

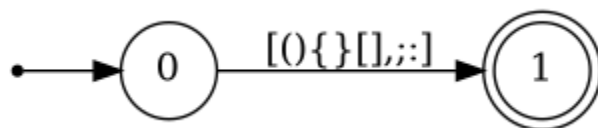
q1	#	q2
----	---	----

q2	[^\\n]	q2 (Accept loop)
----	--------	------------------

3.7 Punctuator NFA

Regex: [(){}\\[\\],:;]

Accept States: [1]



State	Input	Next State
-------	-------	------------

-------	--	--

q0	[((){\[\],,:]	q1 (Accept)
----	---------------	-------------

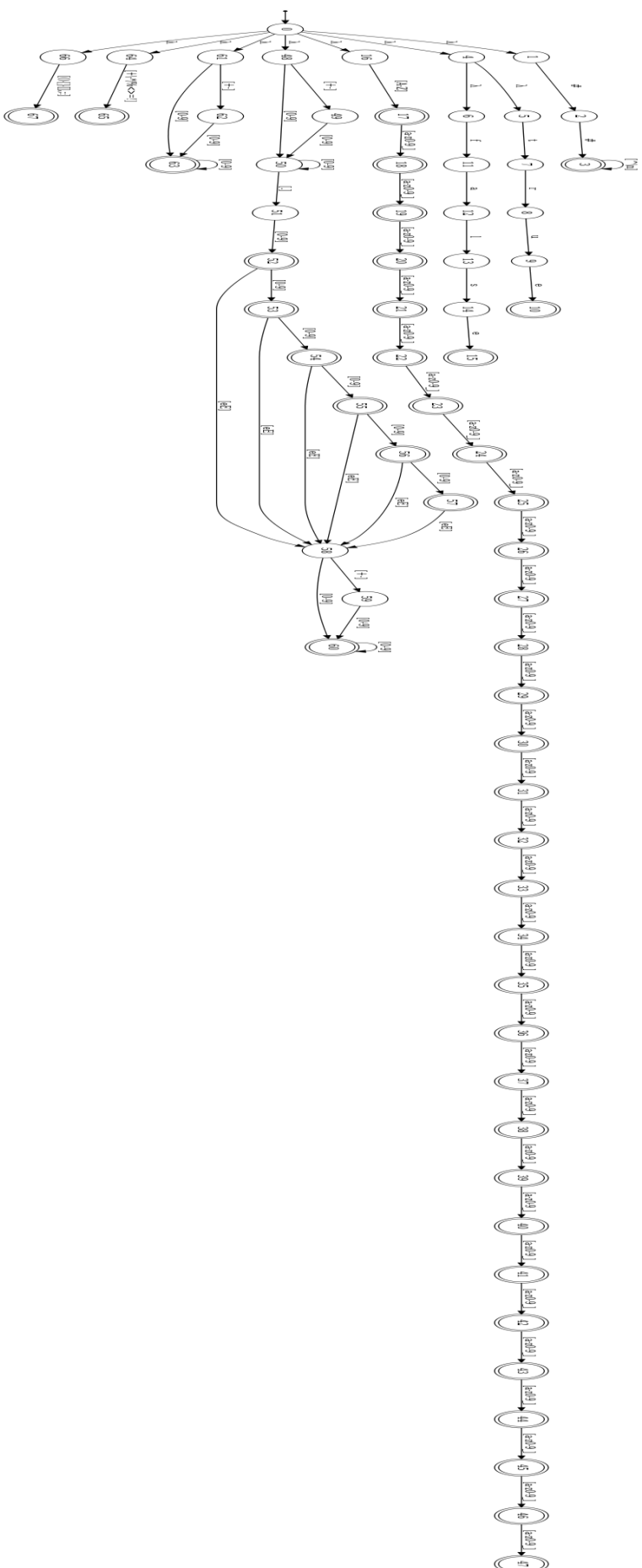
4. Combined NFA

The individual NFAs are combined into a single machine using epsilon (λ) transitions from a new start state.

