

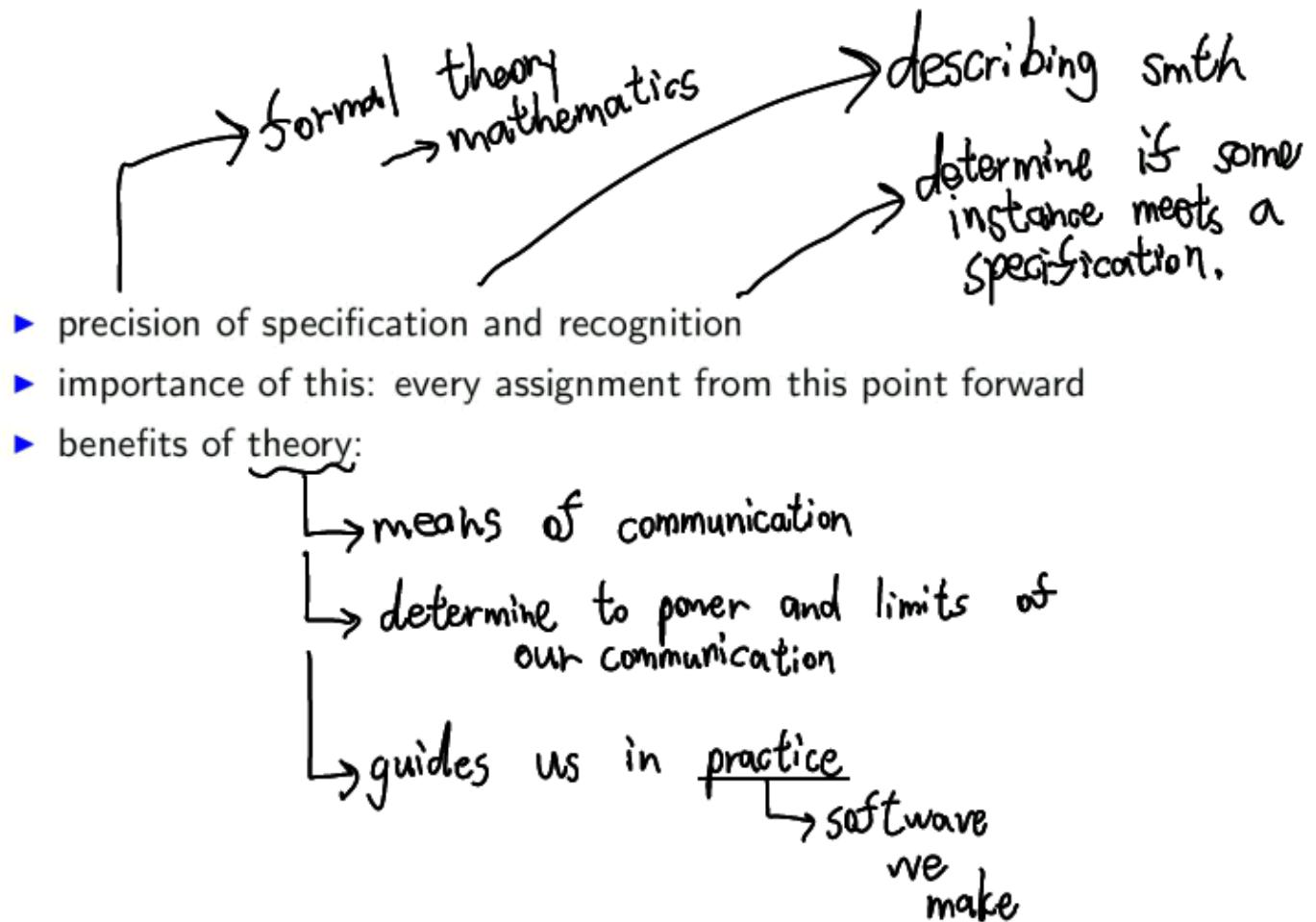
Lecture 9

Introduction to Formal Languages

CS 241: Foundations of Sequential Programs
Fall 2014

Troy Vasiga et al
University of Waterloo

Motivation



Terminology

- ▶ rooted in set theory $\left[\cdot \right]$ contains unique objects
- ▶ alphabet: a finite set of symbols

$$\Sigma_1 = \{a, b, c\} \quad \Sigma_2 = \{\square, \blacktriangle, \cancel{\star}, 173, \text{blue}\}$$

