

Lexical Anaylsis

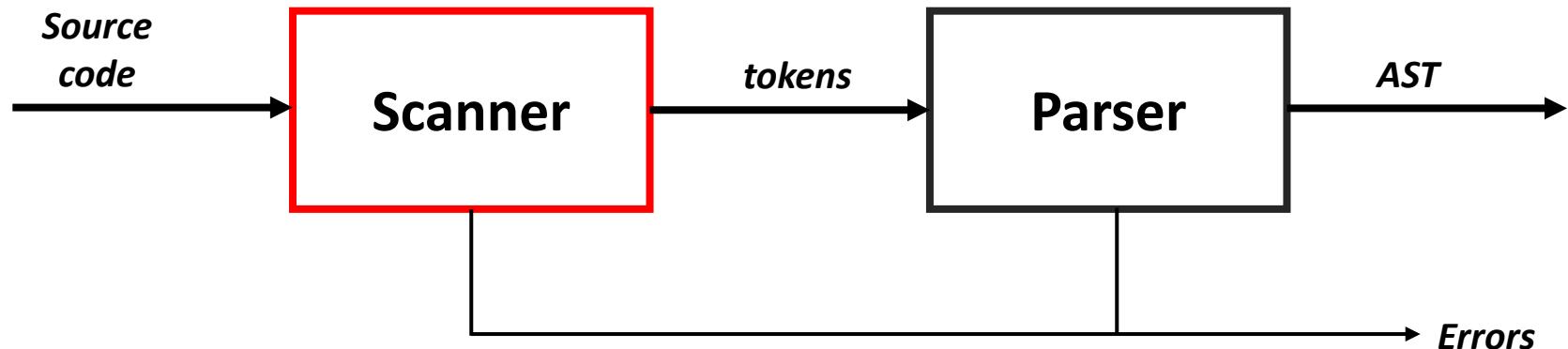
Jeonggeun Kim
Kyungpook National University



Outline

- The role of a scanner
- Scanner concepts
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata
 - Definitions of REs, DFAs and NFAs
 - REs → NFA
 - (Thompson's construction, Algorithm 3.3., Red Dragon book, Algorithm 3.23, Purple Dragon book)
 - NFA → DFA
 - (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon book)
 - DFA → minimal-state DFA
 - (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
- Scanner generators

The Role of the Scanner



Scanner

- Maps stream of characters into words called **tokens**.
 - Basic units of syntax
 - $x = y + 21;$ becomes...
 $\langle id, x \rangle \langle = \rangle \langle id, y \rangle \langle + \rangle \langle intliteral, 21 \rangle \langle ; \rangle$
- Characters that form a word are its **lexeme**.
- Formally, a token is a tuple $\langle Token_type, value \rangle$.
- Informally, we often say “token x” when we mean $\langle id, x \rangle$
- Tokens are to programming languages what words are to natural languages.

The Role of the Scanner

- Lexeme: a sequence of characters in the source program that matches the pattern for a token
 - a lexeme is identified by the lexical analyzer as an instance of that token.
- Token: a pair consisting of a token name and an optional attribute value.
 - token name is an abstract symbol representing a kind of lexical unit, e.g., a particular keyword or sequence of input characters (e.g., IDENTIFIER, ASSIGN)

[Token]	[Informal Description]	[Sample Lexemes]
if	characters i, f	if
else	characters e, l, s, e	else
comparison	< or > or <= or >= or == or !=	<=, !=
id	letter followed by letters and digits	pi, score, D2
number	any numeric constant	3.14159, 0, 6.02e23
literal	anything but ", surrounded by "'s	"core dumped"

x = a + b * 2

Lexemes: {x, =, a, +, b, *, 2}

Tokens: { <id, 0>, <=>, <id, 1>, <+>, <id, 2>, <*>, <id, 3> }

Tokens

- The tokens of our MiniC language are classified into **token types**:
 - identifiers: i, j, initial, position, ...
 - keywords: if, for, while, int, float, bool, ...
 - operators: + - * / <= && ...
 - separators: { } () [] ; ,
 - literals
 - integer literals: 0, 1, 22, ...
 - float literals: 1.25 1. .01 1.2e2 ...
 - bool literals: true, false
 - string literals: “hello\n”, “string literal”, ...
- The exact token set depends on the given programming language. Pascal and Ada use “:=” for assignment, C uses “=”.
- Natural languages also contain different kinds of tokens (words): verbs, nouns, articles, adjectives, ... The exact token set depends on the natural language in question.

Tokens can be described by patterns

- Pattern: a rule describing the set of lexemes of a particular token type.

Token type	Pattern (informal)	Set of lexemes
INTLITERAL	a string of decimal digits	{0, 1, 324, ...}
ID	a string of letters, digits and underscores, beginning with a letter or underscore	{sum, _sum, facto, KNU_2023, ...}
+	the character “+”	{+}
if	the characters “i”, “f”	{if}

- The pattern is said to match each string in the set.
- We want a formal notation for patterns.
 - allows us to specify the tokens of programming languages.

Tokens can be described by patterns

- Pattern: a rule describing the **set of lexemes** (=strings) of a particular token type.

Pattern name	Pattern	Set of lexemes (strings)
BinaryDigit	0 1	{0, 1}
DecimalDigit	0 1 2 3 4 5 6 7 8 9	{0,1,2,3,4,5,6,7,8,9}
TwoBinaryDigits	(0 1) (0 1)	{00, 01, 10, 11}
TwoBinaryDigits	BinaryDigit BinaryDigit	{00, 01, 10, 11}
Banana	b a (n a)*	{ba, bana, banana, bananaa, ...}
Register	r (0 1 2 3 4 5 6 7 8 9)	{r0, r1, r2, ... r9}

- The set of lexemes (=strings) described by pattern **X** is called the **language of X**.
- We write **L(X)** to denote the language of X.

String Concatenation

- If x and y are strings, then xy is the string formed by appending y to x .
- Examples:

x	y	xy
key	board	keyboard
java	script	javascript

Set Operations (Review)

<i>Operation</i>	<i>Definition</i>
Union of L and M , written $L \cup M$,	$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$
Concatenation of L and M , written LM ,	$LM = \{st \mid s \in L \text{ and } t \in M\}$
Kleene star (or closure) of L , written L^* ,	$L^* = \bigcup_{0 \leq i \leq \infty} L^i$ $L^0 = \{\epsilon\}, \quad L^i = \underbrace{LL\dots L}_{i \text{ times concatenation}}$ <p>ϵ denotes the empty string</p>
Positive Kleene star (or positive closure) of L , written L^+ ,	$L^+ = \bigcup_{1 \leq i \leq \infty} L^i$

Examples: Operations on Languages

- $L = \{a, \dots, z, A, \dots, Z\}$
- $D = \{0, \dots, 9\}$

Example	Language
$L \cup D$	letters and digits
L^3	all 3-letter strings, e.g., abc, xyz, ...
LD	strings consisting of a letter followed by a digit, e.g., a1, b2, x1
L^*	all strings of letters, including the empty string ε
$L(L \cup D)^*$	all strings of letters and digits, starting with a letter, e.g., x, y, a1, aa, c999ccc
D^+	all strings of one or more digits

Regular Expressions (inductive definition)

- **Regular expressions (REs)** can be used to describe patterns.
 - A regular expression r is a pattern that describes a set of lexemes $L(r)$.
 - Lexemes are made from characters of a finite set Σ called **alphabet**.
 - Notation: we underline the characters from Σ .
- Inductive definition of regular expressions over alphabet Σ :
- **R1:** ϵ is a RE denoting the set $\{\epsilon\}$
- **R2:** If a is in Σ , then a is a RE denoting $\{a\}$
- If x and y are REs denoting $L(x)$ and $L(y)$ then
 - **R3:** $x \mid y$ is a RE denoting $L(x) \cup L(y)$
 - **R4:** xy is a RE denoting $L(x)L(y)$
 - **R5:** x^* is a RE denoting $L(x)^*$

Precedence is closure,
then concatenation,
then alternation

called “alternation”

called “concatenation”

called “repetition” or “closure”

Example REs

- $\Sigma = \{\underline{0}, \underline{1}\}$

Rules R2-R5 are from
the previous slide...

RE	Language $L(Re)$
$\underline{1}$	$L(\underline{1}) \stackrel{R2}{=} \{\underline{1}\}$
$\underline{0} \underline{1}$	$L(\underline{0} \underline{1}) \stackrel{R3}{=} L(\underline{0}) \cup L(\underline{1}) \stackrel{R2}{=} \{\underline{0}\} \cup \{\underline{1}\} = \{\underline{0}, \underline{1}\}$
$\underline{1}^*$	$L(\underline{1}^*) \stackrel{R5}{=} L(\underline{1})^* \stackrel{R2}{=} \{\underline{1}\}^* = \{\varepsilon, \underline{1}, \underline{11}, \underline{111}, \dots\}$
$\underline{1}^* \underline{1}$	$L(\underline{1}^* \underline{1}) \stackrel{R4}{=} L(\underline{1}^*)L(\underline{1}) = \dots = \{\underline{1}, \underline{11}, \underline{111}, \dots\}$
$\underline{0} \underline{1}^*\underline{0}$	the set containing $\underline{0}$ and all strings consisting of zero or more $\underline{1}$'s followed by a $\underline{0}$: $\{\underline{0}, \underline{10}, \underline{110}, \dots\}$
$(\underline{111})^*$	the set of strings that contain zero or more sequences of $\underline{111}$: $\{\varepsilon, \underline{111}, \underline{111111}, \underline{11111111}, \dots\}$

Example REs (cont.)

- $\Sigma = \{\underline{a}, \underline{b}, \underline{c}\}$

RE	Language $L(\text{RE})$
$(\underline{ab})^*$	zero or more concatenations of the string \underline{ab} $\{\epsilon, \underline{ab}, \underline{abab}, \underline{ababab}, \dots\}$
$(\underline{a} \mid \underline{b})^*$	zero or more concatenations of $(\underline{a} \mid \underline{b})$ $L((\underline{a} \mid \underline{b})^*) = (L(\underline{a} \mid \underline{b}))^* = (L(\underline{a}) \cup L(\underline{b}))^* = \{\underline{a}, \underline{b}\}^* = \{\epsilon, \underline{a}, \underline{b}, \underline{aa}, \underline{ab}, \underline{ba}, \underline{bb}, \underline{aaa}, \underline{aab}, \underline{aba}, \underline{abb}, \underline{baa}, \underline{bab}, \underline{bba}, \underline{bbb}, \underline{aaaa}, \dots\}$
$(\underline{a} \mid \underline{b})^* \underline{c}$	strings from above, each string concatenated with character \underline{c} $\{\underline{c}, \underline{ac}, \underline{bc}, \underline{aac}, \underline{abc}, \underline{bac}, \underline{bbc}, \underline{aaac}, \underline{aabc}, \dots\}$

Precedence Rules

- Precedence rules are needed to disambiguate regular expressions.
 - Like with arithmetic: $a + b * c = a + (b * c)$
 - Use parenthesis if you want $(a + b) * c$
- Precedence of regular expression operators is **closure**, then **concatenation**, then **alternation**.
- Examples:

$ab^* = a(b^*)$	use parenthesis if you want $(ab)^*$
$ab c = (ab) c$	use parenthesis if you want $a(b c)$

Example Token Specifications using REs

- digit: $\underline{0} \mid \underline{1} \mid \underline{2} \mid \underline{3} \mid \underline{4} \mid \underline{5} \mid \underline{6} \mid \underline{7} \mid \underline{8} \mid \underline{9}$
 - Example digits: $\underline{0}, \underline{1}, \underline{2}, \dots, \underline{9}$
- letter: $\underline{a} \mid \underline{b} \mid \dots \mid \underline{z} \mid \underline{A} \mid \underline{B} \mid \dots \mid \underline{Z}$
 - Example letters: $\underline{a}, \underline{b}, \underline{c}, \dots, \underline{Z}$
- Simple 2-character identifiers, starting with a letter, followed by a letter or digit:
simple_id: letter (letter | digit)
 - Example for simple_id: $\underline{a}\underline{1}, \underline{A}\underline{1}, \underline{a}\underline{b}, \underline{Y}\underline{D}$
- 2-character identifiers, each character can be either a letter or a digit:
twochar_id: (letter | digit) (letter | digit)
 - Example for twochar_id: $\underline{a}\underline{1}, \underline{A}\underline{1}, \underline{1}\underline{A}, \underline{x}\underline{Y}, \underline{1}\underline{2}, \underline{0}\underline{0}$
- What's the problem with **identifiers** consisting of **digits only?**

Notation: digit and letter are defined as symbolic names for REs.

Notation: simple_id uses the symbolic names for digit and letter (like a shorthand notation).

Example Token Specifications (cont.)

- Identifiers with at least one character, the first character must be a letter, the subsequent characters can be letters or digits:
 - **id: letter (letter | digit)***
 - Example for id: KNU, K, Daegu, cse2023, NeWJeaNS
- Integer numbers: first character is a digit, followed by zero or more digits:
 - **integer: digit digit***
 - Examples: 0, 1, 2, 01, 101, 999, 2023
- Signed integer numbers: an integer that can optionally have “+” or “-” in front:
 - **signed_integer: (+ | - | ε) integer**
 - Examples: 0, 1, +1, -1, -1557, +1557

Outline

- The role of a scanner ✓
- Scanner concepts ✓
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata (to be continued...)
 - Definitions of REs (✓), DFAs and NFAs

Writing Regular Expressions

- Write a regular expression for each of the following sets of tokens:
 - Ruby binary literals consisting of “0b” followed by the binary number.
Examples: 0b001011, 0b01.
 - Ruby binary literals, with an optional underscore (“_”) between binary digits.
Examples: 0b0_101, 0b11_01, but not 0b_1 or 0b1_.
 - Ada identifiers: a letter followed by any number of letters, digits, and underlines. An identifier must not end in an underline or have two underlines in a row.
 - A floating-point number: one or more digits (whole-number part) followed by a decimal point (“.”) and one or more digits (fractional part).
Example: 3.14, 0.1234567
 - Floating point number in scientific notation: same as above, but optionally followed by “e” or “E”, and a signed integer exponent.
Example: 1.2e-2, 2.3E+34, 2.3E34.