Start State: 0

Accept States: [3, 10, 15, 17...56] (Union of all NFA accept states)

Priority Order: Comment > Boolean > Identifier > Float > Integer > Operator > Punctuator



State	Input	Next State
-------	-------	------------

q0	λ	q1 (Comment Start)
----	-----------	--------------------

q0	λ	q4 (Boolean Start)
----	-----------	--------------------

q0	λ	q16 (Identifier Start)
----	-----------	------------------------

q0	λ	q48 (Float Start)
----	-----------	-------------------

q0	λ	q61 (Integer Start)
----	-----------	---------------------

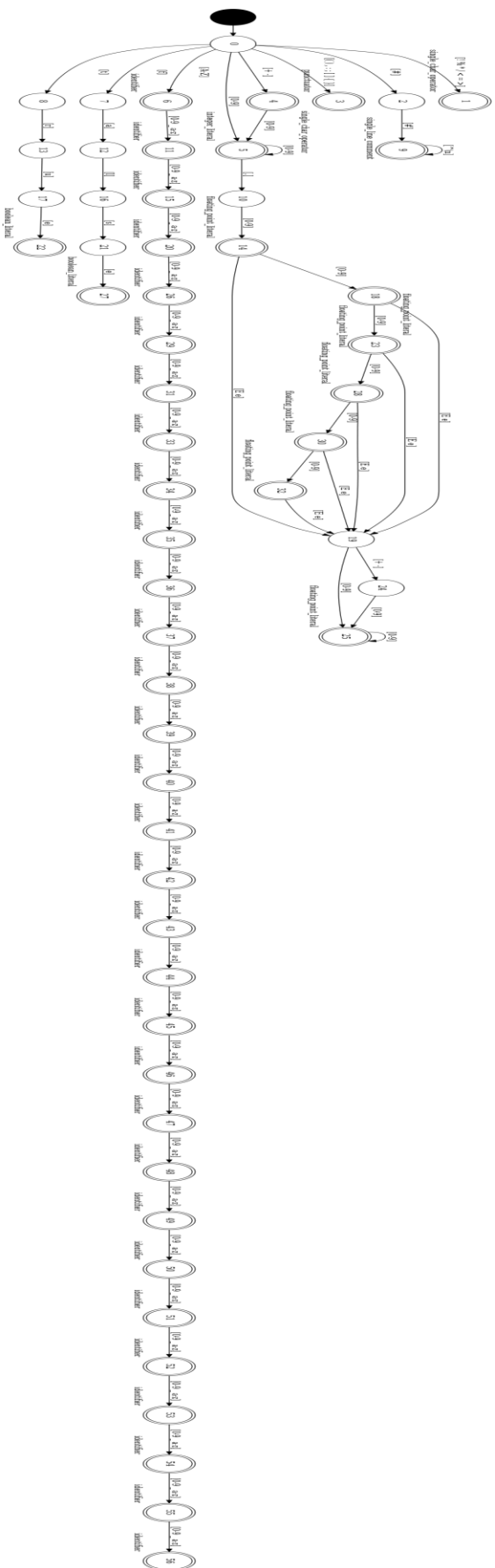
q0	λ	q64 (Operator Start)
----	-----------	----------------------

q0	λ	q66 (Punctuator Start)
----	-----------	------------------------

... (Internal transitions from individual NFAs) ...

5. Minimized DFA

After subset construction (NFA \rightarrow DFA)



and minimization, we achieve the optimal state machine.

