

YAFSM Language Specification

Compilers Project

Group 6

Authors

Vishal Vijay Devadiga (CS21BTECH11061)

Satpute Aniket Tukaram (CS21BTECH11056)

Mahin Bansal (CS21BTECH11034)

Harshit Pant (CS21BTECH11021)



भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

Table of Contents

- Table of Contents
- Introduction
 - What is YAFSM?
 - Why YAFSM?
- Language Specifications
 - Data Types
 - * Primitive Data Types
 - Integer:
 - Character:
 - Float:
 - Boolean:
 - * Composite Data Types
 - Strings:
 - Finite-Sets:
 - Structs:
 - Regular Expressions:
 - Context Free Grammars:
 - DFAs:
 - NFAs:
 - PDAs:
 - * Comments
 - Operators
 - * Arithmetic Operators
 - * Logical Operators
 - * Comparison Operators
 - * Assignment Operators
 - * Set Operators
 - * Automaton Operators
 - * Misc Operators
 - * Operator Precedence
 - Set Operators Precedence
 - Automaton Operators Precedence:
 - Control Flow
 - * If-Else
 - * Loops
 - Constants

- Keywords
- Identifiers
- Statements
 - * Declaration Statement
 - * Assignment Statement
 - * Function Declaration Statement
 - * Function Call Statement
 - * IO Statements
- In-built Functions
 - * Set Functions
 - * CFG Functions
 - * DFA Functions
 - * NFA Functions
 - * PDA Functions
 - * Regular Expression Functions
- References

Introduction

What is YAFSM?

YAFSM is a domain specific language that simplifies working with Finite State Machines(FSMs). Finite State Machines include Deterministic Finite Automata(DFAs), Non-Deterministic Finite Automata(NFAs), Pushdown Automata(PDAs).

It supports the following features:

- Defining Finite State Machines
- Regular Expressions
- Context Free Grammars

Why YAFSM?

Finite State Machines are used in many applications, such as:

- Regular Expressions
- Lexical Analysis
- Compilers
- Network Protocols
- Digital Logic
- Artificial Intelligence
- Natural Language Processing
- etc.

Finite State Machines are used in many applications, but the syntax for defining a Finite State Machine is not very intuitive. YAFSM aims to simplify the syntax for defining a Finite State Machine, making it easier for programmers to work with Finite State Machines.

Language Specifications

YAFSM follows, making it easier for programmers to pick up YAFSM easily and keep their focus on the logic rather than YAFSM.

- YAFSM is a **statically typed** language
- YAFSM is a **strongly typed** language
- YAFSM is a **procedural** language
- YAFSM is case sensitive.

YAFSM does not support Object Oriented Programming(**OOPs**).

Data Types

YAFSM uses common data types found in most programming languages.

Primitive Data Types

Integer: Signed Integers are represented by the `int_x` keyword, where `x` is the number of bits used to represent the integer. YAFSM supports 8, 16, 32 and 64 bit integers.

```
int_8 a;  
int_8 b = 10;  
int_16 c = 20;  
int_32 d = 30;  
int_64 e = 40;
```

Unsigned Integers are represented by the `uint_x` keyword, where `x` is the number of bits used to represent the integer. YAFSM supports 8, 16, 32 and 64 bit integers.

```
uint_8 a;  
uint_8 b = 10;  
uint_16 c = 20;  
uint_32 d = 30;  
uint_64 e = 40;
```

Character: Characters are represented by the `char` keyword. YAFSM supports 8 bit characters.

```
char a;  
char b = 'a';
```

Float: Floats are represented by the `float_x` keyword, where `x` is the number of bits used to represent the float. YAFSM supports 32 and 64 bit floats.

```
float_32 a;  
float_32 b = 10.5;  
float_64 c = 20.5;
```

Boolean: Booleans are represented by the `bool` keyword, which is similar to the `bool` keyword in C, C++, Java and Python.

```
bool a;  
bool b = true;  
bool c = false;
```

Composite Data Types

Strings: Strings are represented by the `string` keyword. Strings are immutable, and can be indexed using the `[]` operator.

```
string temp;  
string a = "Hello World";  
char b = a[0];
```

Finite-Sets: Sets are collections of elements of the same data type. YAFSM supports two types of sets: Ordered Sets and Unordered Sets.

