

Finite Automata

COP-3402 Systems Software

Paul Gazzillo

1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

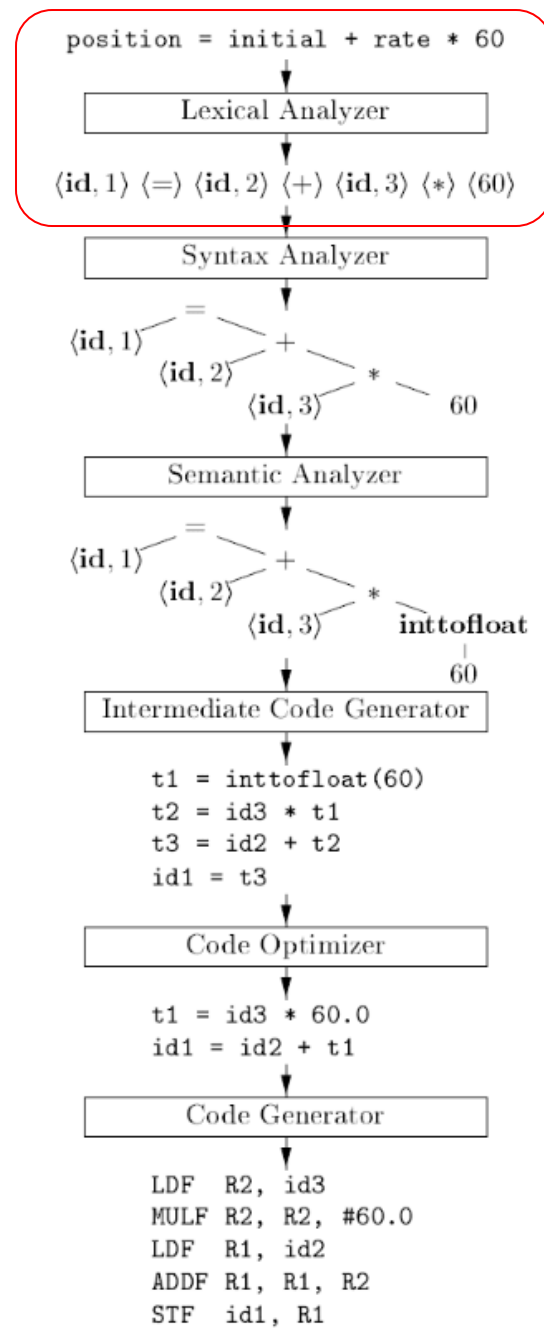


Figure 1.7: Translation of an assignment statement

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
if	characters <code>i</code> , <code>f</code>	<code>if</code>
else	characters <code>e</code> , <code>l</code> , <code>s</code> , <code>e</code>	<code>else</code>
comparison	<code><</code> or <code>></code> or <code><=</code> or <code>>=</code> or <code>==</code> or <code>!=</code>	<code><=</code> , <code>!=</code>
id	letter followed by letters and digits	<code>pi</code> , <code>score</code> , <code>D2</code>
number	any numeric constant	<code>3.14159</code> , <code>0</code> , <code>6.02e23</code>
literal	anything but <code>"</code> , surrounded by <code>"</code> 's	<code>"core dumped"</code>

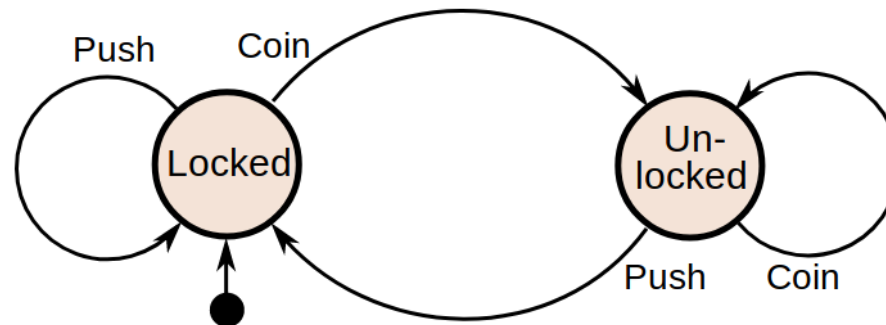
Figure 3.2: Examples of tokens

Core Problem: Recognizing Tokens

- Regular expressions *specify* tokens
- How do we recognize them?
- Use *finite automata*, a.k.a., finite-state machines