\uparrow
sigma

- ▶ word (aka string, sentence): finite sequence of symbols from the alphabet

word w_1 over Σ_1 : $w_1 = baab$

w_2 over Σ_1 : $w_2 = ccc$

w_3 over Σ_1 : $w_3 = \epsilon$ ← the empty word
 \uparrow
epsilon

Terminology (continued)

- ▶ language: set of words

\uparrow possibly infinite

$$L_1 \text{ over } \Sigma_1: L_1 = \{ba, baa, baaa\}$$

$$L_2 \text{ over } \Sigma_1: L_2 = \{a, aa, aaa, \dots\}$$

- ▶ $|W|$: the size of W (where W is a word or language)

$$|L_1| = 3 \quad |L_2| = \infty$$

$$|w_1| = |baab| = 4$$

$$|\epsilon| = 0$$

Uses of Formal Languages: Specification

A statement of what
is in the language

Ideally, we have 3 goals of a specification:
mandatory → ① precise
 ② easy-to-use
 ③ automatable

Uses of Formal Languages: Recognition

determine if a given word
meets the specification.

Formally: Given a language L and a word w
over alphabet Σ_i , is $w \in L$?

w in L

In practice: We want more than just "yes/no"

certificate
of validity

reason for
invalid:
⇒ where?
⇒ how?
⇒ how to fix?

WLP4

Let's try to formally specify WLP4.

- ▶ alphabet:

Valid WLP4 tokens

$\langle \text{IF} \rangle, \langle \text{WHILE} \rangle, \dots$

- ▶ word:

sequence of WLP4 tokens
 \Rightarrow WLP4 program

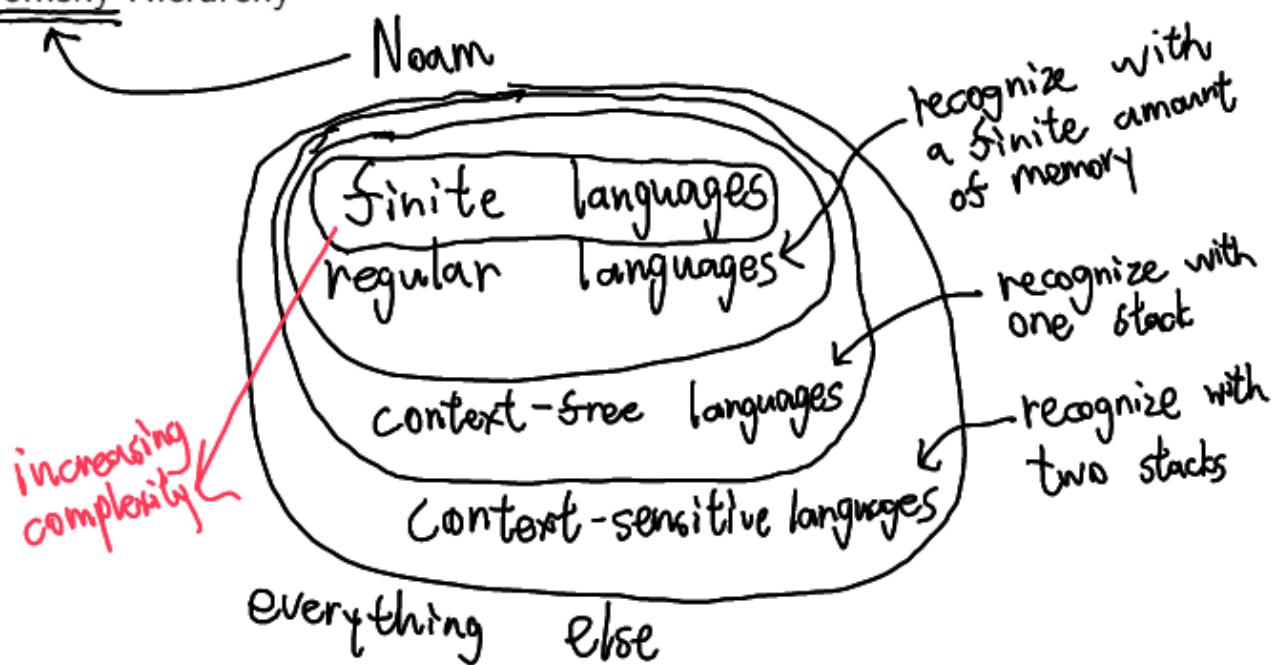
- ▶ language:

{valid WLP4 programs}

Waterloo
Language
plus
Pointers
plus
Procedures

Language Classes

- ▶ a set of languages may share common characteristics
- ▶ Chomsky Hierarchy



Uses of Formal Languages: Organization of Compilation

- ▶ lexical analysis
 - ▶ syntactic analysis
 - ▶ context-sensitive analysis (semantic) analysis
 - ▶ synthesis (code generation)
- find each lexical element
syntax correct \Rightarrow context-free languages
- regular languages
- context-sensitive languages
-
- The diagram illustrates the organization of compilation processes. It shows four main steps: lexical analysis, syntactic analysis, context-sensitive analysis (semantic analysis), and synthesis (code generation). A bracket groups syntactic analysis and context-sensitive analysis under the heading 'syntax correct \Rightarrow context-free languages'. Another bracket groups all four steps under the heading 'context-sensitive languages'. An arrow points from the first two steps to the first part of the bracket. A handwritten note 'regular languages' is associated with the first step. A handwritten note 'find each lexical element' is placed above the first two steps.

Which language level?

Let $\Sigma = \{\text{ASCII characters}\} - \{\text{CR}\}$.

- $L_1 = \text{all valid assembly language programs}$

regular: binary

context-sensitive: assembly

- $L_2 = \{\$0, \$1, \$2, \dots, \$31\}$

\Rightarrow finite language

- $L_3 = \text{valid labels in MIPS}$

\Rightarrow regular language

- $L_4 = \text{valid load word (lw) offsets}$

\Rightarrow finite

lw \$3, \$7
↑ 16-bit integer

- $L_5 = \text{valid line of AL for A3P3}$

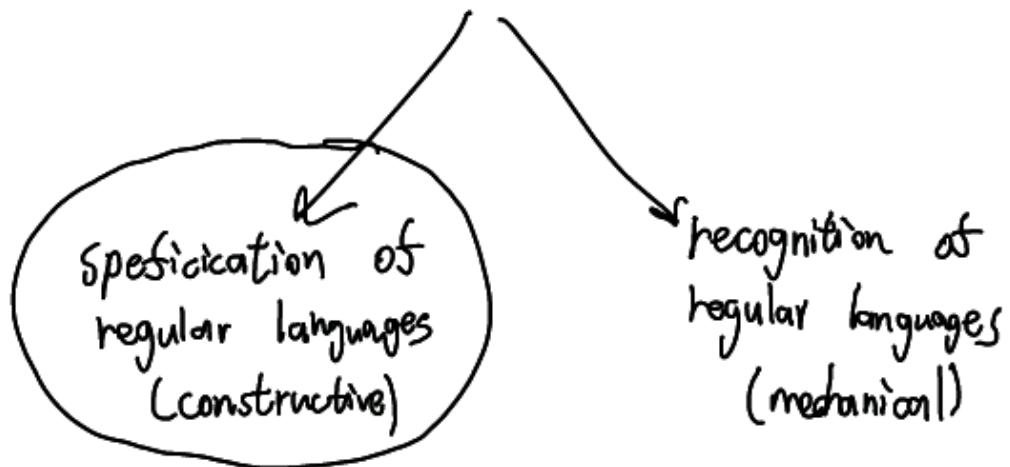
.word int < 4 billion
.word hexint < 4 billion \Rightarrow finite language

- $L_6 = \text{valid line of AL for A3P4}$

label: .word blat \Rightarrow context-sensitive language

Regular Languages

- ▶ Defining regular languages: two approaches



- ▶ We will see that these are equivalent

Specifying Regular Languages

Basic building blocks

1. finite languages
2. union
3. concatenation
4. repetition

Union

- ▶ Definition: If T_1 & T_2 are regular languages, then
 $T_1 \cup T_2 = \{x : x \in T_1 \text{ or } x \in T_2\}$
is a regular language.