Shorthand Notations

- One or more concatenations: $r^+ = rr^*$
 - Denotes the language $(L(r))^+$
 - Same precedence and associativity as *
- Zero or one instance: $r? = \epsilon \mid r$
 - Denotes the language $L(r) \cup \{\epsilon\}$
 - Written as $(r)?$ to indicate grouping, e.g., $(12)?$
- Character classes:
 - $[a\underline{z} A\underline{Z}]$

MiniC RE Examples

Token	RE
letter	[a-zA-Z_]
identifier	letter (letter digit)*
integer	digit⁺

- Note that with MiniC, letter includes the underscore character “_”.
- In Java, letters and digits may be drawn from the entire Unicode character set.

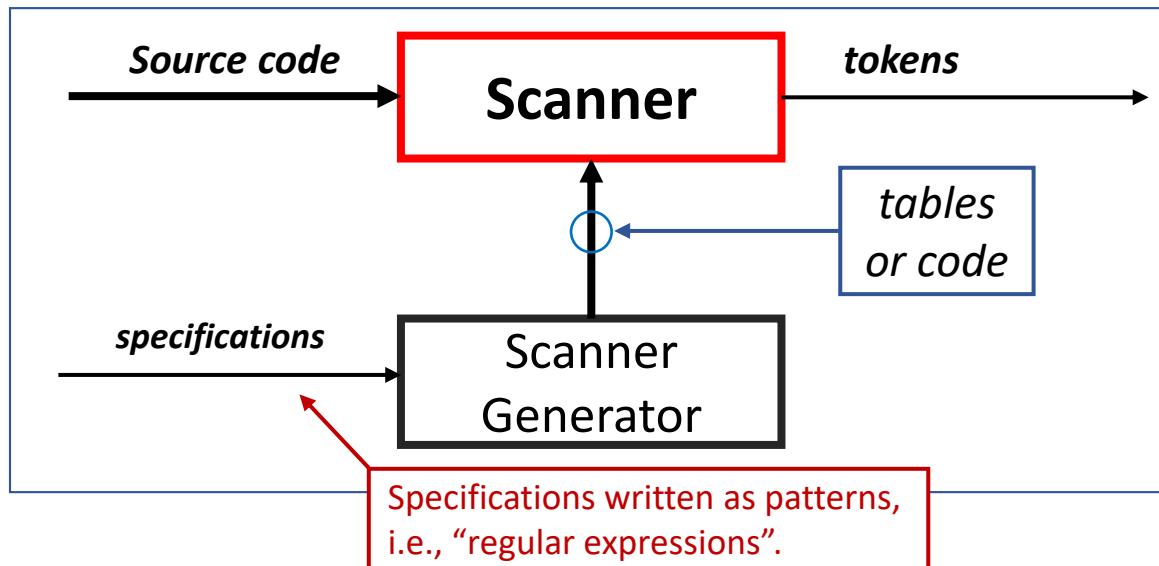
– Examples of Java identifier are

abc リルバ

컴파

The Big Picture

- Why are we doing this?
 - We want to use regular expressions to specify scanners.
 - We want to harness the theory from classes like “Automata”.



- Goals:
 - To simplify specification & implementation of scanners
 - To understand the underlying techniques and technologies

Regular Expressions

- We use regular expressions to specify the mapping of lexemes to tokens.
 - Using results from automata theory of algorithms, we can automatically build recognizers from regular expressions.
- ➔ We study REs and associated theory to automate scanner construction!
- ➔ Fortunately, the automatic techniques produce fast scanners. Used with text editors, URL filtering software, ...

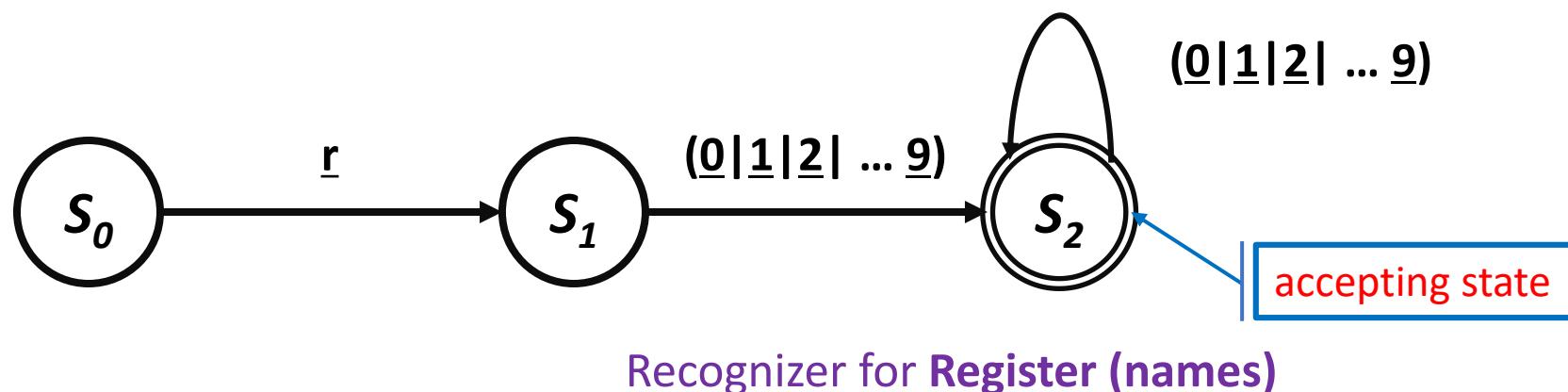
Example

- Consider the problem of recognizing register names

Register: $r (0|1|2| \dots |9) (0|1|2| \dots |9)^*$

- Allows registers of arbitrary number
- Requires at least one digit
- $\Sigma = \{r, 0, 1, \dots, 9\}$

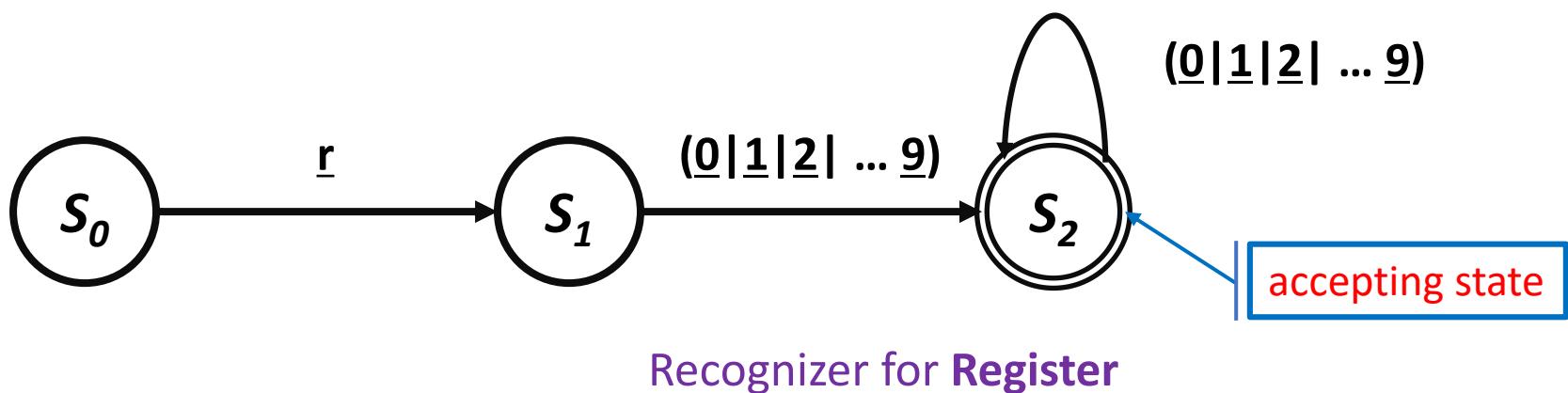
- RE corresponds to a recognizer (or Deterministic Finite Automation, DFA)



Transitions on other inputs go to an error state, s_e

Example (cont.)

- DFA operation
 - Starts in state S_0 and take transitions on each input character
 - DFA **accepts** a word x iff x leaves it in a final state (S_2)



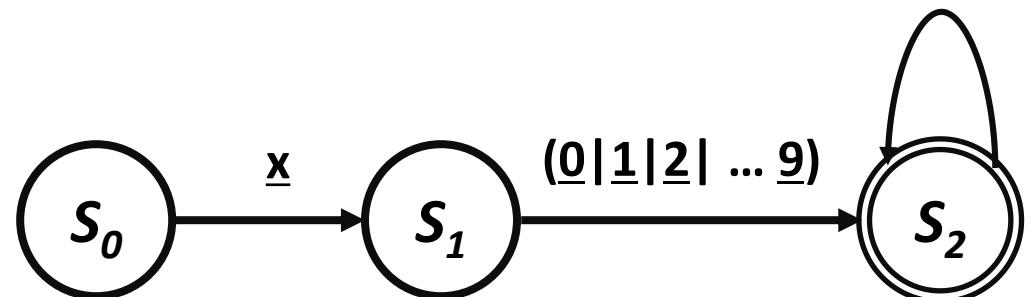
- So,
 - $\underline{r17}$ takes it through s_0 , s_1 , s_2 and accepts
 - \underline{r} takes it through s_0 , s_1 and fails
 - $\underline{1}$ takes it straight to s_e

RISC-V Example

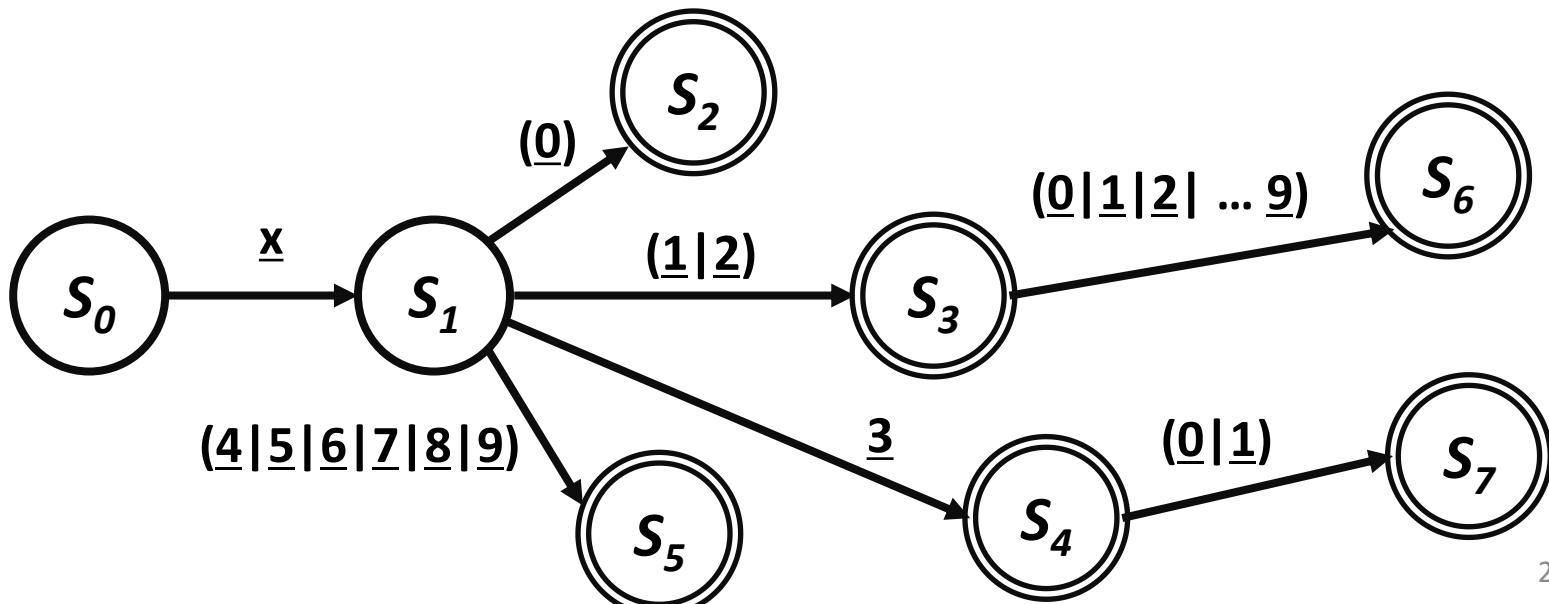
- What about RISC-V architecture's registers?