Finite Automata are a Model of Computation

- A model of computation: transitions between states
 - Finite alphabet of symbols (just like regular languages)
 - Finite set of states
 - Set of accepting states
 - An initial state
 - A transition function: $f(\text{current_state}, \text{symbol}) \rightarrow \text{next_state}$



Finite Automata Have Wide Application

- Vending machines: count coins
- Elevators: sequence of stops
- Traffic lights: order of changes
- Combination lock: numbers in correct order

Representation with Transition Diagrams

- States: circles
- Transitions: labeled arrows
- Starting state: arrow
- Accepting states: double circles

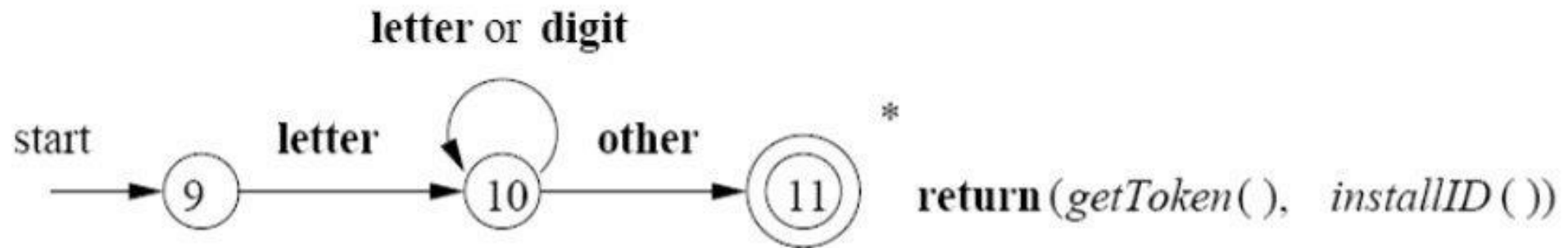
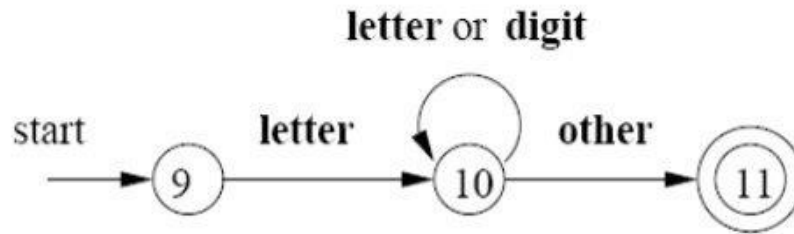


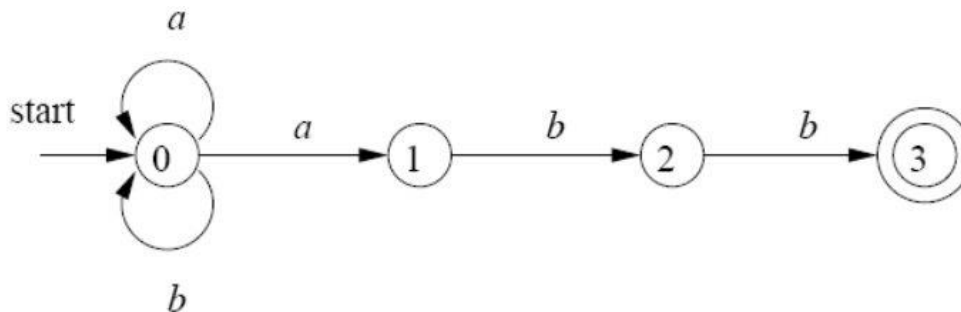
Figure 3.14: A transition diagram for **id**'s and keywords

Deterministic vs Nondeterministic Finite Automata (DFA vs NFA)

- *Deterministic*: one state at-a-time
 - We can uniquely "determine" which state we are in



- *Nondeterministic*: in multiple states at once
 - Must track all states *simultaneously*