- Ordered Sets are represented by the `o_set` keyword.
- Unordered Sets are represented by the `u_set` keyword.

```
o_set<int_8> a;  
o_set<int_8> b = {1, 2, 3};
```

Structs: Structs are represented by the `struct` keyword. Structs can contain any data type supported by YAFSM.

```
struct Point {  
    int_8 x;  
    int_8 y;  
    float_32 z;  
    string name;  
}
```

Additionally structs can contain other data types which have already been defined as structs

```
Point a,b;

struct set_of_coords {
    o_set<Point> = {a,b};
}
```

Regular Expressions: Regular Expressions are represented by the `regex` keyword. Regular Expression can contain definitions of other Regular Expressions, and can be used to define Finite State Machines.

```
regex alphabet = r'[a-z]';
regex Letter = r'{alphabet}';
regex Digit = r'[0-9]';
regex a = r'[ab]{2}';
regex b = r'{a}*';
regex c = r'{a}+';
regex d = r'{a}?';
regex e = r'^[ab]';
regex f = r'[ab]$';
regex g = r'{a}|{b}';
```

Context Free Grammars: Context Free Grammars are represented by the `cfg` keyword. Context Free Grammars are defined by a 4-tuple:

$$(N, \Sigma, P, S)$$

where:

- N is a set of non-terminal symbols
- Σ is a set of terminal symbols
- P is a set of production rules
- S is the start symbol

A production rule is represented as:

- $A \rightarrow \alpha$ where α is a string of terminal and non-terminal symbols
- $A \rightarrow \{ \alpha_1, \alpha_2, \dots \}$ where $\alpha_1, \alpha_2, \dots$ are strings of terminal and non-terminal symbols.

```

cfg a;
a.T = {a:"a", demo:"b"};
a.N = {A, B};
a.S = A;
a.P = {
    A -> <a>B,
    B -> <demo>A,
    A -> {<a>A, \e}
};

```

DFAs: DFAs are represented by the `dfa` keyword.

A DFA is defined by a 5-tuple:

$$(Q, \Sigma, \delta, q_0, F)$$

where:

- Q is a `o_set` of states
- Σ is a `set` of input symbols
- δ is the transition function, which maps $Q \times \Sigma$ to Q
- q_0 is the initial state
- F is a set of final states

A transition can be represented as:

state1, input_symbol -> state2

In case of multiple transitions from the same state on different input symbols to the same state, the transitions can be represented as:

state1, {input_symbol1, input_symbol2, ...} -> state2

This can also be done as:

state1, <regex> -> state2

state1, <set> -> state2

δ is either a set of such transitions or it can be represented as a matrix of size $|Q| \times |\Sigma|$, where each element of the matrix is a state.

```

dfa a;
a.Q = {q0, q1, q2};
a.Sigma = {0,1,2};

```



```

a.delta = {
    q0, 0 -> q1,
    a.Q[0] , 1 -> a.Q[1],
    q1 , r'[12]' -> q2,
    q2, {0,2} -> q1
};
a.q0 = q0;
a.F = {q1,q2};

```

NFAs: NFAs are represented by the `nfa` keyword.

A NFA is defined by a 5-tuple:

$$(Q, \Sigma, \delta, q_0, F)$$

where:

- Q is a `o_set` of states
- Σ is a `set` of input symbols
- δ is the transition function, which maps $Q \times \Sigma$ to 2^Q
- q_0 is the initial state
- F is a set of final states

Here 2^Q represent the power set of Q .

A transition can be represented as:

- `state1, input_symbol -> {state2, state3, ...}`
- `state1, {input_symbol1, input_symbol2, ...} -> {state2, state3, ...}`
- `state1, <regex> -> {state2, state3, ...}`
- `state1, <set> -> {state2, state3, ...}`

Here `input_symbols` can include ϵ which is represented by `'\e'`.

```

nfa a;
a.Q = {q0, q1, q2};
a.Sigma = {0,1,2};
a.delta = {
    q0, 0 -> {q1, q2},
    a.Q[0] , 1 -> a.Q[1],
    q1 , r'[12]' -> q2,
    q2, {0,2} -> q1,
    q2, \e -> q0
}

```

```
};
a.q0 = q0;
a.F = {q1,q2};
```

PDAs: PDAs are represented by the `pda` keyword.

A PDA is defined by a 6-tuple:

$$(Q, \Sigma, \Gamma, \delta, q_0, F)$$

where:

- Q is a `o_set` of states
- Σ is a `set` of input symbols
- Γ is a `set` of stack symbols
- δ is the transition function, which maps $Q \times \Sigma_{\epsilon} \times \Gamma_{\epsilon}$ to $2^{Q \times \Gamma_{\epsilon}}$
- q_0 is the initial state
- F is a set of final states

Here $2^{Q \times \Gamma_{\epsilon}}$ represent the power set of $Q \times \Gamma_{\epsilon}$.