$(\underline{0}|\underline{1}|\underline{2}| \dots \underline{9})$

Register	ABI Name	Description
x0	zero	Hard-wired zero
x1	ra	Return address
x2	sp	Stack pointer
x3	gp	Global pointer
x4	tp	Thread pointer
x5	t0	Temporary/alternate link register
x6-7	t1-2	Temporaries
x8	s0/fp	Saved register/frame pointer
x9	s1	Saved register
x10-11	a0-1	Function arguments/return values
x12-17	a2-7	Function arguments
x18-27	s2-11	Saved registers
x28-31	t3-6	Temporaries

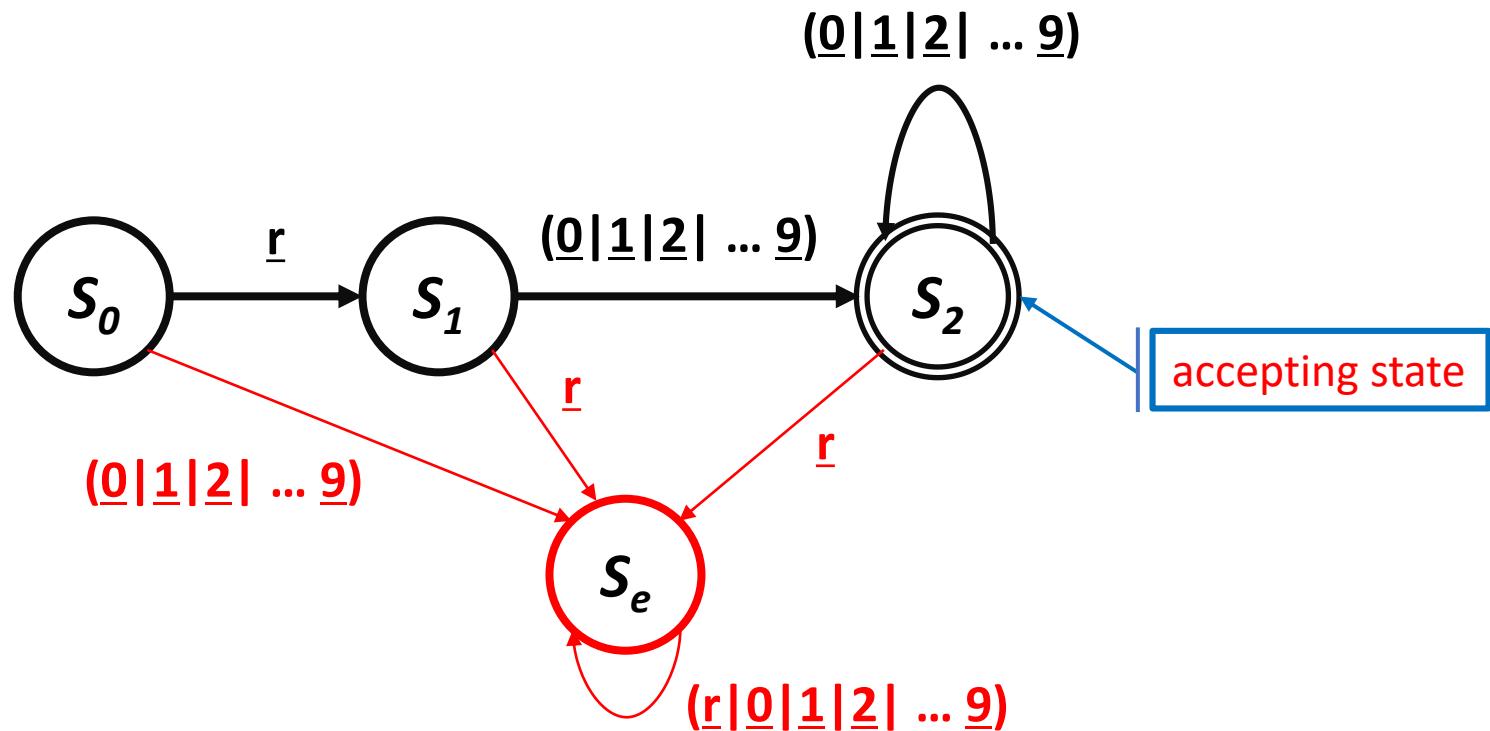


“How about x99??”



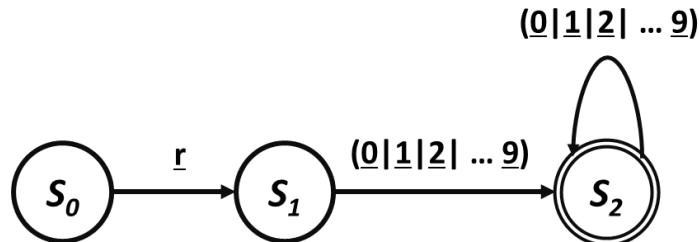
Error State

- Per convention, the error state, s_e , and transitions to it, are not drawn.
- Drawing the error state s_e , and transitions to it would give us:



Remember: $\Sigma = \{r, \underline{0}, \underline{1}, \dots, \underline{9}\}$

Finite Automata (FA)



- A FA consists of a 5-tuple

$$\langle \Sigma, S, \delta, F, I \rangle$$

- where

- Σ is an alphabet
- S is a finite set of states
- δ is a state transition function
- F is a set of **final** or **accepting** states
- I is the **start** state

Example:

$$\Sigma = \{r, 0, 1, \dots, 9\}$$

$$S = \{s_0, s_1, s_2\}$$

δ = takes a state and a symbol as input and returns the next state (see next slide).

$$F = \{s_2\}$$

$$I = \{s_0\}$$

State Transition Function and Code

- To be useful, recognizer must turn into code

δ	<u>r</u>	<u>0,1,2,3,4,5,6,7,8,9</u>	All others
S_0	S_1	S_e	S_e
S_1	S_e	S_2	S_e
S_2	S_e	S_2	S_e
S_e	S_e	S_e	S_e

Table encoding RE

```

Char ← next character
State ←  $S_0$ 

while (Char ≠ EOF)
    State ←  $\delta(\text{State}, \text{Char})$ 
    Char ← next character

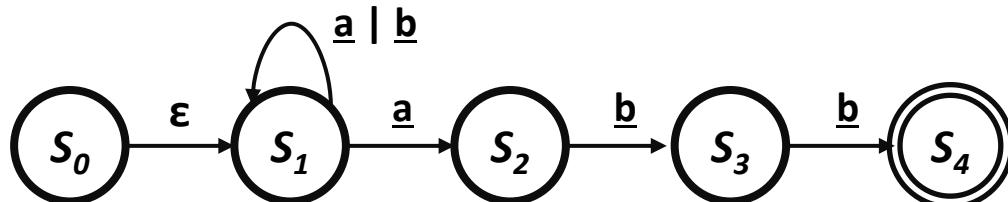
if (State is a final state)
    then report success
    else report failure

```

Skeleton recognizer

Non-deterministic Finite Automata (NFAs)

- We can construct an FA for each RE, but...
- What about an RE such as $(\underline{a} \mid \underline{b})^* \underline{abb}$?



- This is a little different than the FA before:
 - S_0 has a transition on ϵ
 - ϵ does not consume any input!
 - S_1 has **two** transitions on \underline{a}
- This is a non-deterministic finite automaton (NFA)

Problem: state s_1 gives us two choices for character a : staying in s_1 or going to s_2 . We have to guess the correct transition if we want to reach the accepting state s_4 .

* E.g., for string “abb”:
 $s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4$ (accept).

Guessing wrongly causes us to reject a valid string.

* E.g., for string “abb”:
 $s_0 \rightarrow s_1 \rightarrow s_1 \rightarrow s_1 \rightarrow s_1$ (fail).

Non-deterministic Finite Automata

- An NFA accepts a string x iff \exists a path through the transition graph from s_0 to a final state such that the edge labels spell x
- Transitions on ϵ consume no input
- To “run” the NFA, start in s_0 and **guess** the right transition at each step (if there is more than one transition for a given symbol)
 - Always guess correctly
 - If some sequence of correct guesses accepts x then accept
- Why study NFAs?
 - They are the key to automating the RE \rightarrow DFA construction
 - We can glue together NFAs with ϵ -transitions



Relationship between NFAs and DFAs

- DFA is a special case of an NFA
 - DFA has no ϵ transitions
 - DFA's transition function is single-valued
 - not possible to have 2 transitions from state s on a symbol a
 - Same rules will work
- DFA can be simulated with an NFA
 - *Obviously*
- NFA can be simulated with a DFA (*less obvious*)
 - Simulate sets of possible states
 - Possible exponential blowup in the state space
 - Still, one state per character in the input stream

Summary: NFAs and DFAs

- A FA is a **DFA** if
 - no state has an **ϵ -transition**, i.e., there is no transition on input ϵ
 - for each state **s** and input symbol **a**, there is **at most** one edge labeled **a** leaving state **s**.
- A FA is an **NFA**
 - if it contains **ϵ -transitions** or
 - if it has **several possible transitions** from a state **s** on a given input symbol **a**.

Automating Scanner Construction

- To convert a specification to code:
 - 1) Write down the RE for the input language
 - 2) Build a big NFA
 - 3) Build the DFA that simulates the NFA
 - 4) Minimize the number of states in the DFA
 - 5) Generate the scanner code
- Scanner generators
 - Lex and Flex work along these lines
 - Algorithms are well-known and well-understood
 - Key issue is interface to parser
 - You could build one in a weekend!

Outline

- The role of a scanner ✓
- Scanner concepts ✓
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata (to be continued...)
 - Definitions of REs, DFAs and NFAs ✓
 - REs→NFA
(Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)
 - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)
 - DFA→minimal-state DFA
(state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
- Scanner generators

Automating Scanner Construction

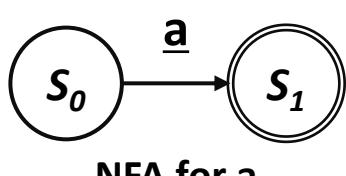
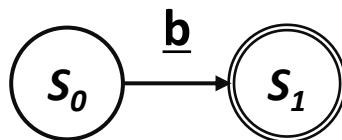
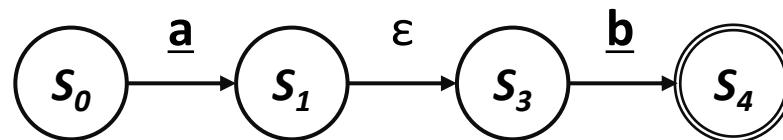
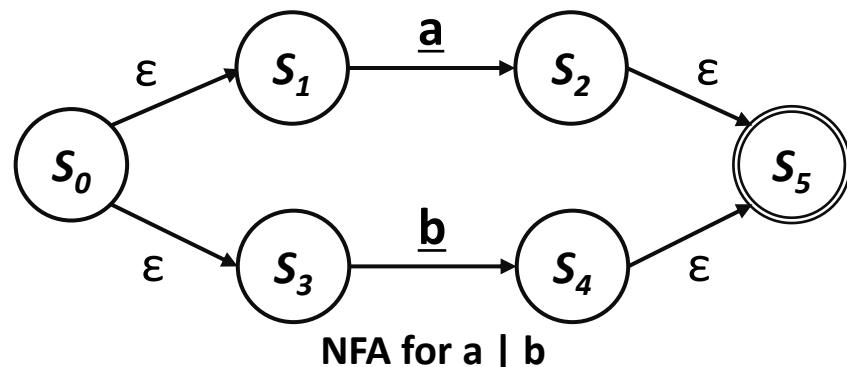
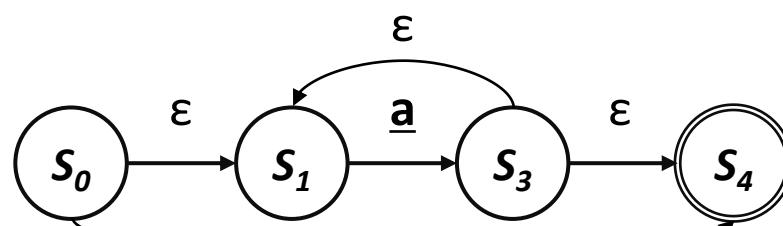
- **RE \rightarrow NFA** (*Thompson's construction*)
 - Build an NFA for ϵ and each symbol $a \in \Sigma$ occurring in the RE
 - Combine them with ϵ -moves
- **NFA \rightarrow DFA** (*subset construction*)
 - Build the simulation of the NFA
- **DFA \rightarrow Minimal DFA**
 - Hopcroft's algorithm
- **DFA \rightarrow RE** (*Not part of the scanner construction*)
 - All pairs, all paths problem
 - Take the union of all paths from s_0 to an accepting state

The Cycle of Constructions



RE → NFA using Thompson's Construction

- Key idea:
 - construct NFA for each symbol and each operator of the RE
 - Join them with ϵ moves in precedence order of RE operators

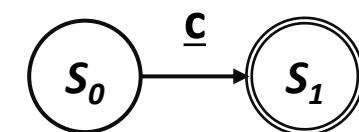
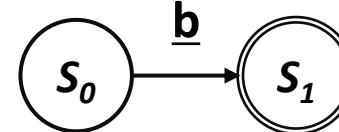
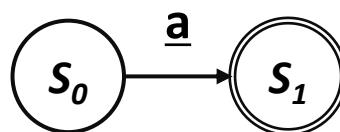
NFA for aNFA for bNFA for abNFA for a | bNFA for a^{*}

Note: when joining two NFAs, we renumber states so that s_0 is again the start state, and all state names are unique.