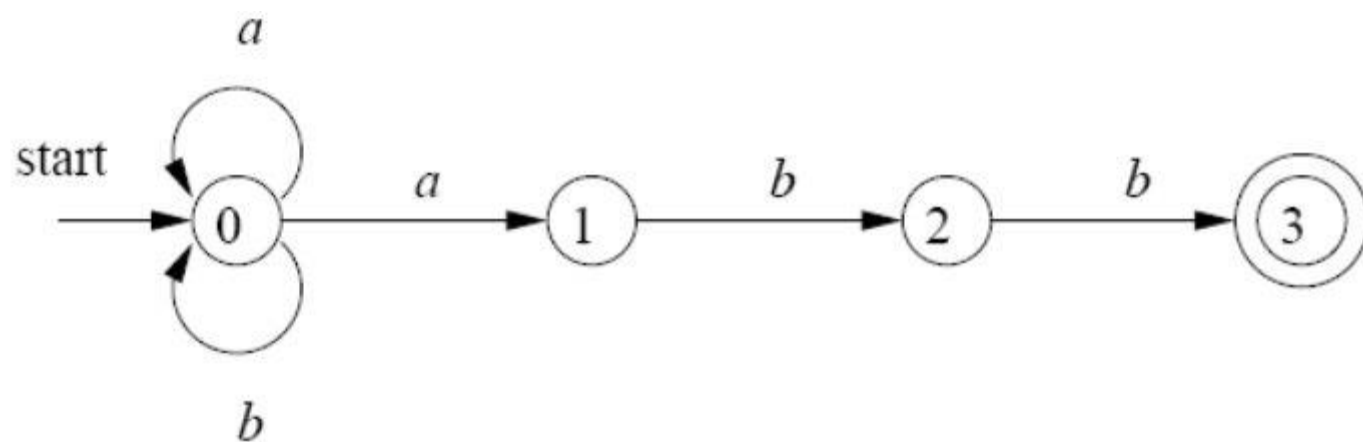


Figure 3.24: A nondeterministic finite automaton

STATE	a	b	ϵ
0	$\{0, 1\}$	$\{0\}$	\emptyset
1	\emptyset	$\{2\}$	\emptyset
2	\emptyset	$\{3\}$	\emptyset
3	\emptyset	\emptyset	\emptyset