A transition can be represented as:

- `state1, input_symbol, stack_symbol -> state2, stack_symbol`
- `state1, {(input_symbol1, stack_symbol1), (input_symbol2, stack_symbol2), ...} -> state2, stack_symbol`
- `state1, input_symbol, stack_symbol -> {(state2, stack_symbol2), (state3, stack_symbol3), ...}`

```
pda a;
a.Q = {A, B, C};
a.S = {a, b, c};
a.G = {a, b, d};
a.delta = {
    A, a, a -> A, a,
    A, {(b, a), (b, c)} -> {(B, a), (A, d)},
    B, {(b, a), (b, c)} -> B, a,
    B, c, a -> C, \e,
    C, c, a -> {(C, \e), (C, a)}
};
a.q0 = A;
a.F = {C};
```

Comments

YAFSM has only one type of comment, that can act as both single line and multi line comments. The comment starts with `<!--` and ends with `--!>`. Below is an example of a comment:

```
<!-- This is a comment --!>
```

```
<!-- This is a  
multi line comment --!>
```

Operators

Operators supported by YAFSM are similar to the operators supported by C.

Arithmetic Operators

Operator	Description	Associativity
+	Addition	Left to Right
-	Subtraction	Left to Right
*	Multiplication	Left to Right
/	Division	Left to Right
%	Modulo	Left to Right

Logical Operators

Operator	Description	Associativity
&&	Logical AND	Left to Right
	Logical OR	Left to Right
!	Logical NOT	Right to Left

Comparison Operators

Operator	Description	Associativity
==	Equal to	Left to Right
!=	Not equal to	Left to Right
>	Greater than	Left to Right
<	Less than	Left to Right
>=	Greater than or equal to	Left to Right

Operator	Description	Associativity
\leq	Less than or equal to	Left to Right

Assignment Operators

Operator	Description	Associativity
$=$	Assignment	Right to Left
$+=$	Addition Assignment	Right to Left
$-=$	Subtraction Assignment	Right to Left
$*=$	Multiplication Assignment	Right to Left
$/=$	Division Assignment	Right to Left
$\%=$	Modulo Assignment	Right to Left
$\&=$	Logical AND Assignment	Right to Left
$ =$	Logical OR Assignment	Right to Left

Set Operators

Operator	Description	Associativity
$+$	Union	Left to Right
$-$	Difference	Left to Right
$*$	Intersection	Left to Right
2	Power Set	Left to Right

```
o_set<int_8> a = {1, 2, 3};
o_set<int_8> b = {2, 3, 4};
o_set<int_8> c = a + b; <!-- c = {1, 2, 3, 4} --!>
o_set<int_8> d = a - b; <!-- d = {1} --!>
o_set<int_8> e = a * b; <!-- e = {2, 3} --!>
o_set<o_set<int_8>> f = a^2; <!-- f = {{}, {1}, {2}, {3}, {1, 2},
{1, 3}, {2, 3}, {1, 2, 3}} --!>
```

Automaton Operators

Operator	Description	Associativity
$*$	Kleene	Left to Right

Operator	Description	Associativity
@	Concatenation	Left to Right
+	Union	Left to Right
!	Negation	Right to Left

```
dfa a;
dfa b;

dfa c = a*; <!-- Kleene Star --!>
dfa d = a@b; <!-- Concatenation --!>
dfa e = a+b; <!-- Union --!>
dfa f = !a; <!-- Negation --!>
```

Misc Operators

Operator	Description	Associativity
.	Access Struct Member	Left to Right
[]	Access Set Element	Left to Right
()	Function Call	Left to Right

```
struct Point {
    int_8 x;
    int_8 y;
}

Point p;
p.x = 10;

o_set<int_8> a = {1, 2, 3};
int_8 b = a[0];

int_8 func(int_8 a, int_8 b) {
    return a + b;
}

int_8 a = func(10, 20);
```

Operator Precedence

Operator	Description
()	Parentheses
!	Logical NOT
*, /, %	Multiplication, Division, Modulo
+, -	Addition, Subtraction
>, <, >=, <=	Comparison
==, !=	Equality
&&	Logical AND
\ \	Logical OR
=	Assignment
+=, -=, *=, /=, %=, &=, \ =	Assignment

Set Operators Precedence

Operator	Description
()	Parentheses
^2	Power Set
*,+, -	Intersection, Union, Set Difference

Automaton Operators Precedence:

Operator	Description
()	Parentheses
*,!	Kleene Star, Negation
@, +	Concatenation, Union

Control Flow

YAFSM enforce the use of curly braces for all control flow statements. YAFSM does not support the use of indentation for control flow statements. YAFSM supports the following control flow statements:

If-Else

Below is the syntax for the if-else statement:

```

if (condition) {
    statement;
}
elif (condition) {
    statement;
}
else {
    statement;
}

```

Loops

YAFSM only supports the **while** loop. Below is the syntax for the **while** loop:

```

while (condition) {
    statements;
}

```

Constants

Constants are represented by the **const** keyword. Constants can be of any data type supported by YAFSM.