- ▶ Examples:

$$T_1 = \{\text{dog, cat}\}$$

$$T_2 = \{\text{red, blue, cat}\}$$

$$T_1 \cup T_2 = \{\text{dog, cat, red, blue}\}$$

Concatenation

- Definition: If T_1 & T_2 are regular languages, then

$T_1 \cdot T_2 = \{xy : x \in T_1 \text{ and } y \in T_2\}$
is regular.

- Examples:

$$T_1 = \{\text{dog, cat}\}$$

$$w \cdot \epsilon = w$$

$$T_2 = \{\text{fish, } \epsilon\}$$

$$T_1 \cdot T_2 = \{\text{dogfish, catfish, cat, dog}\}$$

Repetition

- ▶ Definition: If T is a reg. lang., then

$$T^* = \{\epsilon\} \cup \{xy : x \in T, y \in T^*\}$$

is regular.

- ▶ Alternate Definition:

$$\begin{aligned} T^0 &= \{\epsilon\} & T^1 &= T & T^2 &= T \cdot T = \{xy : x \in T, y \in T\} \\ T^3 &= T \cdot T^2, \dots, & T^* &= \boxed{\bigcup_{k=0}^{\infty} T^k} \end{aligned}$$

- ▶ Examples:

$*$ ← Kleene
Closure

$$T = \{\text{man}, \text{dog}\}$$

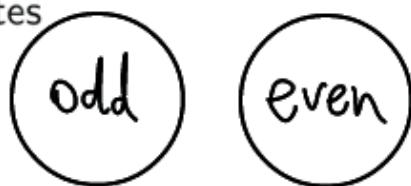
$$T^* = \{ \epsilon, \text{man}, \text{dog}, \text{mandog}, \text{dogman}, \text{dogmandog}, \text{dogdog}, \text{mannan}, \dots \}$$

Recognizers: Finite Automata

Regular languages can be recognized by *finite automata*.

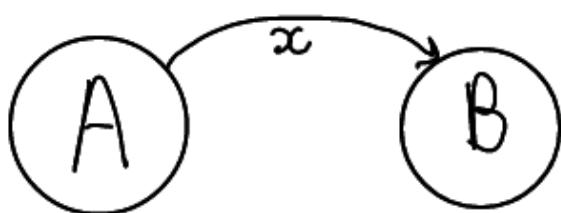
We begin with *deterministic finite automata*, also called DFAs.