Figure 3.25: Transition table for the NFA of Fig. 3.24

Remember: Epsilon (ϵ) is Empty String

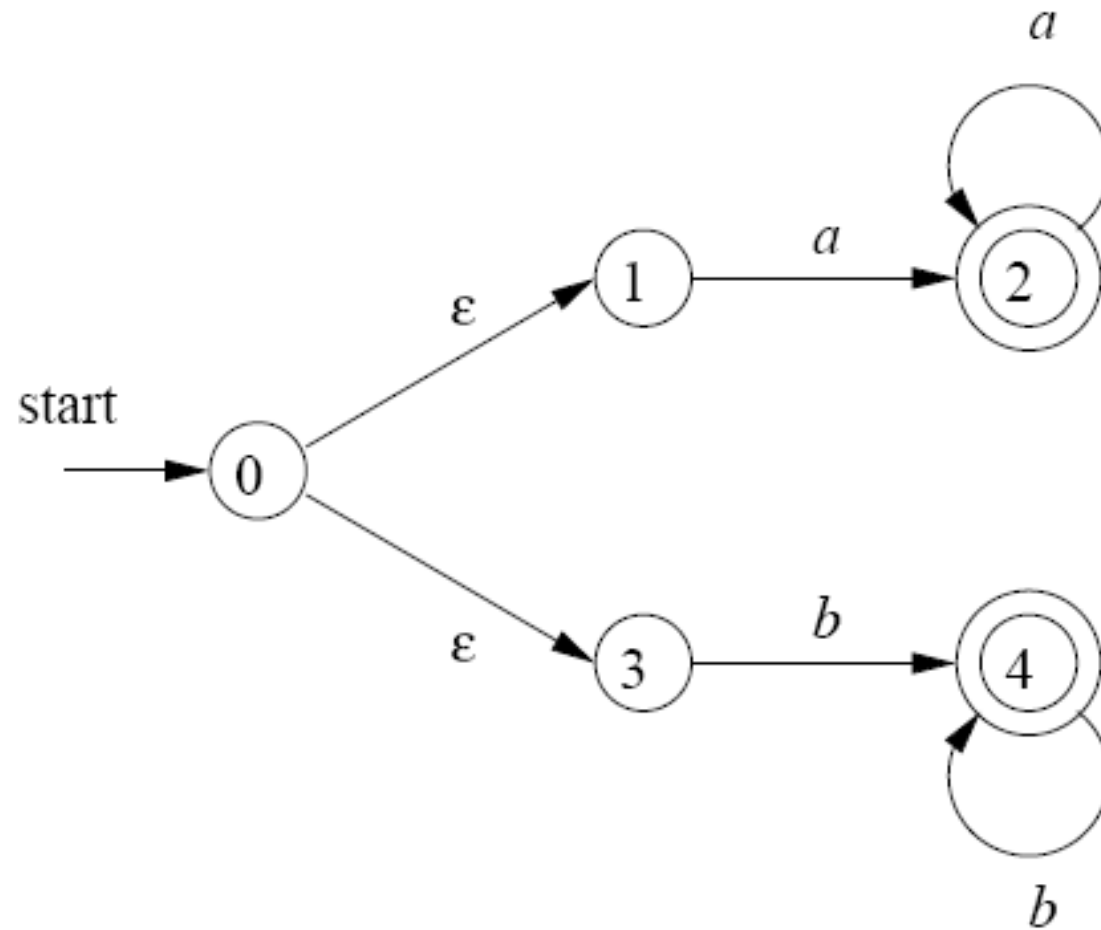
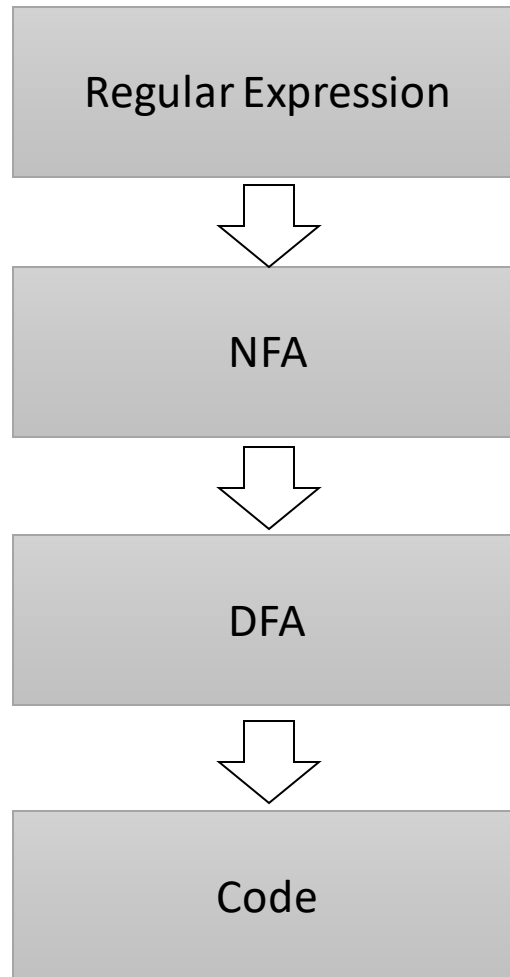


Figure 3.26: NFA accepting **$aa^*|bb^*$**

Using Finite Automata to Implement a Lexer



Regular Expression to NFA

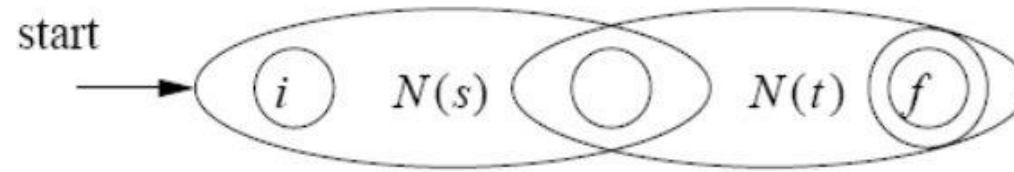


Figure 3.41: NFA for the concatenation of two regular expressions

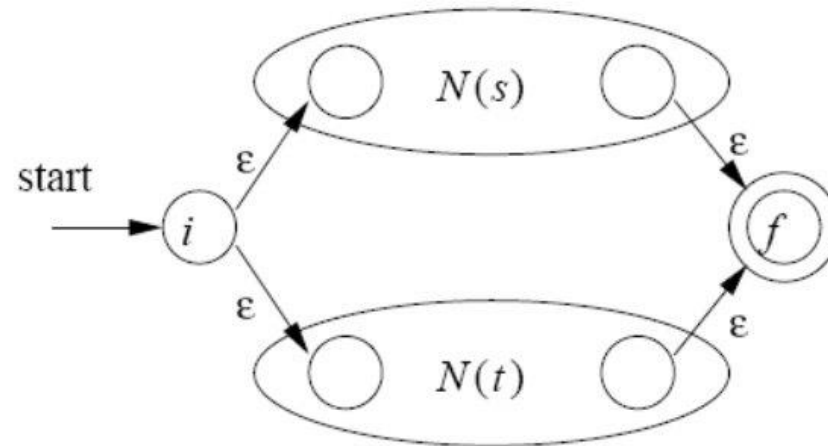


Figure 3.40: NFA for the union of two regular expressions

Regular Expression to NFA (continued)