Keywords

Keyword	Description
<code>int_x</code>	Integer
<code>uint_x</code>	Unsigned Integer
<code>char</code>	Character
<code>float_x</code>	Float
<code>bool</code>	Boolean
<code>const</code>	Constant
<code>struct</code>	Struct
<code>o_set</code>	Ordered Set
<code>u_set</code>	Unordered Set
<code>string</code>	String
<code>regex</code>	Regular Expression
<code>dfa</code>	DFA

Keyword	Description
<code>nfa</code>	NFA
<code>pda</code>	PDA
<code>cfg</code>	CFG
<code>if</code>	If
<code>elif</code>	Else If
<code>else</code>	Else
<code>while</code>	While
<code>break</code>	Break
<code>continue</code>	Continue
<code>return</code>	Return
<code>true</code>	True
<code>false</code>	False
<code><!--</code>	Start of comment
<code>--!></code>	End of comment

Identifiers

YAFSM uses the following rules for identifiers:

- Identifiers can only contain alphanumeric characters and underscores.
- Identifiers cannot start with a number.
- Identifiers cannot be a keyword.
- Identifiers cannot contain spaces.
- Identifiers cannot contain special characters.

Regular Expressions for Identifiers:

```
[a-zA-Z_][a-zA-Z0-9_]*
```

Statements

YAFSM supports the following statements:

Declaration Statement

Declaration statements are used to declare variables. Below is the syntax for declaration statements:

```
data_type identifier;
```

Multiple variables of the same data type can be declared in a single statement:


```
data_type identifier1, identifier2, ...;
```

Assignment Statement

Assignment statements are used to assign values to variables. Below is the syntax for assignment statements:

```
identifier = expression;
```

Function Declaration Statement

Function declaration statements are used to declare functions. Below is the syntax for function declaration statements:

```
data_type function_name(data_type1 arg1, data_type2 arg2, ...) {  
    statements;  
}
```

Function Call Statement

Function call statements are used to call functions. Below is the syntax for function call statements:

```
function_name(arg1, arg2, ...);
```

In case the function returns a value, the function call statement can be used as an expression:

```
data_type variable = function_name(arg1, arg2, ...);
```

IO Statements

Print statements are used to print values to the console. Below is the syntax for print statements:

```
out(expression);
```

Input statements are used to take input from the console. Below is the syntax for input statements:

```
inp(identifier);
```

In case multiple variables need to be inputted, the input statement can be used as:

```
inp(identifier1, identifier2, ...);
```

In-built Functions

YAFSM supports multiple in-built functions to help with working with Finite State Machines.

Set Functions

Function	Description	Return Type
<code>size(o_set<T> S)</code>	returns the size of the <code>o_set</code>	<code>int_8</code>
<code>size(u_set<T> S)</code>	returns the size of the <code>u_set</code>	<code>int_8</code>
<code>empty(o_set<T> S)</code>	returns true if S is empty otherwise returns false	<code>bool</code>
<code>empty(u_set<T> S)</code>	returns true if S is empty otherwise returns false	<code>bool</code>
<code>find(o_set<T> S, T x)</code>	finds x in S and returns true if $x \in S$, otherwise returns false	<code>bool</code>
<code>find(u_set<T> S, T x)</code>	finds x in S and returns true if $x \in S$, otherwise returns false	<code>bool</code>
<code>insert(o_set<T> S, T x, ...)</code>	inserts x into S (no duplicates)	-
<code>insert(u_set<T> S, T x, ...)</code>	- inserts x into S (no duplicates)	-
<code>remove(o_set<T> S, T x, ...)</code>	- remove x from S if $x \in S$	-
<code>remove(u_set<T> S, T x, ...)</code>	- remove x from S if $x \in S$	-
<code>delete(o_set<T> S)</code>	- deletes all elements from the <code>o_set</code>	-
<code>delete(u_set<T> S)</code>	deletes all elements from the <code>u_set</code>	-
<code>out(o_set<T> S)</code>	prints all elements of the <code>o_set</code>	-
<code>out(u_set<T> S)</code>	prints all elements of the <code>u_set</code>	-