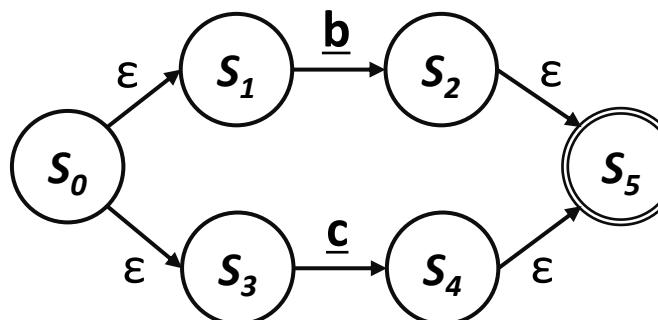
Example of Thompson's Construction

- Let's try $\underline{a} (\underline{b} | \underline{c})^*$

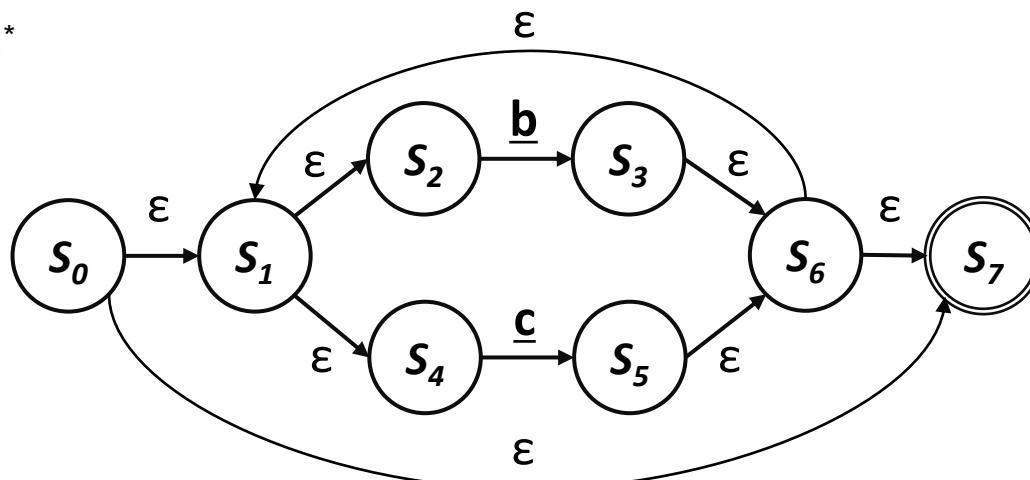
- 1. \underline{a} , \underline{b} , and \underline{c} :



- 2. $\underline{b} | \underline{c}$

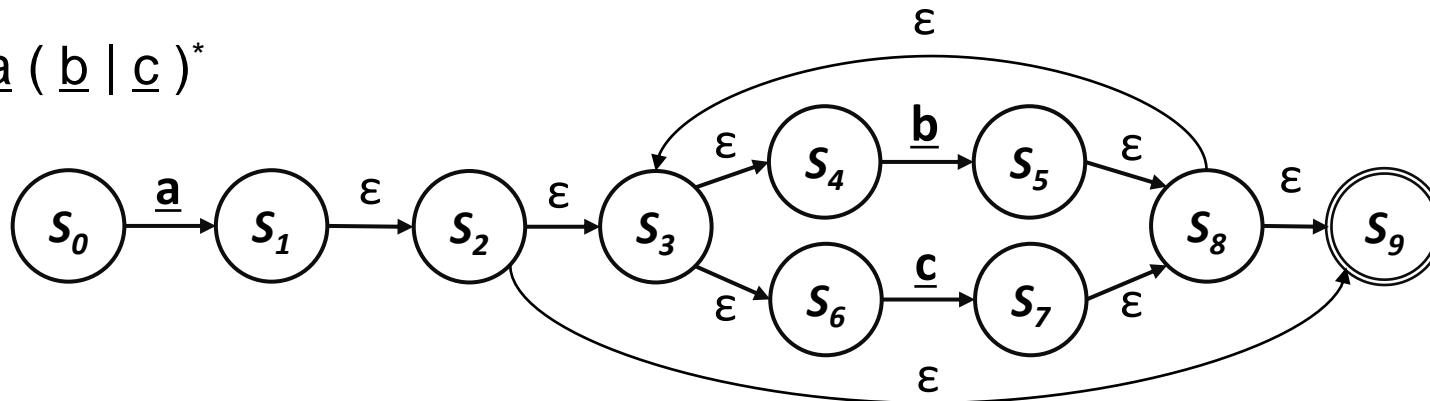


- 3. $(\underline{b} | \underline{c})^*$

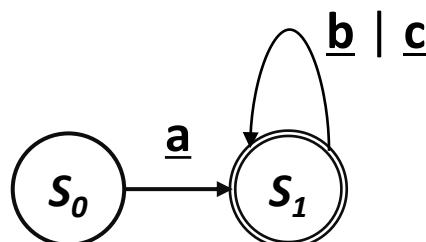


Example of Thompson's Construction (cont.)

- 4. a (b | c)^{*}



Of course, a human would design something simpler ...



But, we can automate the production of the more complex automaton ...

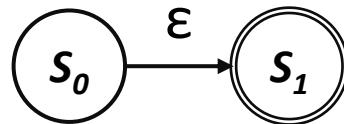
Thompson's Construction

- Syntax-driven
 - follows the structure of REs
- **Inductive:** the cases in the construction of the NFA follow the cases in the definition of REs.
- **Important:** if a symbol a occurs several times in a RE r, then a separate NFA is constructed for each occurrence of a.

Thompson's Construction

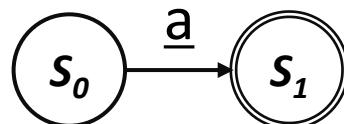
- **Inductive Base:**

- For ϵ , construct the NFA



Thompson's Construction for automata follows the inductive definition of regular expressions from page 68.

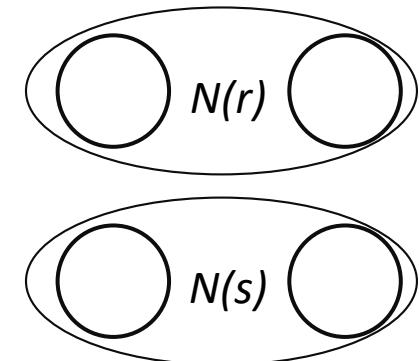
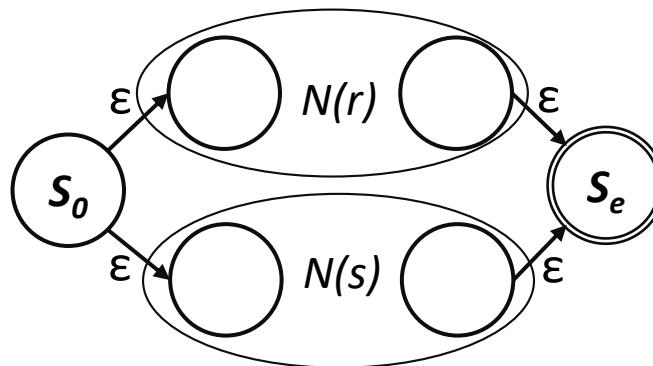
- For $a \in \Sigma$, construct the NFA



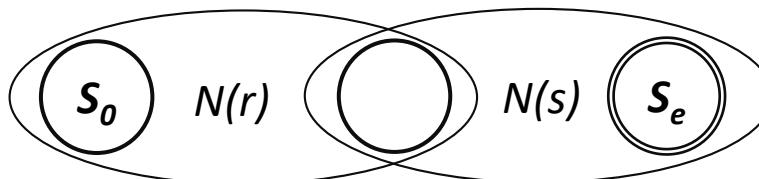
- **Inductive step:** suppose $N(r)$ and $N(s)$ are NFAs for REs r and s .
 - Then... (continued on next slide).

Thompson's Construction

- RE $r \mid s$:

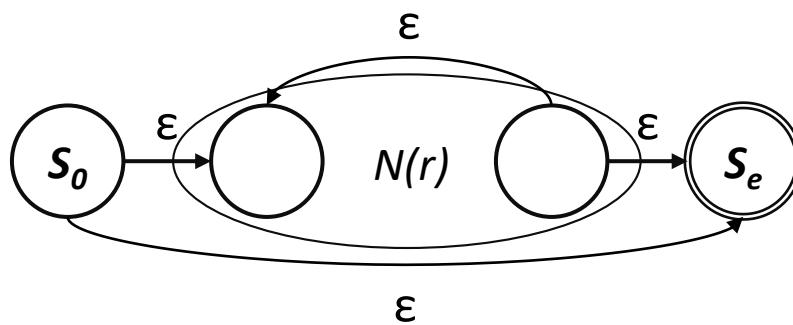


- RE $r s$:



Thompson's Construction for automata follows the inductive definition of regular expressions from page 68.

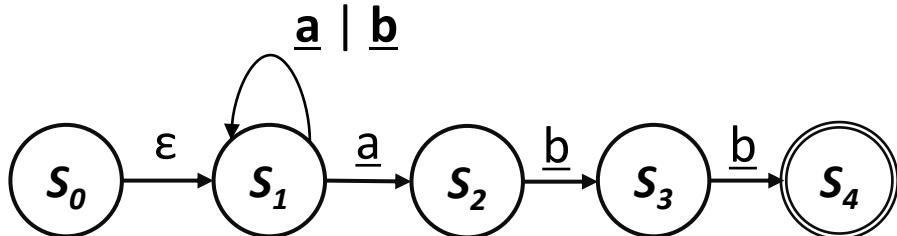
- RE r^* :



Outline

- The role of a scanner ✓
- Scanner concepts ✓
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata (to be continued...)
 - Definitions of REs, DFAs and NFAs ✓
 - REs→NFA ✓
(Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)
 - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)
 - DFA→minimal-state DFA
(state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
- Scanner generators

Subset Construction (NFA → DFA)



Instead of guessing in state s_1 on symbol a, we can follow both transitions ($s_1 \rightarrow s_1$ and $s_1 \rightarrow s_2$) in **parallel**.

- We introduce a “virtual state” that combines the states that we reach from s_1 on symbol a.
- The new virtual state contains states s_1 and s_2 , we write $\{s_1, s_2\}$ for the virtual state.
- **A question:** if we are in the virtual state $\{s_1, s_2\}$, where can we go when we read symbol b?
- **Answer:** the virtual state takes us anywhere that one of its member states takes us on symbol b (either s_1 or s_3 in the above NFA). So we introduce a second virtual state $\{s_1, s_3\}$.
(continued on next slide.)

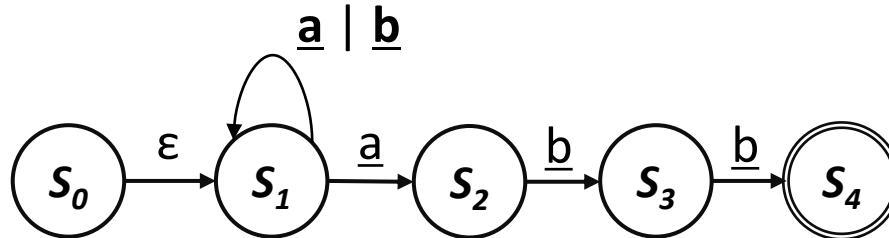
Problem: state s_1 gives us two choices for character a: staying in s_1 or going to s_2 . We have to guess the correct transition if we want to reach the accepting state s_4 .

E.g., for string “abb”: $s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4$ (accept).

Guessing wrongly causes us to reject a valid string.

E.g., for string “abb”: $s_0 \rightarrow s_1 \rightarrow s_1 \rightarrow s_1$ (fail).

Subset Construction (NFA \rightarrow DFA)



- A “virtual state” is a **subset** of the set of states $S=\{S_0, S_1, S_2, S_3, S_4\}$ of the NFA.

δ	<u>a</u>	<u>b</u>
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	$\{s_1, s_4\}$
$\{s_1, s_4\}$	$\{s_1, s_2\}$	$\{s_1\}$

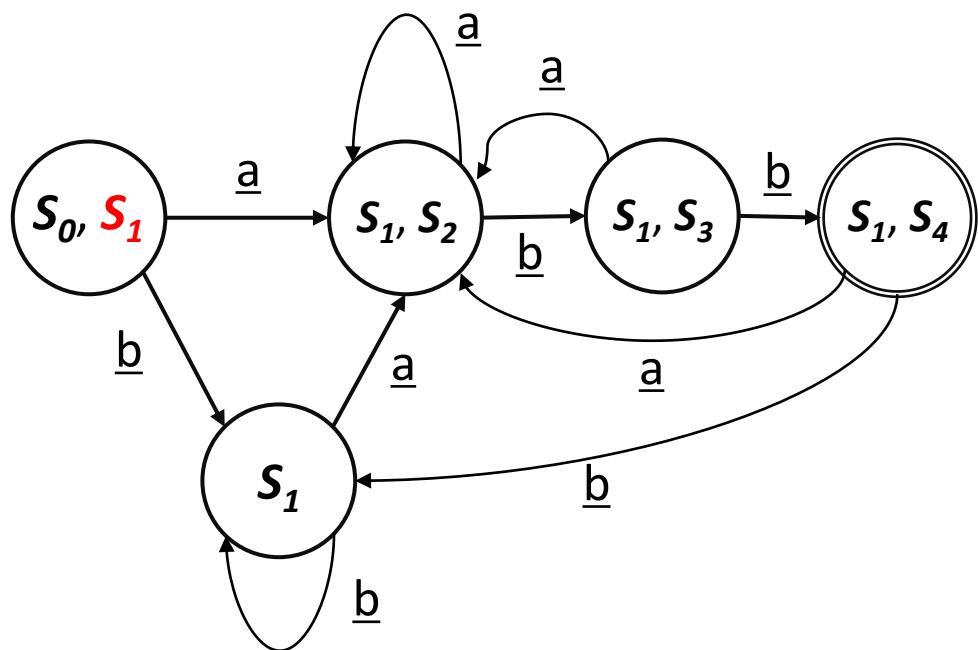


Table encoding DFA

Algorithm: NFA→DFA with Subset Construction

We need to build a simulation of the NFA

Two key functions

- $\text{Move}(s_i, a)$: gives set of states reachable from set s_i by a
- $\varepsilon\text{-closure}(s_i)$: gives set of states reachable from set s_i by ε

The algorithm:

- Start state derived from s_0 of the NFA:
 - Take its ε -closure $S_0 = \varepsilon\text{-closure}(\{s_0\})$
- Take the image of S_0 , $\text{Move}(S_0, a)$, for each $a \in \Sigma$, and take its ε -closure
- Iterate until no more states are added

Sounds more complex than it is...

Algorithm: NFA→DFA with Subset Construction

The algorithm:

```

 $Dstates \leftarrow \{ \};$ 
add  $\epsilon$ -closure( $s_0$ ) as an
    unmarked state to  $Dstates$ ;
while ( there is an unmarked
    state  $T$  in  $Dstates$  ) {
    mark  $T$ ;
    for each  $\alpha \in \Sigma$  {
         $U \leftarrow \epsilon$ -closure( $Move(T, \alpha)$ )
        if (  $U \notin Dstates$  ) then
            add  $U$  as an unmarked
                state to  $Dstates$ ;
         $\delta [T, \alpha] \leftarrow U$ 
    }
}

```

Let's think about why this works

The algorithm halts:

1. $Dstates$ contains no duplicates
(test before adding)
2. $2^{|S|}$ is finite
3. while loop adds to $Dstates$, but does
not remove from $Dstates$ (*monotone*)
→ the loop halts

$Dstates$ contains all the reachable NFA
states:

It tries each character $\alpha \in \Sigma$ in each
state T .

**It builds every possible NFA
configuration.**

→ $Dstates$ and δ form the DFA

Any DFA state containing a final state from the
NFA becomes a final state in the DFA.

NFA→DFA with Subset Construction

Example of a fixed-point computation

- Monotone construction of some finite set
- Halts when it stops adding to the set
- Proofs of halting & correctness are similar
- These computations arise in many contexts

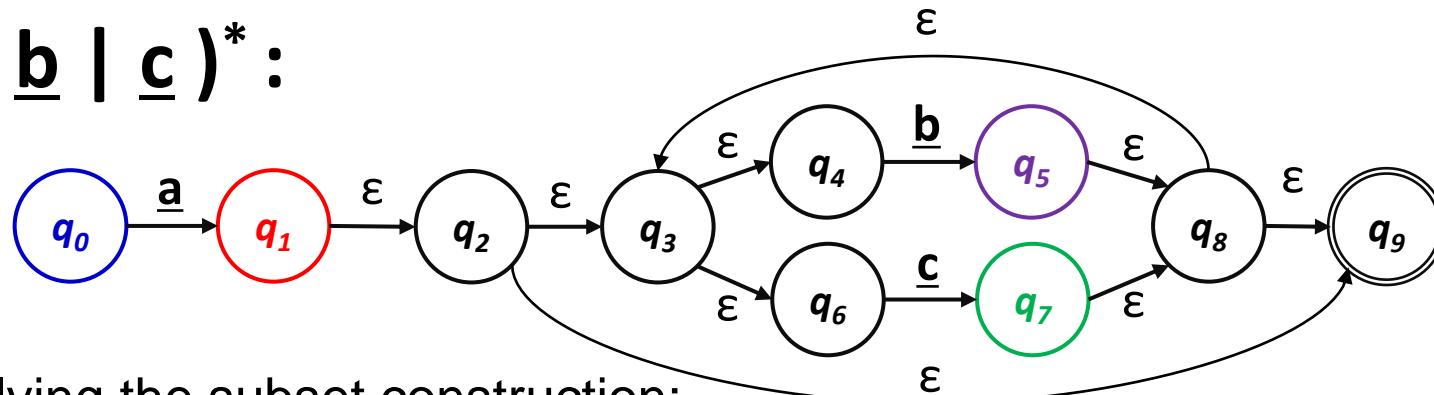
Other fixed-point computations

- Classic data-flow analysis (& Gaussian Elimination)
 - Solving sets of simultaneous set equations

We will see many more fixed-point computations

Example: NFA → DFA with Subset Construction

a (b | c)^{*}:

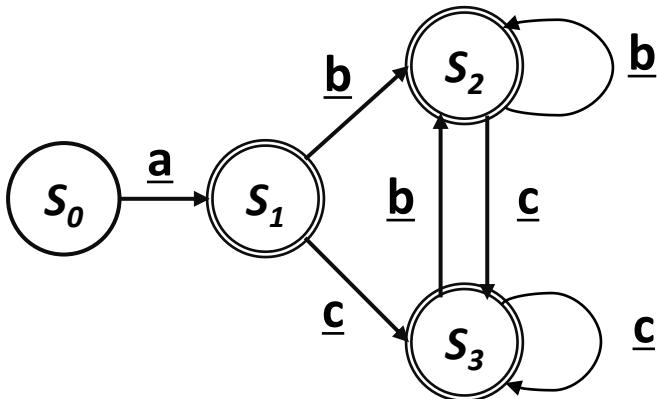


- Applying the subset construction:

	NFA states	ε-closure (Move (s, *))		
		<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$	none	s_2	s_3
s_3	$q_7, q_8, q_9, q_3, q_4, q_6$	none	s_2	s_3
Final states				

Example: NFA → DFA with Subset Construction

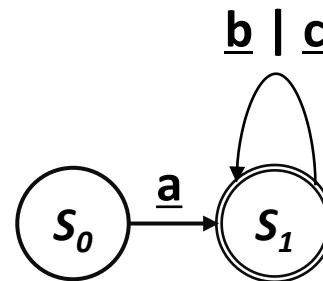
- The DFA for $a (\underline{b} \mid \underline{c})^*$



δ	<u>a</u>	<u>b</u>	<u>c</u>
S_0	S_1	-	-
S_1	-	S_2	S_3
S_2	-	S_2	S_3
S_3	-	S_2	S_3

- Ends up smaller than the NFA
- All transitions are deterministic
- Use the same code skeleton as before

But, remember our goal:



Outline

- The role of a scanner ✓
- Scanner concepts ✓
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata (to be continued...)
 - Definitions of REs, DFAs and NFAs ✓
 - REs→NFA ✓
(Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)
 - NFA→DFA ✓
(subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)
 - DFA→minimal-state DFA
(state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
- Scanner generators

DFA Minimization Overview

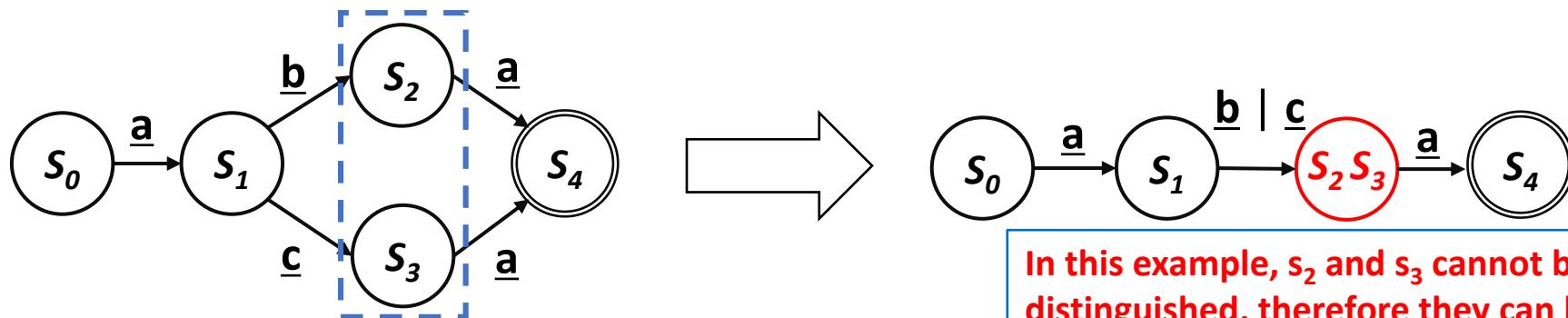
The Big Picture:

- Discover **distinguishable** states
- States that cannot be distinguished can be represented by a **single** state.

Two states p and q are **indistinguishable** by input string w iff the DFA starting in p accepts w and starting in q accepts w, and the DFA starting in p does not accept w and starting in q does not accept w.

Two states p and q are **distinguishable by input string w** iff the DFA starting in p accepts w but starting in q does not accept w.

Two states p and q are **distinguishable** if they are distinguishable by some input string w.



In this example, s_2 and s_3 cannot be distinguished, therefore they can be represented by a single state s_2s_3 .

DFA Minimization Overview

- The set of states is divided into subsets of states that cannot be distinguished.
- We say that the set of states S is partitioned into **partition P** :
 - Each state $s \in S$ is in exactly one set $p_i \in P$
 - States in the same set have not been distinguished yet.
 - States from different sets are known to be distinguishable.

Step 1:

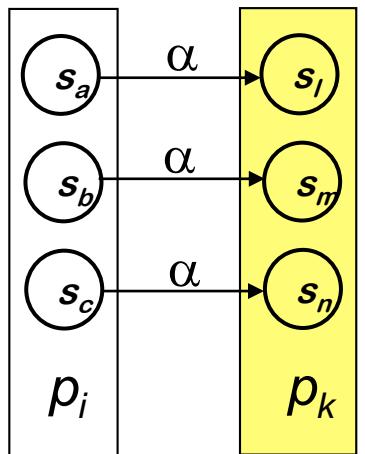
- **Initially** the partition P consists of **2 sets**:
 - the accepting states F and the non-accepting states $S - F$.
 - States from F and $S-F$ are distinguishable by ϵ .

Step 2:

- The minimization algorithm repeatedly picks a set $p_i \in P$ and tries to distinguish between its states by some symbol $\alpha \in \Sigma$.
→ split p_i along α (see next slide).

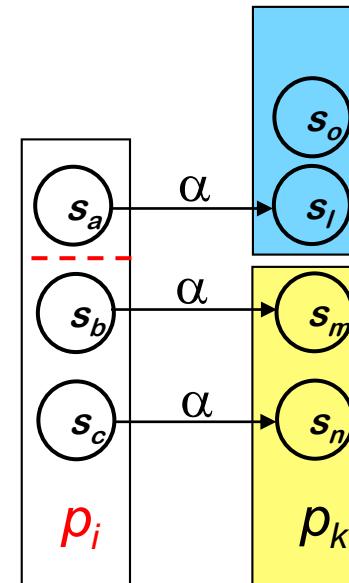
Splitting of a State Set along Symbol α

- The minimization algorithm repeatedly picks a set $p_i \in P$ and tries to distinguish between its states by some symbol $\alpha \in \Sigma$:



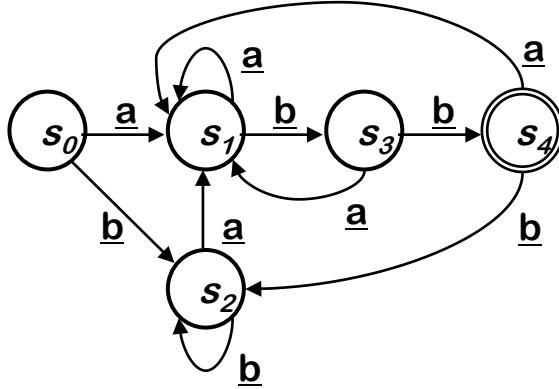
α does not split p_i

α does not split p_i



α splits p_i into $\{s_a\}$ and $\{s_b, s_c\}$
 $\{s_a\}$ and $\{s_b, s_c\}$ are distinguishable by α ...

- Eventually no more set can be split and the algorithm terminates.



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

Start off with **2 subsets of S: F and S-F**
($p \in 2^S$ (powerset))

```

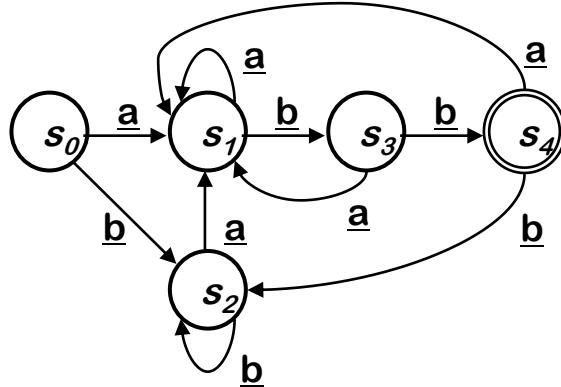
 $P \leftarrow \{ F, S-F \}$  ←
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split ( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

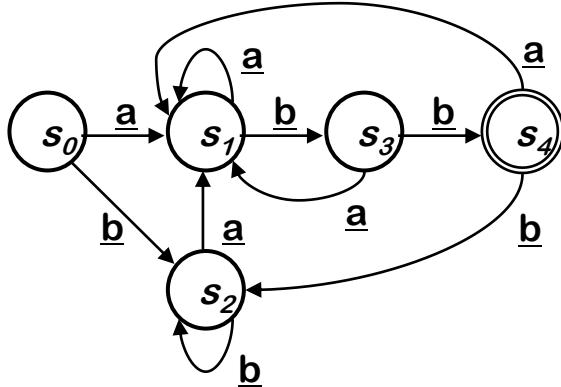
$$T = \{ \}$$

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 
    
```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;
    