Start State: 0

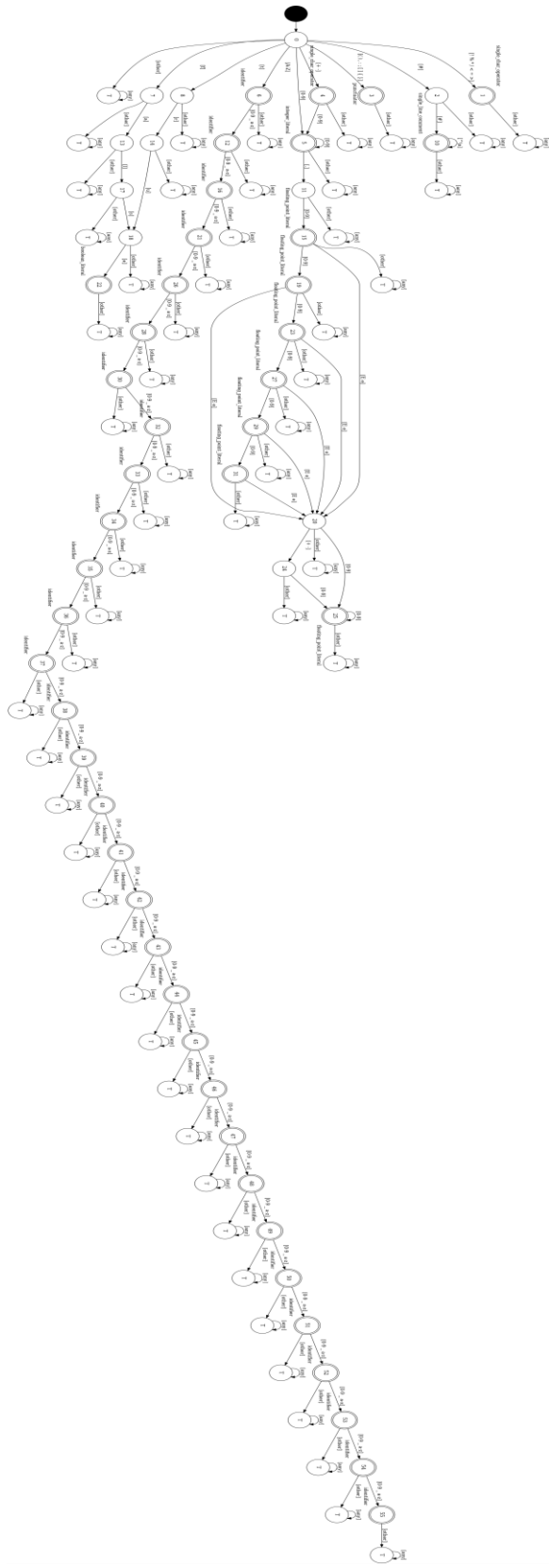
Total States: 56

Accept States: [1, 3, 4, 5, 6, 10, 12, 15, 16, 19, 21...55]

5.1 Acceptance State Labels

State	Token Recognized
q1	Single Char Operator
q3	Punctuator
q5	Integer Literal
q6, q12, q16...	Identifier
q10	Single Line Comment
q15, q19, q23...	Floating Point Literal
q22	Boolean Literal (true or false)

5.2 Transition Table (Min DFA)



State	Input	Next State
-------	-------	------------

q0	[!%*/<>=]	q1
----	-----------	----

q0	[#]	q2
----	-----	----

q0	[(){}[],:;]	q3
----	-------------	----

q0	[+-]	q4
----	------	----

q0	[0-9]	q5
----	-------	----

q0	[A-Z]	q6
----	-------	----

q0	[f]	q7
----	-----	----

q0	[t]	q8
----	-----	----

q0	[other]	q9 (Dead State)
----	---------	-----------------

q2	[#]	q10 (Comment Start)
----	-----	---------------------

q4	[0-9]	q5 (Signed Int)
----	-------	-----------------

q5	[.]	q11 (Float Start)
----	-----	-------------------

q6	[a-z0-9_]	q12 (Identifier)
----	-----------	------------------

q7	[a]	q13 (Start of 'false')
----	-----	------------------------

q8	[r]	q14 (Start of 'true')
----	-----	-----------------------

q13	[l]	q17
-----	-----	-----

q14	[u]	q18
-----	-----	-----

q17	[s]	q18 (Merge 'false'/'true' suffix)
-----	-----	-----------------------------------

q18	[e]	q22 (Boolean Accept)
-----	-----	----------------------

q11	[0-9]	q15 (Float digit 1)
-----	-------	---------------------

q15	[0-9]	q19 (Float digit 2)
-----	-------	---------------------

...

q31 [0-9] q9 (Precision limit exceeded -> Error)

6. Optimization Summary

1. **Boolean Literal Merging:** The suffixes of true and false were effectively merged by the minimization algorithm. State q18 handles the e transition for both.
2. **Precision Enforcement:** The DFA structure (states q11 → q31) explicitly enforces the 6-digit limit for floating-point numbers. Any digit after the 6th transitions to the dead state (q9) or implies the end of the token.
3. **Ambiguity Resolution:**
 - + leads to q4 (Operator Accept).
 - + followed by [0-9] leads to q5 (Integer Accept).
 - This structure enforces the "Longest Match" rule natively.

For Clarity Purposes, pictures of Identifier's NFA, Combined NFA, Combined DFA, and Minimized DFA are attached at the end:



