

# Creating a Lexical Scanner

COMP360

*“Programming today is a race between software engineers striving to build bigger and better idiot-proof programs, and the universe trying to build bigger and better idiots.*

*So far, the universe is winning.”*

Rick Cook

# Stages of a Compiler

- Source preprocessing
- Lexical Analysis (scanning)
- Syntactic Analysis (parsing)
- Semantic Analysis
- Optimization
- Code Generation
- Link to libraries

# Source Preprocessing

- In C and C++, preprocessor statements begin with a #
- The preprocessor edits the source code based on the preprocessor statements
- **#include** is the same as copying the included file at that point with the editor
- The output of the preprocessor is expanded source code with no # statements
- Old C compilers had a separate preprocessor program

# Lexical Analysis

- Lexical Analysis or scanning reads the source code (or expanded source code)
- It removes all comments and white space
- The output of the scanner is a stream of tokens
- Tokens can be words, symbols or character strings
- A scanner can be a finite state automata (FSA)

# Syntactic Analysis

- Syntactic Analysis or parsing reads the stream of tokens created by the scanner
- It checks that the language syntax is correct
- The output of the Syntactic Analyzer is a parse tree
- The parser can be implemented by a context free grammar stack machine

# Semantic Analysis

- The Semantic Analysis inputs the parse tree from the parser
- Language requirements not checked by the syntax are enforced
- This stage determines what the program is to do
- The output of the Semantic Analysis is an intermediate code. This is similar to assembler language, but may include higher level operations

# Simple Program

```
/* This is an example program */  
bull = cow + cat * dog;
```



# After Lexical Scan

bull

=

cow

+

cat

\*

dog

;

# Output of Each Stage

- Source preprocessing – expanded source code
- Lexical Analysis – List of tokens
- Syntactic Analysis – Parse Tree
- Semantic Analysis – Intermediate code
- Optimization – Intermediate code
- Code Generation – Object file
- Link to libraries – Executable program

# Machines of a Compiler

- Source preprocessing – simple editing
- Lexical Analysis – Finite State Automata
- Syntactic Analysis – Push Down Automata
- Semantic Analysis
- Optimization
- Code Generation
- Link to libraries

# State Tables

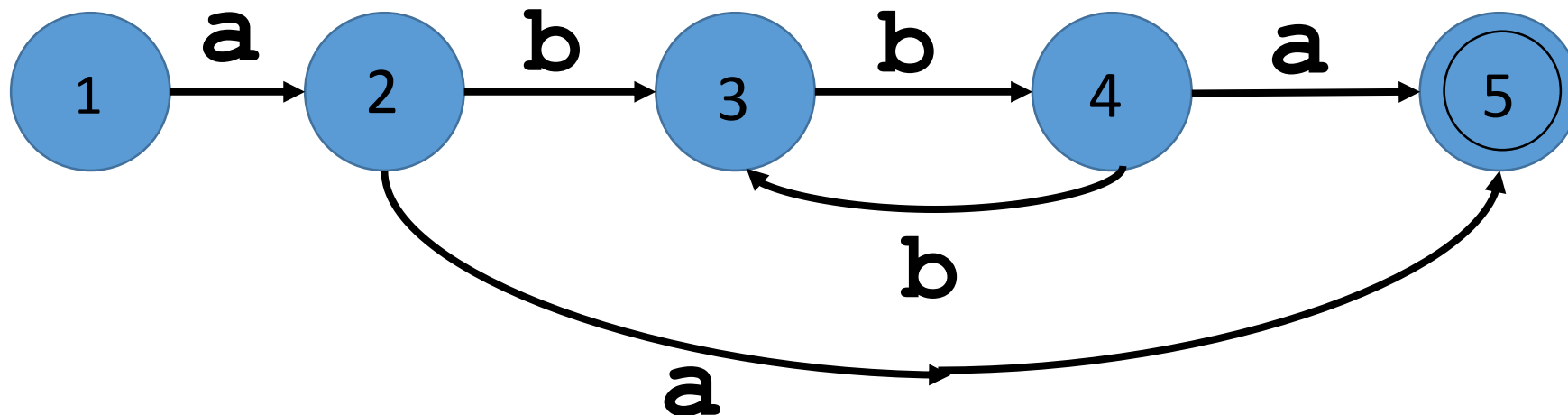
- An FSA graph can be converted to a table
- The table has cells for each state and each input symbol
- In the cell goes the next state if the DFA is in that state and receives that input symbol
- You can consider the state table to be an adjacency table for the graph

# Converting an Example FSA

- Consider the regular expression that begins and ends with an **a** and can have an even number of **b**'s between them

**a (bb)\* a**

- It can be recognized by the FSA



# Convert the Graph to a Table

	1	2	3	4	5
a	2	5	0	5	0
b	0	3	4	3	0

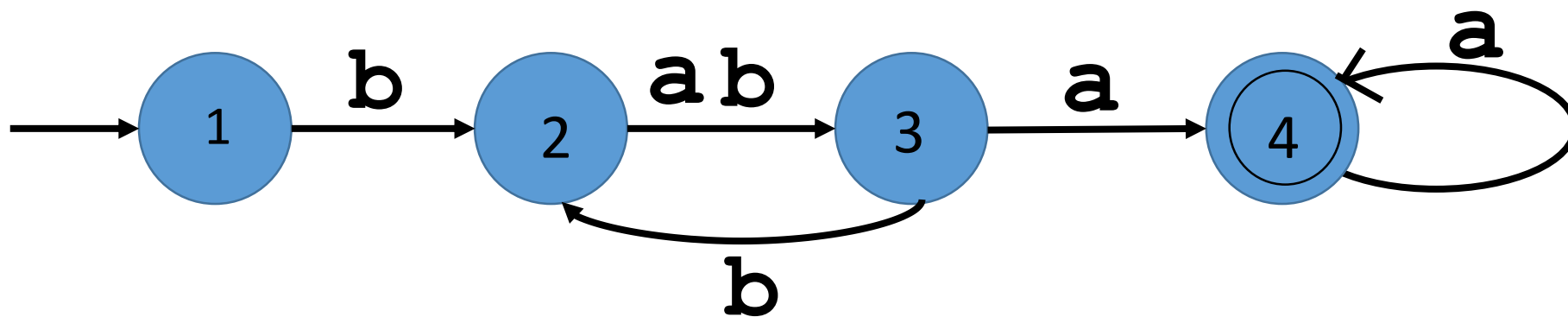
- The states are listed along the top
- The input symbols are along the side
- For that symbol while in that state, the DFA will go to the new state given in the table
- State zero represents a final error state

Draw a DFA for this Regular Language

**(bb | ba) a<sup>+</sup>**

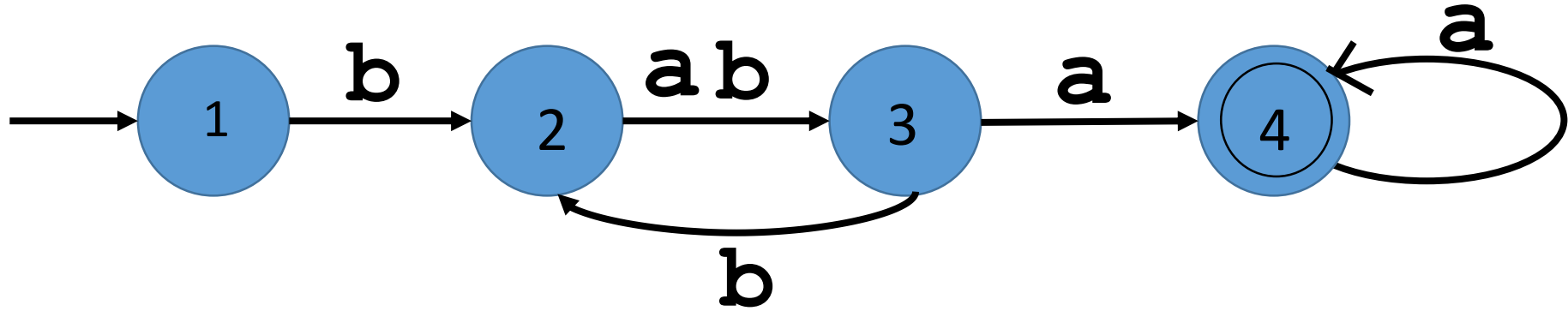
# Possible Solution

**(bb | ba) a<sup>+</sup>**





# Create a state table for the FSA



	1	2	3	4
a	0	3	4	4
b	2	3	2	0

# Programming a FSA

- It is relatively simple to implement a Finite State Automata in a modern programming language
- This program can be used to recognize if a string conforms to a regular language

# FSA Program

```
state = 1
```

```
while not end of file {
```

```
    symbol = next input character
```

```
    state = stateTable[ symbol, state ]
```

```
    if state = 0 then error
```

```
}
```

```
if state is a terminating state, success
```

# Grouping Input Symbols

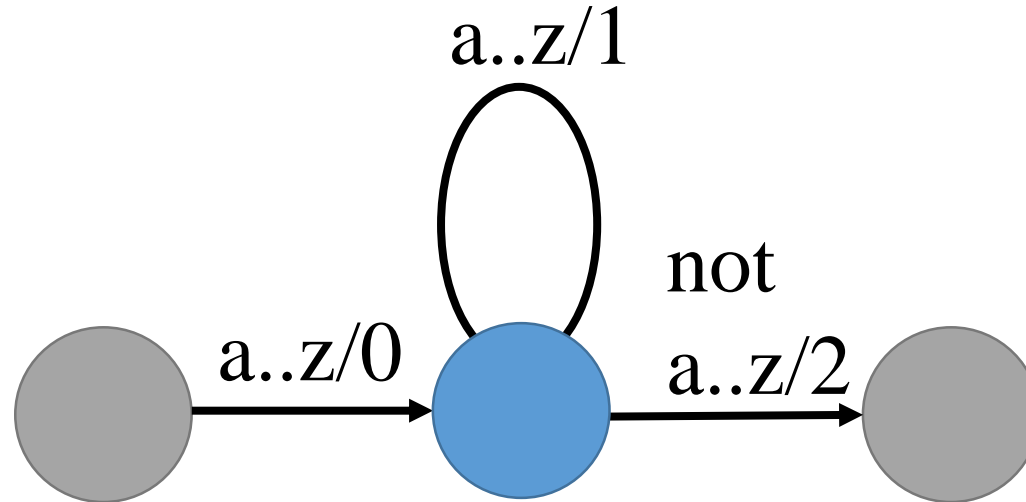
- For many FSA programs, it is useful to create an index value for the input symbols, i.e.  $a = 0$ ,  $b = 1$
- Often you can have groups of symbols use the same index value, i.e. all letters have the index 2 and all numbers have the index 3

# Mealy and Moore Machines

- The FSA we have discussed so far simply determine if the input is valid for the specified language
- An FSA can also produce an output
- A Mealy machine has an output or function associated with each transition or edge of the graph
- A Moore machine has an output or function associated with each state or node of the graph

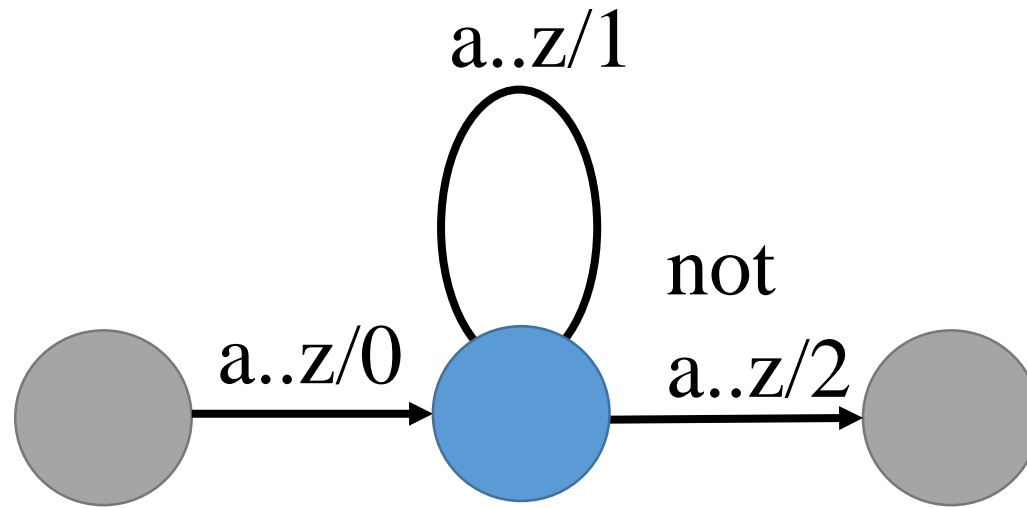
# Using Mealy Machines to Create Tokens

- Consider an FSA reading text and creating token of words separated by spaces or punctuation



- 0 = Save character as first letter in a string
- 1 = Save character as next letter in a string
- 2 = Save the string as a token, handle other symbol

# Mealy Machine State Table



	1	2	3
a .. z	2/0	2/1	2/0
not a .. z	1/x	3/2	3/x

New state / Output function

# Lexical Analysis with a Mealy Machine

- Compilers can use a Mealy machine to scan the source code
- The FSA recognizes and discards comments and white space
- Names, numbers, strings and punctuation are each output as a list of tokens
- A token is an object that contains one unit of the input specifying the value and type



# Coding a Mealy Machine

- You can use two 2D arrays of integers indexed by input symbol type and current state OR
- You can use a 2D array of objects indexed by input symbol type and current state
- One value of each array cell is the new state
- The other value tells the program what action to take at this transition

# Common Transition Actions

- Create a token with the saved character
- Save another character in the token
- Do nothing

# Design a DFA

- Create a DFA that removes // comments and saves strings of letters

# Token Object

- A token created by the lexical scanner should contain:
- String value; // text of the object
- int type; // type of value
  - name
  - number
  - punctuation
- int line number, column // position in source code