```



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$T = \{ \}$$

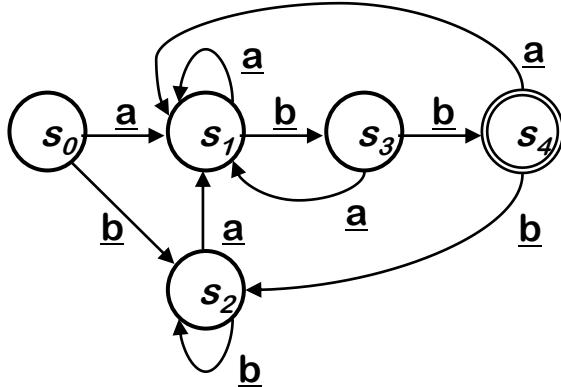
$$p = \{s_4\}$$

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$  ←
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$  ←
     $P \leftarrow T$ 

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

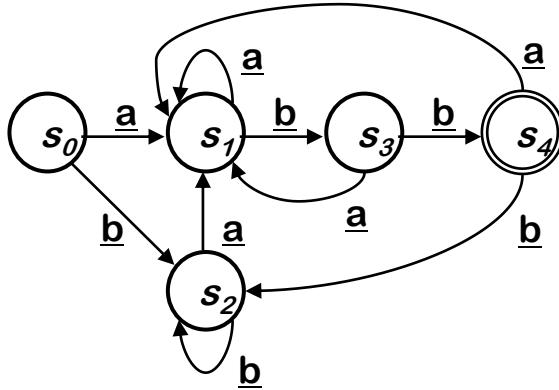
$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

- Invoke Split(p) function:
 - Determines whether p can be split into p_1 and p_2
 - How?
 - If these two states can be distinguishable by a symbol α ...

If X and Y are two states in a DFA,
we can combine these two states into $\{X, Y\}$
when they are not distinguishable.



Minimization of DFA

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

Split( $p$ )
    for each  $\alpha \in \Sigma$  ←
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

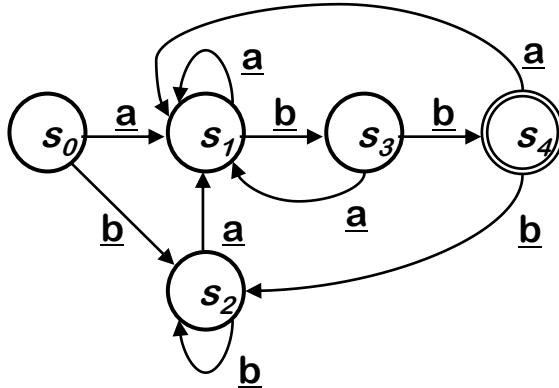
$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

$$\alpha = \underline{a}$$

If X and Y are two states in a DFA,
we can combine these two states into {X, Y}
when they are not distinguishable.



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

}

return p ;

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

$$\alpha = \underline{a}$$

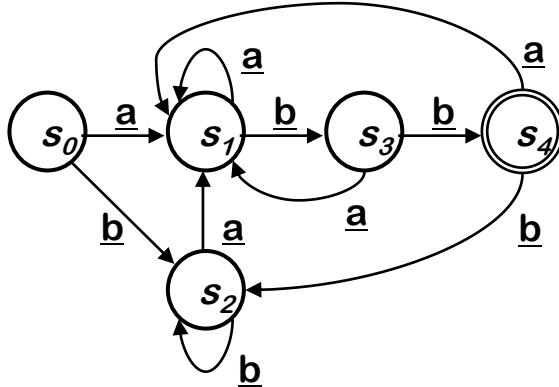
- Does \underline{a} can split $p = \{s_4\}$ to p_1 and p_2 ?

↔ Are there any two subsets can be distinguishable in set p by symbol \underline{a} ?

→ Nope, not this time.

(we cannot split a set which has a single element...)





Minimization of DFA

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

    Split( $p$ )
        for each  $\alpha \in \Sigma$  ↪
        {
            if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
            then return  $\{p_1, p_2\}$ 
        }
        return  $p$ ;
    
```

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

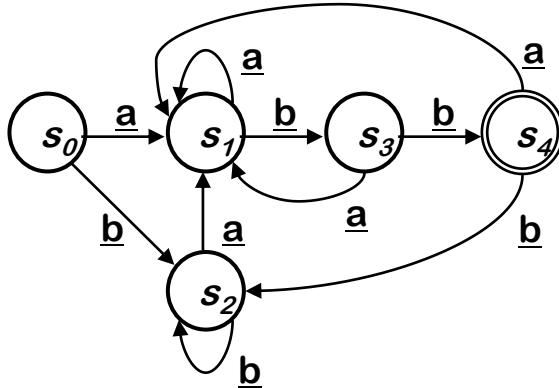
$$\alpha = \underline{b}$$

- Does \underline{b} can split $p = \{s_4\}$ to p_1 and p_2 ?

↔ Are there any two subsets can be distinguishable in set p by symbol \underline{b} ?

→ Nope.

(we cannot split a set which has a single element...)



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

$$\alpha =$$

```

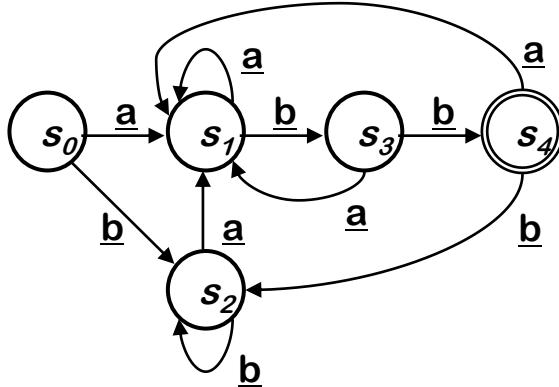
 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ; 

```



Minimization of DFA

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$  ←
     $P \leftarrow T$ 

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

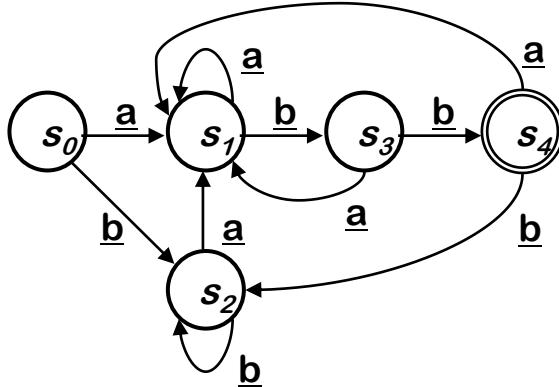
$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \}$$

$$p = \{s_4\}$$

$$\alpha =$$

$$T = \{ \} \cup \text{Split}(\{s_4\}) = \{s_4\}$$



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \{s_4\} \}$$

$$p = \{s_4\}$$

$$\alpha =$$

```

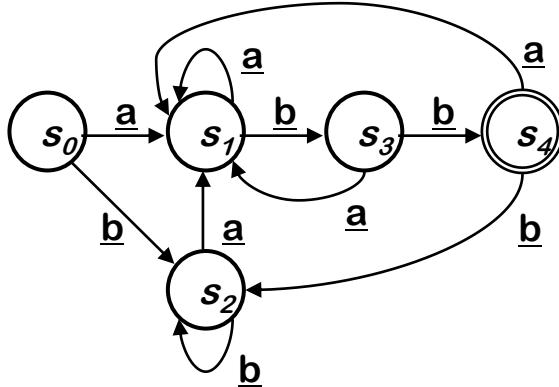
 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$  ←
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha =$$

```

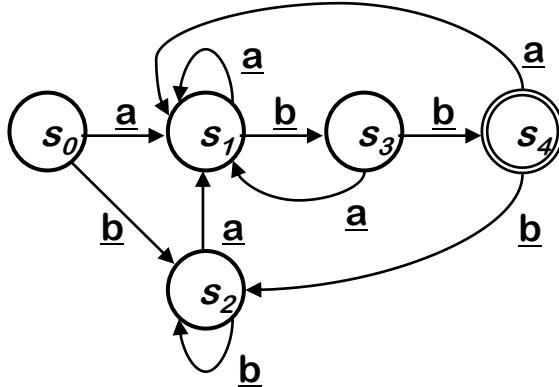
 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$  ←
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha =$$

```

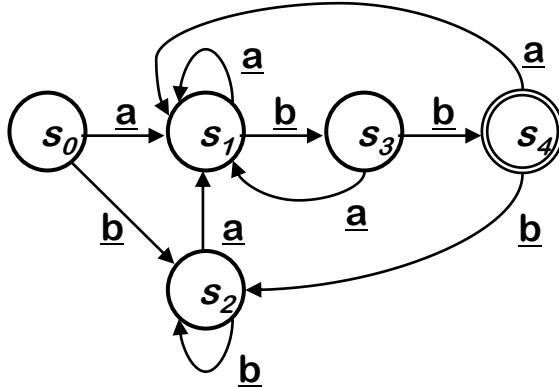
 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$  ←
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha = \underline{a}$$

```

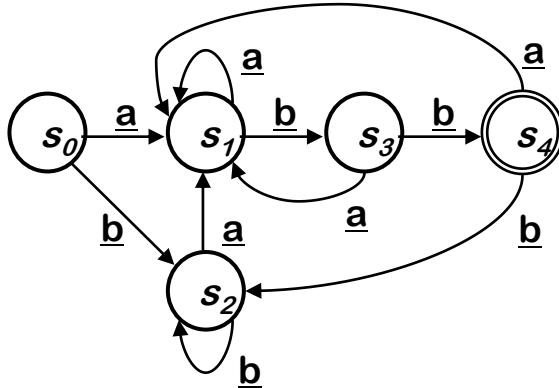
 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$  ←
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

 }

return p ;

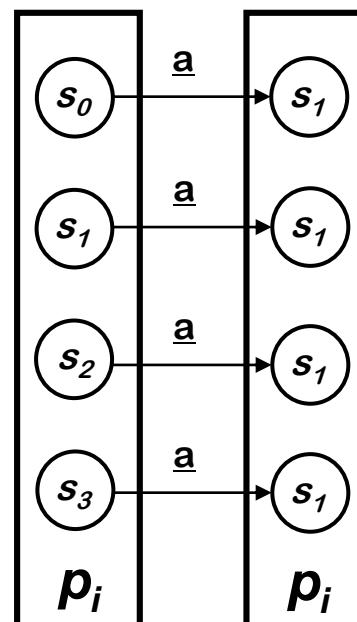
$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

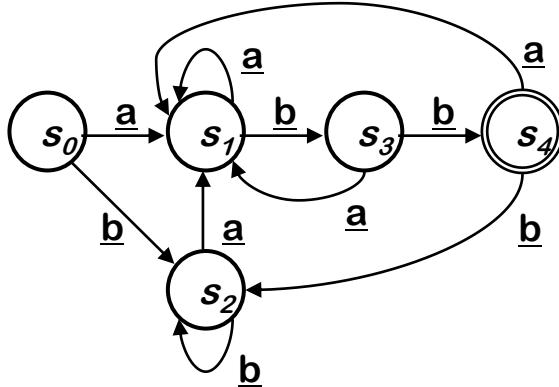
$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha = \underline{a}$$



- Split on \underline{a} :
 - None...



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

 }

return p ;

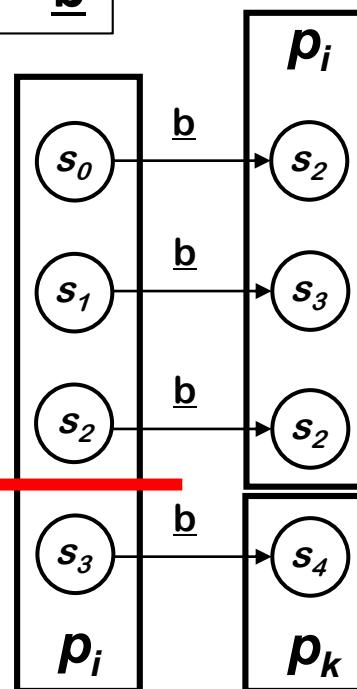
$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \{s_4\} \}$$

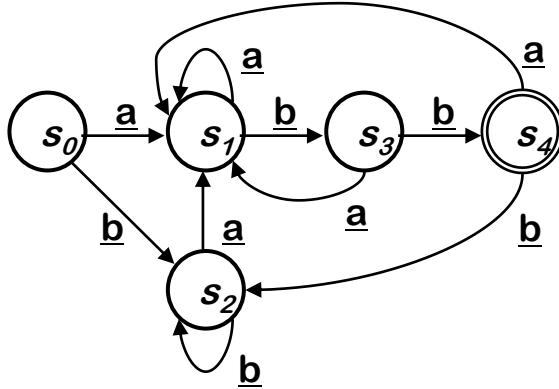
$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha = \underline{b}$$



- Split on \underline{b} :

- α splits p_i into $\{s_0, s_1, s_2\}$ and $\{s_3\}$



```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$  ←
    }
    return  $p$ ;

```

Minimization of DFA

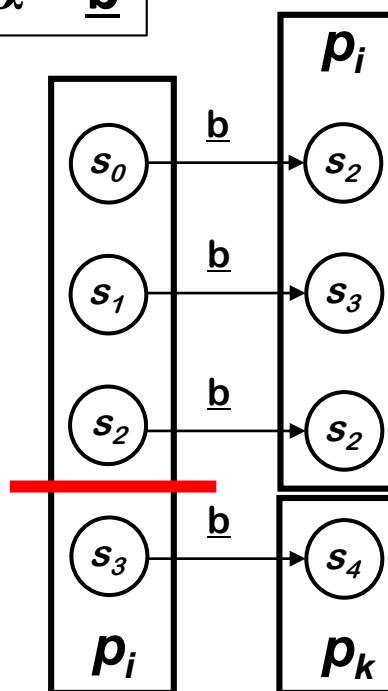
$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

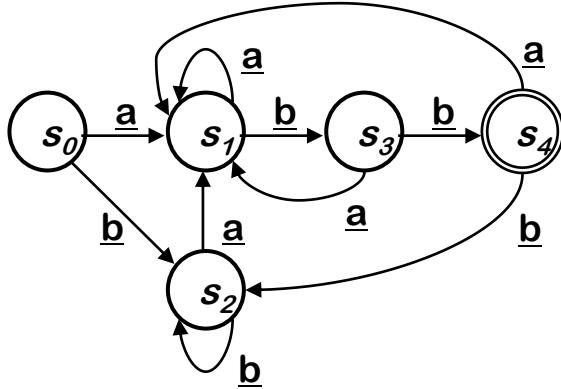
$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha = \underline{b}$$