- states



number of states

- transitions

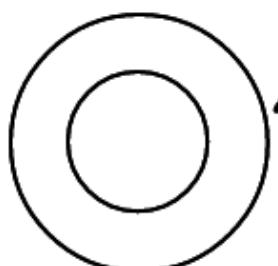


input x moves from state A to state B

- start state



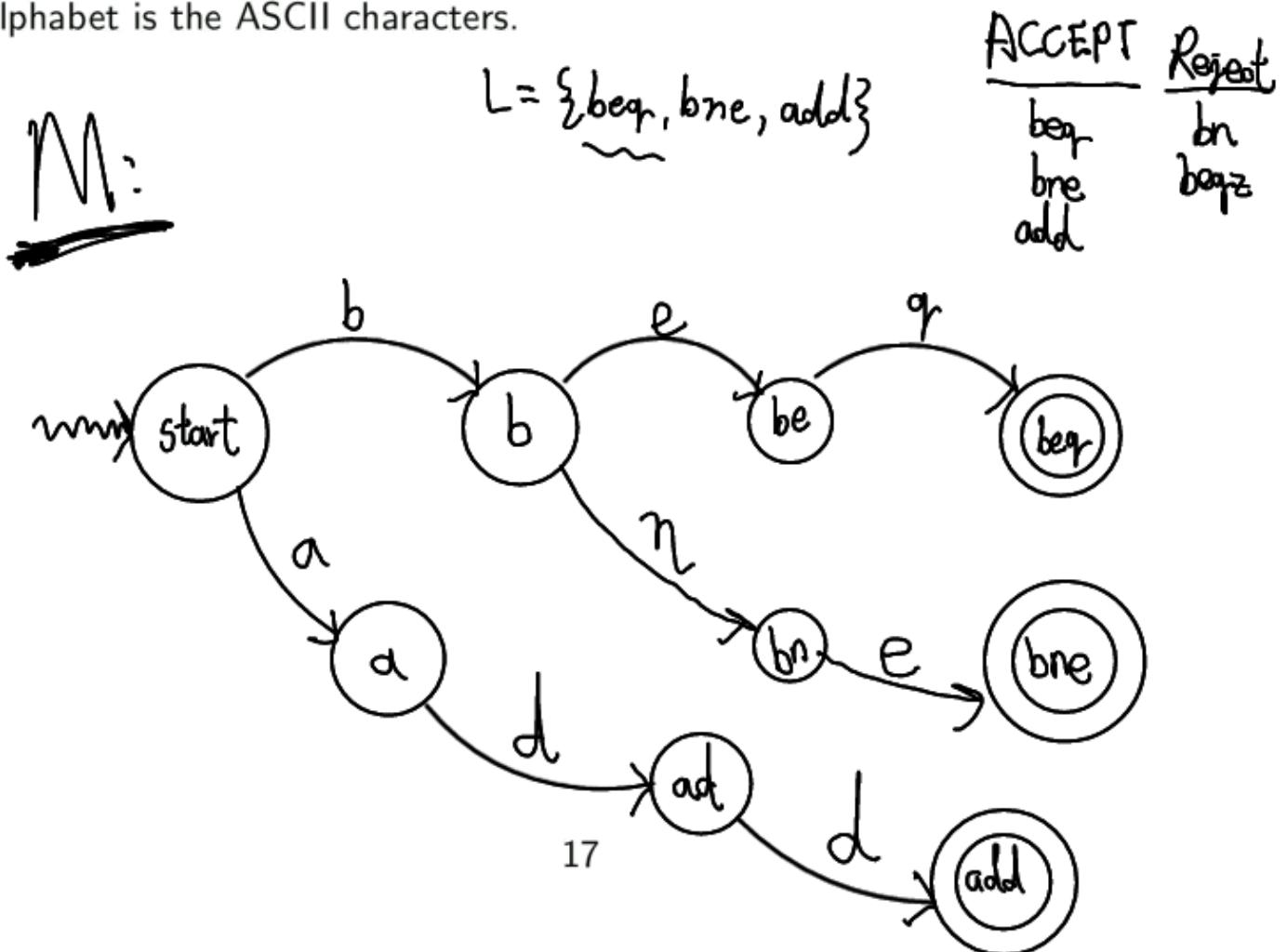
- final states



where I want to be at the end

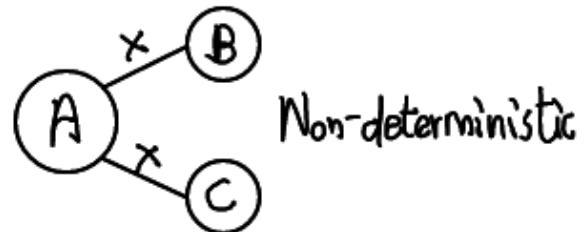
Finite Automata Example 1

Example: selected opcodes from MIPS assembly language, where alphabet is the ASCII characters.



Observations About Finite Automata

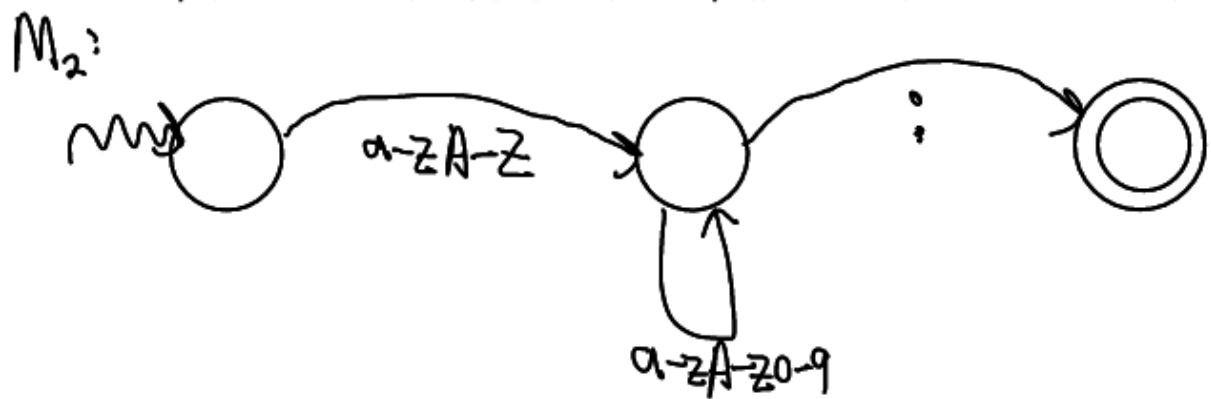
- ▶ ability to trace
 - ↳ Remember which state I am in.
- ▶ transitions out of a state are unique
 - ↳ deterministic
- ▶ errors
 - ↳ implicit error state \Rightarrow non-accepting, never leave it.
- ▶ size of this language
 - 3