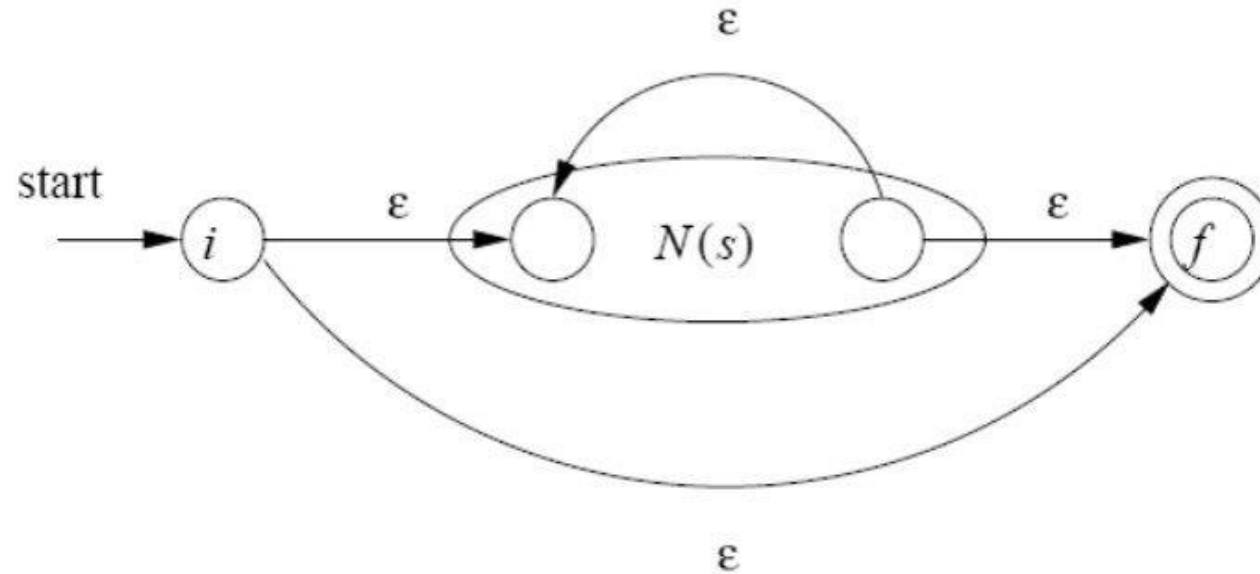


Figure 3.42: NFA for the closure of a regular expression

NFA to DFA with Subset Construction

- Construct DFA from NFA systematically
- Each DFA state created from subset of NFA states
 - Remember: can be in multiple states
- "Simulate" being in multiple states using a single state
- Dragon book 3.7

Example

Converting regular expressions to nondeterministic and deterministic automata

Hand-Crafted Finite Automata in Code

- Lexer reads one character at a time, e.g., `fgetc`
- Concatenation is sequence of statements
- Union is a conditional
- Closure is a while loop

Concatenation is sequence of statements, e.g., ab

```
char c;  
c = fgetc(lexerin);  
assert 'a' == c  
c = fgetc(lexerin);  
assert 'b' == c
```

Union is a conditional, e.g., $a|b$

```
char c;  
c = fgetc(lexerin);  
if ('a' == c) {  
    // ...  
} else if ('b' == c) {  
    // ...  
} else {  
    error(...)  
}
```

Closure is a while loop, e.g., a^*

```
char c;  
c = fgetc(lexerin);  
while('a' == c) {  
    c = fgetc(lexerin);  
}  
// c is now the next character
```

EBNF Standardizes Language Specifications

- EBNF = Extended Backus-Naur Format
- Different notation for regular expressions
 - a^* $\rightarrow \{ a \}$
 - $a|\epsilon$ $\rightarrow [a]$
 - ab $\rightarrow a b$
 - $a | b$ $\rightarrow a | b$

Demo

Lexing PL/0