- **Return**
 $\{\{s_0, s_1, s_2\}, \{s_3\}\}$



Minimization of DFA

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$   $\leftarrow$ 
     $P \leftarrow T$ 

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

$$\Sigma = \{a, \underline{b}\}$$

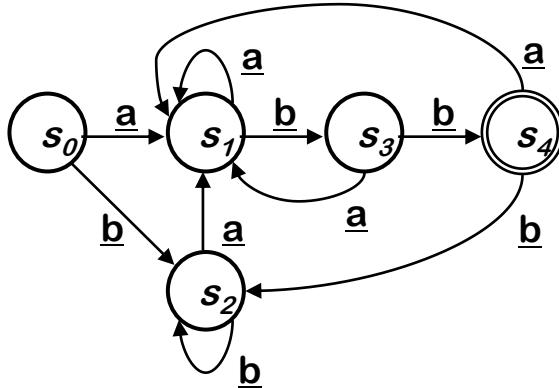
$$T = \{ \{s_4\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha =$$

- **Split(p)** = $\{\{s_0, s_1, s_2\}, \{s_3\}\}$

$$T = \{ \{s_0, s_1, s_2\}, \{s_3\} \}$$



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$ ←

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into $p1$ and $p2$

then return $\{p1, p2\}$

 }

return p ;

$$P_0 = \{ \{s_4\}, \{s_0, s_1, s_2, s_3\} \}$$

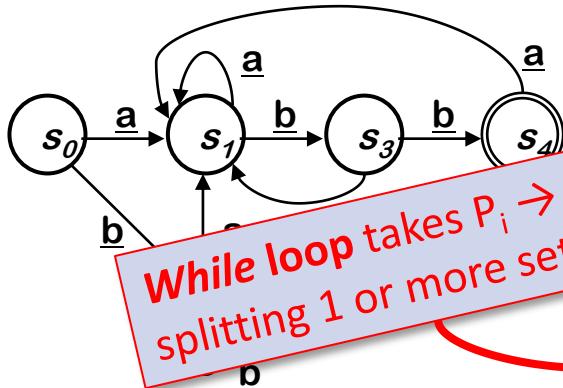
$$\Sigma = \{a, \underline{b}\}$$

$$T = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\alpha =$$

Minimization of DFA



While loop takes $P_i \rightarrow P_{i+1}$ by splitting 1 or more sets

$$P \leftarrow \{ F, S-F \}$$

while (P is still changing)

$$T \leftarrow \{ \}$$

for each set $p \in P$

$$T \leftarrow T \cup \text{Split}(p)$$

$$P \leftarrow T$$

Split (p)

for each $\alpha \in \Sigma$

{

if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

}

return p ;

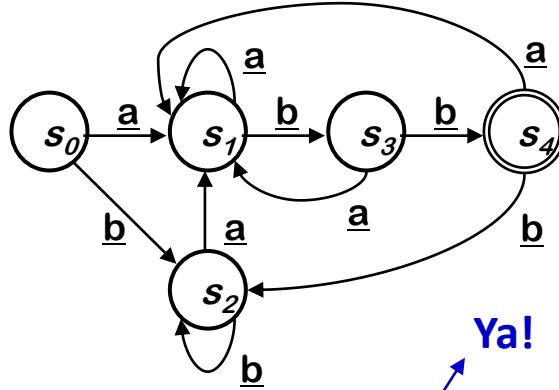
$$P_1 = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$T = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$p = \{s_0, s_1, s_2, s_3\}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$\alpha =$$



```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )  $\leftarrow$ 
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

Minimization of DFA

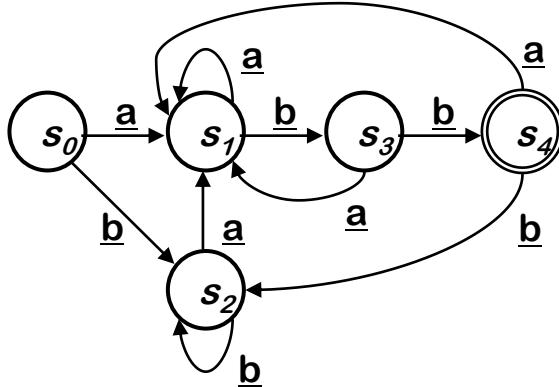
$$P_1 = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$T = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$p = \{ \}$$

$$\alpha =$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$ ←

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

}

return p ;

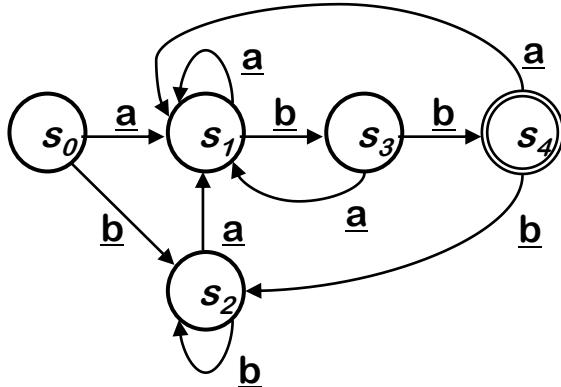
$$P_1 = \{ \{s_4\}, \{s_0, s_1, s_2\}, \{s_3\} \}$$

$$T = \{ \}$$

$$p = \{ \}$$

$$\alpha =$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$



Minimization of DFA

$$P_1 = \{ \{s_0, s_1, s_2\}, \{s_3\}, \{s_4\} \}$$

Let's sort it by its subscript for readability!

$$T = \{ \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

```

P ← { F, S-F }
while ( P is still changing )
  T ← {}
  for each set p ∈ P
    T ← T ∪ Split (p)
  P ← T
  ↙
Split (p)
  for each α ∈ Σ
  {
    if α splits p into p1 and p2
    then return {p1, p2}
  }
  return p;
  
```

- **Split(p)**

$$\alpha = \underline{a}$$

$$\alpha = \underline{b}$$

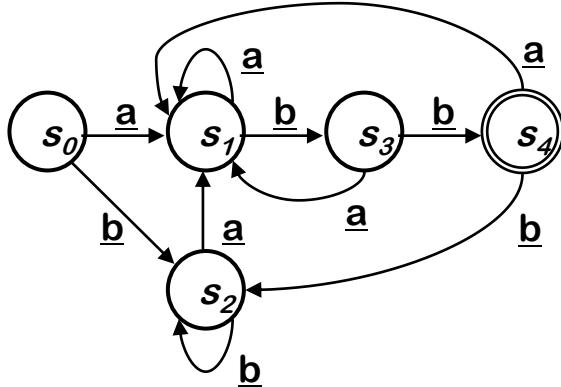
$$p = \{s_0, s_1, s_2\}$$

none

$$\{s_0, s_2\}, \{s_1\}$$

$$p = \{s_3\}$$

$$p = \{s_4\}$$



Minimization of DFA

$$P_1 = \{ \{s_0, s_1, s_2\}, \{s_3\}, \{s_4\} \}$$

$$T = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$p = \{ \}$$

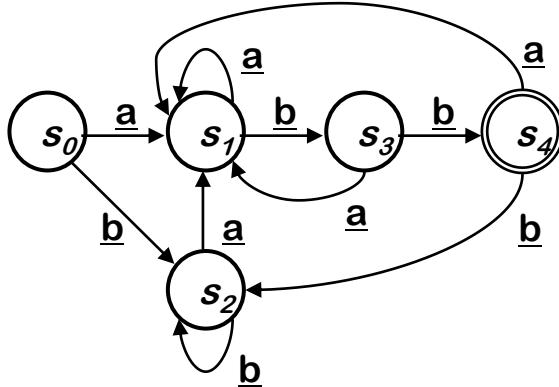
```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
   $T \leftarrow \{ \}$ 
  for each set  $p \in P$ 
     $T \leftarrow T \cup \text{Split}(p)$ 
   $P \leftarrow T$ 
  
```

←

```

Split( $p$ )
  for each  $\alpha \in \Sigma$ 
  {
    if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
    then return  $\{p1, p2\}$ 
  }
  return  $p$ ;
  
```



Minimization of DFA

$$P_2 = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$T = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

$$p = \{ \}$$

```

P ← { F, S-F }
while ( P is still changing )
  T ← {}
  for each set p ∈ P
    T ← T ∪ Split (p)
  P ← T
  
```

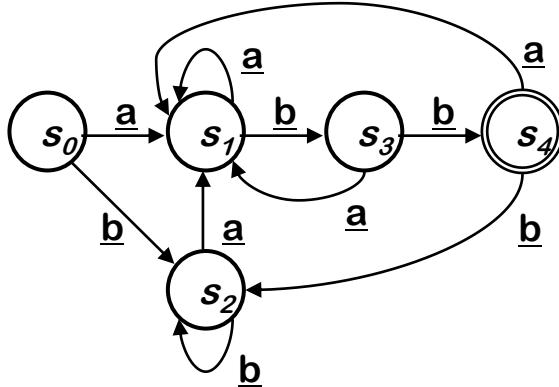


```

    for each α ∈ Σ
    {
      if α splits p into p1 and p2
      then return {p1, p2}
    }
    return p;
  
```

```

Split (p)
  for each α ∈ Σ
  {
    if α splits p into p1 and p2
    then return {p1, p2}
  }
  return p;
  
```



Minimization of DFA

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

if α splits p into $p1$ and $p2$

then return $\{p1, p2\}$

}

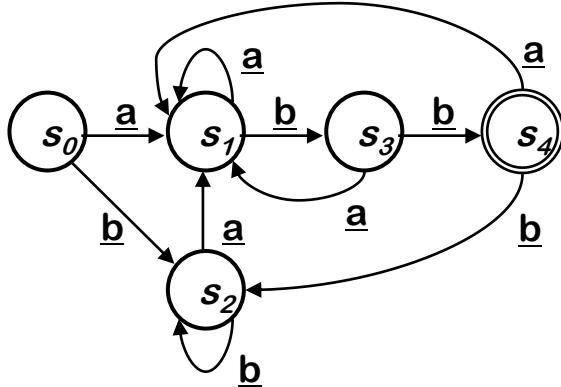
return p ;

$$P_2 = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$T = \{ \}$$

$$p = \underbrace{\{s_0, s_2\}}_{1^{\text{st iter}}}, \underbrace{\{s_1\}}_{2^{\text{nd}}}, \underbrace{\{s_3\}}_{3^{\text{rd}}}, \underbrace{\{s_4\}}_{4^{\text{th}}}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$



Minimization of DFA

$$P_2 = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$T = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$p = \underbrace{\{s_0, s_2\}}_{1^{\text{st}} \text{ iter}}, \underbrace{\{s_1\}}_{2^{\text{nd}}}, \underbrace{\{s_3\}}_{3^{\text{rd}}}, \underbrace{\{s_4\}}_{4^{\text{th}}}$$

$$\Sigma = \{\underline{a}, \underline{b}\}$$

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

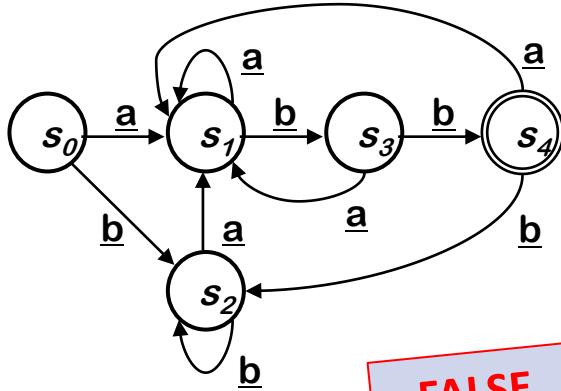
```

←

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p1$  and  $p2$ 
        then return  $\{p1, p2\}$ 
    }
    return  $p$ ;

```



```

 $P \leftarrow \{ F, S-F \}$ 
while ( P is still changing ) ←
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

```

```

Split( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

Minimization of DFA

$$P_2 = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

$$T = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$

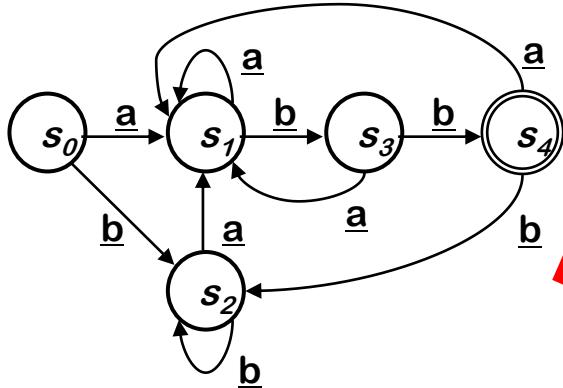
$$p = \underline{\{s_0, s_2\}}, \underline{\{s_1\}}, \underline{\{s_3\}}, \underline{\{s_4\}}$$

1st iter 2nd 3rd 4th

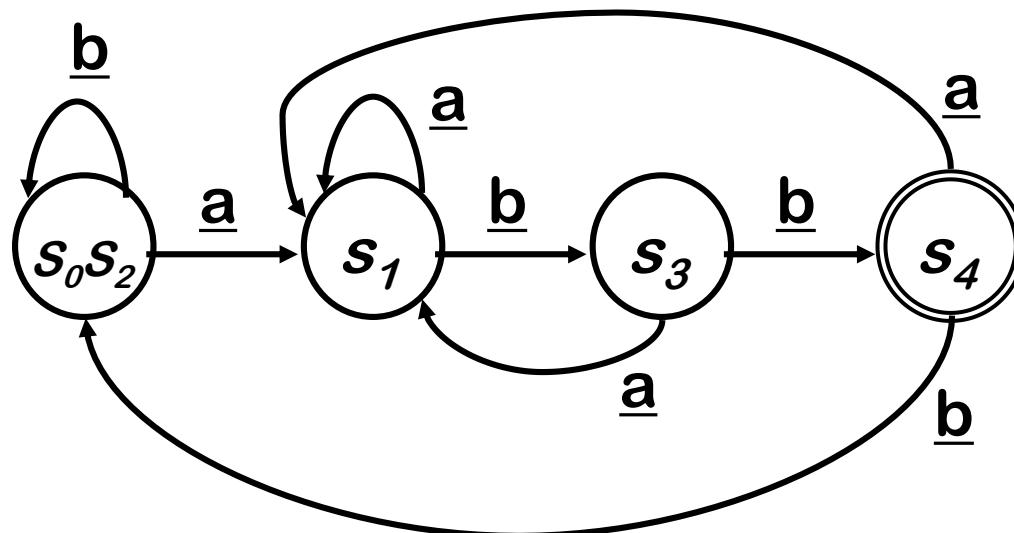
$$\sum = \{\underline{a}, \underline{b}\}$$

- While loop comes to a halt...
 - P does not change by previous step

Minimization of DFA



$$P = \{ \{s_0, s_2\}, \{s_1\}, \{s_3\}, \{s_4\} \}$$



DFA Minimization

The algorithm:

```

 $P \leftarrow \{ F, S-F \}$ 
while (  $P$  is still changing )
     $T \leftarrow \{ \}$ 
    for each set  $p \in P$ 
         $T \leftarrow T \cup \text{Split}(p)$ 
     $P \leftarrow T$ 

Split ( $p$ )
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return  $p$ ;

```

Why does this **terminate**?