- ▶ DFA M and language $L(M)$
 - ↳ language accepted by M
 - $L(M) = \{ \text{beg, bne, add} \}$

Finite Automata Example 2

Example: MIPS labels, where the alphabet is ASCII characters.

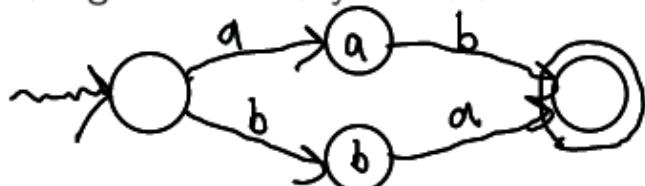


$$|L(M_2)| = \infty$$

More Finite Automata Examples

Let $\Sigma = \{a, b, c\}$.

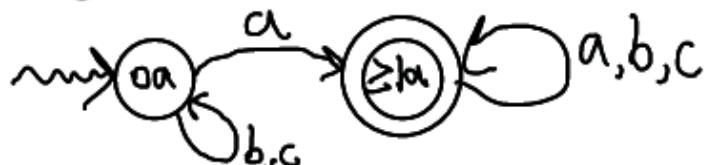
- ▶ strings with exactly one a and exactly one b and no c's



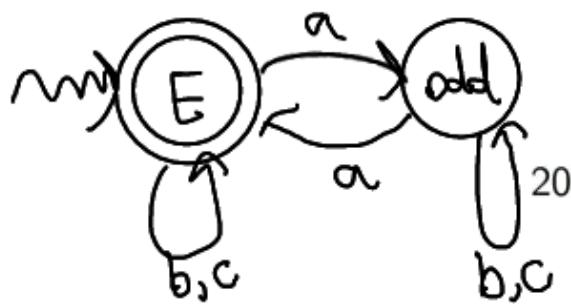
- ▶ strings with one a, one b and one c

Ex:

- ▶ strings with at least one a



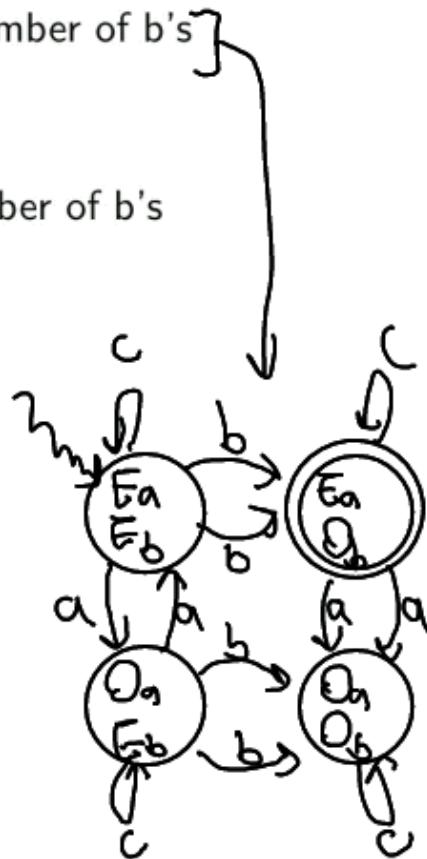
- ▶ string with an even number of a's



More Finite Automata Examples

Let $\Sigma = \{a, b, c\}$.

- ▶ strings with an even number a's and odd number of b's
- ▶ strings with an even number a's or odd number of b's
- ▶ valid MIPS identifiers
- ▶ valid register designations in MIPS
- ▶ valid lines of MIPS assembly language



Killer app for Finite Automata

Scanner: see asm.*