CFG Functions

Function	Description	Return Type
<code>add_T(cfg a, X)</code>	adds a terminal X to the CFG	-
<code>add_NT(cfg a, X:“val”)</code>	adds a non-terminal “val” with name X to the CFG	-
<code>add_P(cfg a, X)</code>	adds a production to the CFG where the production is in one of the allowed forms	-
<code>remove_T(cfg a, X)</code>	removes a terminal from CFG along with all the rules that use X	-
<code>remove_NT(cfg a, X)</code>	removes a non-terminal with name X from CFG along with all the rules that use X	-

Function	Description	Return Type
remove_P(cfg a, X)	removes a production from CFG	-
change_start(cfg a, X)	changes the start symbol of CFG to X where X must be a terminal	-
test_membership(cfg a, X)	returns true if X is a member of the language generated by CFG	bool
cfg_to_PDA(cfg a)	converts the CFG a to a PDA and returns the PDA	pda
delete(cfg a)	deletes the CFG a	-
out(cfg a)	prints the CFG a	-

DFA Functions

Function	Description	Return Type
add_states(dfa a, X, ...)	adds states X, ... to the DFA along with transitions from/to X, ...	-
remove_states(dfa a, X)	removes states X, ... from the DFA along with transitions from/to X, ...	-
add_letters(dfa a, X, ...)	adds letters X, ... to Σ of the DFA a	-
remove_letters(dfa a, X)	removes letters X, ... from Σ of a	-
change_start(dfa a, X)	changes the start state of the DFA to X where X	-
add_final(dfa a, X, ...)	adds states X, ... to the set of final states	-
remove_final(dfa a, X, ...)	removes states X, ... from the set of accepting states	-
add_transition(dfa a, X, ...)	adds a transition X, ... where X is in one of the forms of the transitions mentioned in the DFA section	-
remove_transition(dfa a, X)	removes the transition X from the DFA	-
simulate(dfa a, X)	returns true if X is accepted by the DFA a, otherwise false	bool
delete(dfa a)	deletes the DFA a	-
out(dfa a)	prints the DFA a	-

NFA Functions

Functions that work with DFAs also work with NFAs. Some additional functions are

Function	Description	Return Type
<code>nfa_to_dfa(nfa a)</code>	converts the NFA <code>a</code> to a DFA	<code>dfa</code>

PDA Functions

Function	Description	Return Type
<code>insert_states(pda a, X, ...)</code>	inserts states <code>X, ...</code> into the PDA	-
<code>remove_states(pda a, X, ...)</code>	removes states <code>X, ...</code> from the PDA	-
<code>insert_input_letters(pda a, X, ...)</code>	inserts letters <code>X, ...</code> into the PDA	-
<code>remove_input_letters(pda a, X, ...)</code>	removes letters <code>X, ...</code> from the PDA	-
<code>insert_stack_letters(pda a, X, ...)</code>	inserts letters <code>X, ...</code> into the PDA	-
<code>remove_stack_letters(pda a, X, ...)</code>	removes letters <code>X, ...</code> from the PDA	-
<code>change_start(pda a, X)</code>	changes the start state of the PDA	-
<code>insert_final(pda a, X, ...)</code>	inserts states <code>X, ...</code> into the PDA	-
<code>remove_final(pda a, X, ...)</code>	removes states <code>X, ...</code> from the accepting states of PDA	-
<code>insert_transition(pda a, X, ...)</code>	inserts a transition <code>X, ...</code> into the PDA	-
<code>remove_transition(pda a, X, ...)</code>	removes a transition <code>X, ...</code> from the PDA	-
<code>simulate(pda a, X, ...)</code>	returns true if <code>X</code> is accepted by the PDA <code>a</code> , otherwise	<code>bool</code>
<code>delete(pda a)</code>	deletes the PDA	-
<code>out(pda a)</code>	prints the PDA <code>a</code>	-

Regular Expression Functions

Function	Description	Return Type
<code>simulate(regex a, string s)</code>	returns true if s is accepted by the regex a, otherwise false	bool
<code>delete(regex a)</code>	deletes the regex a	-
<code>regex_to_dfa(regex a)</code>	converts the regex a to a DFA and returns the DFA	dfa
<code>regex_to_nfa(regex a)</code>	converts the regex a to a NFA and returns the NFA	nfa

References

- Wikipedia: FSMs
- Wikipedia: DFAs
- Wikipedia: NFAs
- Wikipedia: PDAs
- Wikipedia: CFGs
- Wikipedia: Regular Expressions
- Michael Sipser: Introduction to the Theory of Computation
- C_Programming Language by Kernighan and Ritchie