- $p \in 2^S$ (powerset)
- Start off with 2 subsets of S : F and $S-F$
- **While loop** takes $P_i \rightarrow P_{i+1}$ by splitting 1 or more sets
- P_{i+1} is at least one step closer to the partition with $|S|$ sets
- Maximum of $|S|$ splits

Note that

- Partitions are never combined, only split.
- algorithm eventually terminates

This is a fixed-point algorithm!

DFA Minimization

The algorithm:

```
P ← { F, S-F }
while ( P is still changing )
    T ← {}
    for each set  $p \in P$ 
        T ← T ∪ Split (p)
    P ← T

Split (p)
    for each  $\alpha \in \Sigma$ 
    {
        if  $\alpha$  splits  $p$  into  $p_1$  and  $p_2$ 
        then return  $\{p_1, p_2\}$ 
    }
    return p;
```

Why does this **work**?

- *The algorithm maintains 2 invariants:*
 - 1) States remaining in the **same set** have not been **distinguished** yet by any string.
 - 2) States winding up in **different sets** are **distinguishable** by some string.
- The sets in the final partition contains the set of **distinguishable** states of the DFA.
- Proof sketch: Dragon book, 2nd ed., Section 3.9.6.

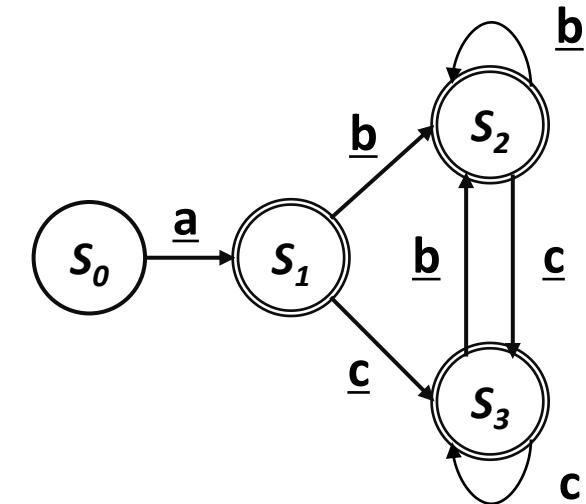
This is a fixed-point algorithm!

DFA Minimization Example

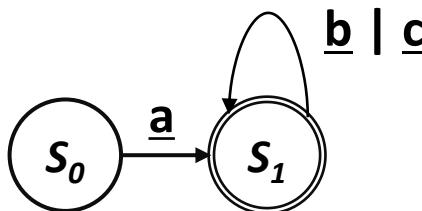
- Applying the minimization algorithm...

	<i>Current Partition</i>	<i>Split on</i>		
		<u>a</u>	<u>b</u>	<u>c</u>
P ₀	{S ₁ , S ₂ , S ₃ } {S ₀ }	<i>none</i>	<i>none</i>	<i>none</i>

accepting states



... to produce the minimal DFA

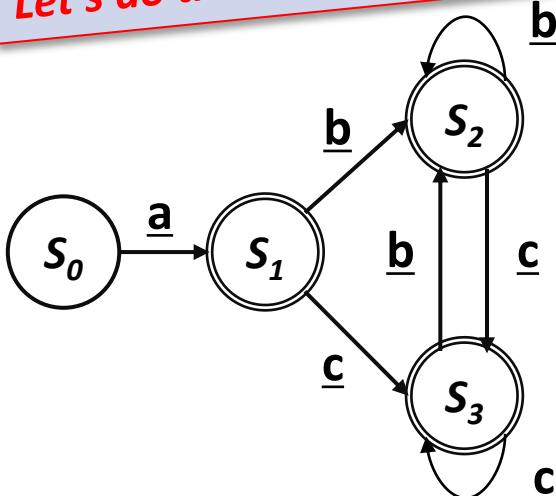


In lecture 5, we observed that a human would design a simpler automaton than Thompson's construction & the subset construction did.

Minimizing that DFA produces the one that a human would design.

Theoretical result: every RE language can be recognized by a minimal-state DFA that is unique up to state names

Let's do the same things to this DFA!



Minimization of DFA

$$P = \{ \} }$$

$$T = \{ \} }$$

$$p = \quad \quad \quad \sum = \{\underline{a}, \underline{b}, \underline{c}\}$$

$P \leftarrow \{ F, S-F \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $p \in P$

$T \leftarrow T \cup \text{Split}(p)$

$P \leftarrow T$

Split (p)

for each $\alpha \in \Sigma$

 {

 if α splits p into p_1 and p_2

then return $\{p_1, p_2\}$

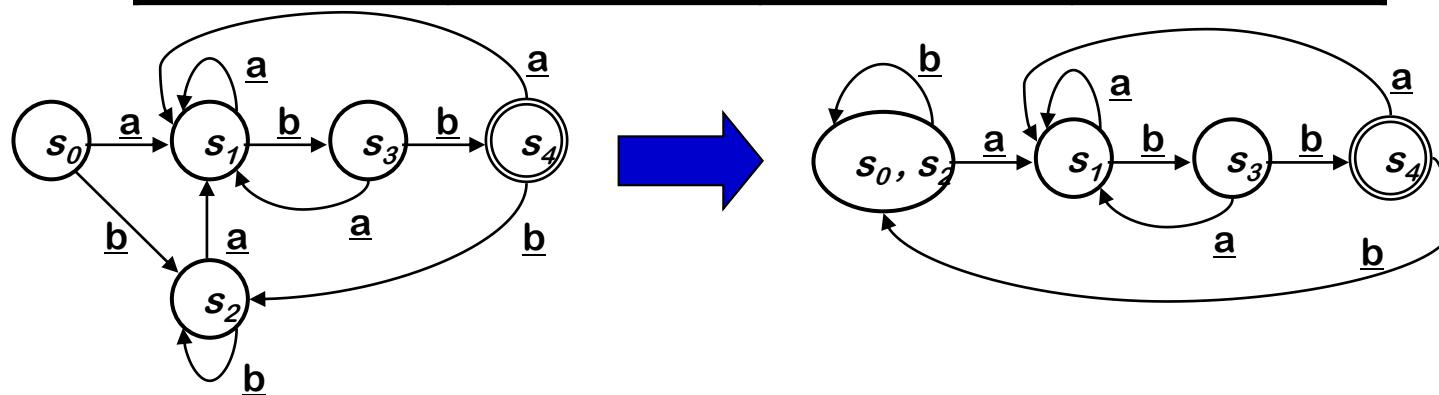
 }

return p ;

A Detailed Example

- Applying the minimization algorithm to the DFA for $(\underline{a}|\underline{b})^*abb$

Partition	Set	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_0, s_1, s_2, s_3\}$	none	$\{s_0, s_1, s_2\}, \{s_3\}$
	$\{s_4\}$	---	---
P_1	$\{s_0, s_1, s_2\}$	none	$\{s_0, s_2\}, \{s_1\}$
	$\{s_3\}$	---	---
	$\{s_4\}$	---	---
P_2	$\{s_0, s_2\}$,	none	none
	$\{s_1\}$	---	---
	$\{s_3\}$	---	---
	$\{s_4\}$	---	---



Outline

- The role of a scanner ✓
- Scanner concepts ✓
 - Tokens, Lexemes, Patterns
- Regular Expression & Automata (to be continued...)
 - Definitions of REs, DFAs and NFAs ✓
 - REs→NFA ✓
(Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)
 - NFA→DFA ✓
(subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)
 - DFA→minimal-state DFA ✓
(state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
- Scanner generators

Scanner Generators

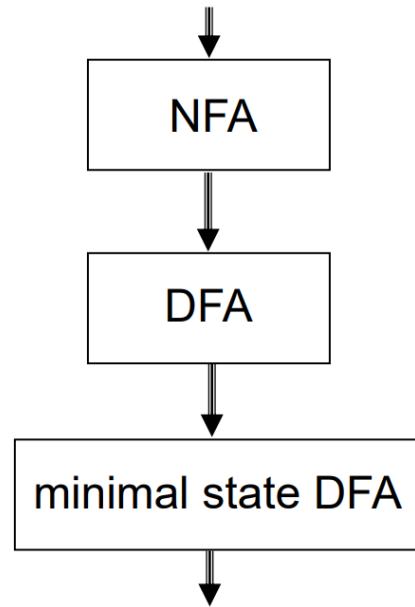
- Scanners generated in C:
 - lex (UNIX)
 - flex (GNU's fast lex, UNIX)
 - Mks lex (MS-DOS, Windows, OS/2)
- Scanners generated in Java:
 - Jlex (Princeton University)
 - JFlex
 - JavaCC (SUN Microsystems → now, Oracle)

The Scanner Spec in JLex

```
user code // copied verbatim to the scanner file
%%
JLEX directives (declarations)
%%
regular expression rules
```

How a Scanner Generator Works

Specification of tokens using REs



- Table-driven code (JLex), simulates a DFA on an input
- Hard-wired code
(like our hand-crafted scanner for Assignment #1)
 - Assignment #1 = Scanner + Parser

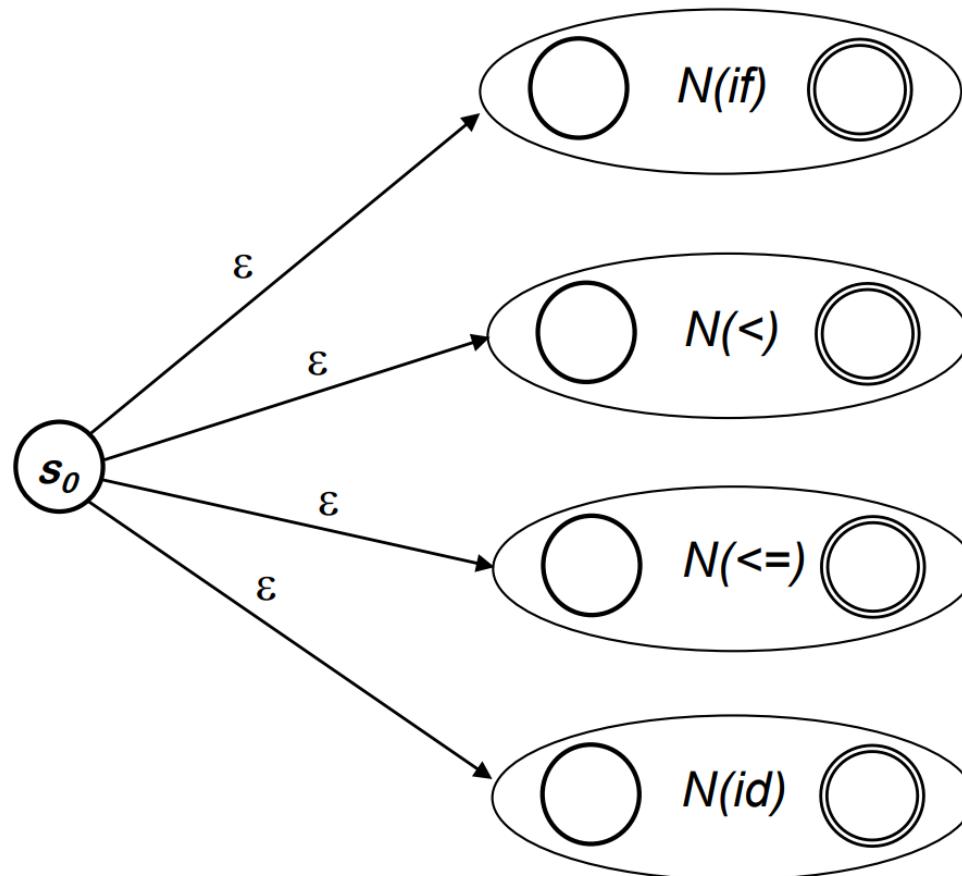
JLex Example Spec

```
%%
LETTER=[a-zA-Z_]
DIGIT=[0-9]
%%
"if" { return new Token(Token.IF, "if", src_pos); }
"<" { return new Token(Token.LESS, "<", src_pos); }
"<="
{ return new Token(Token.LESS_EQ, "<", src_pos); }
{LETTER}({LETTER}|{DIGIT})*
{ return new Token(Token.ID, "spelling", src_pos); }
```

Two rules:

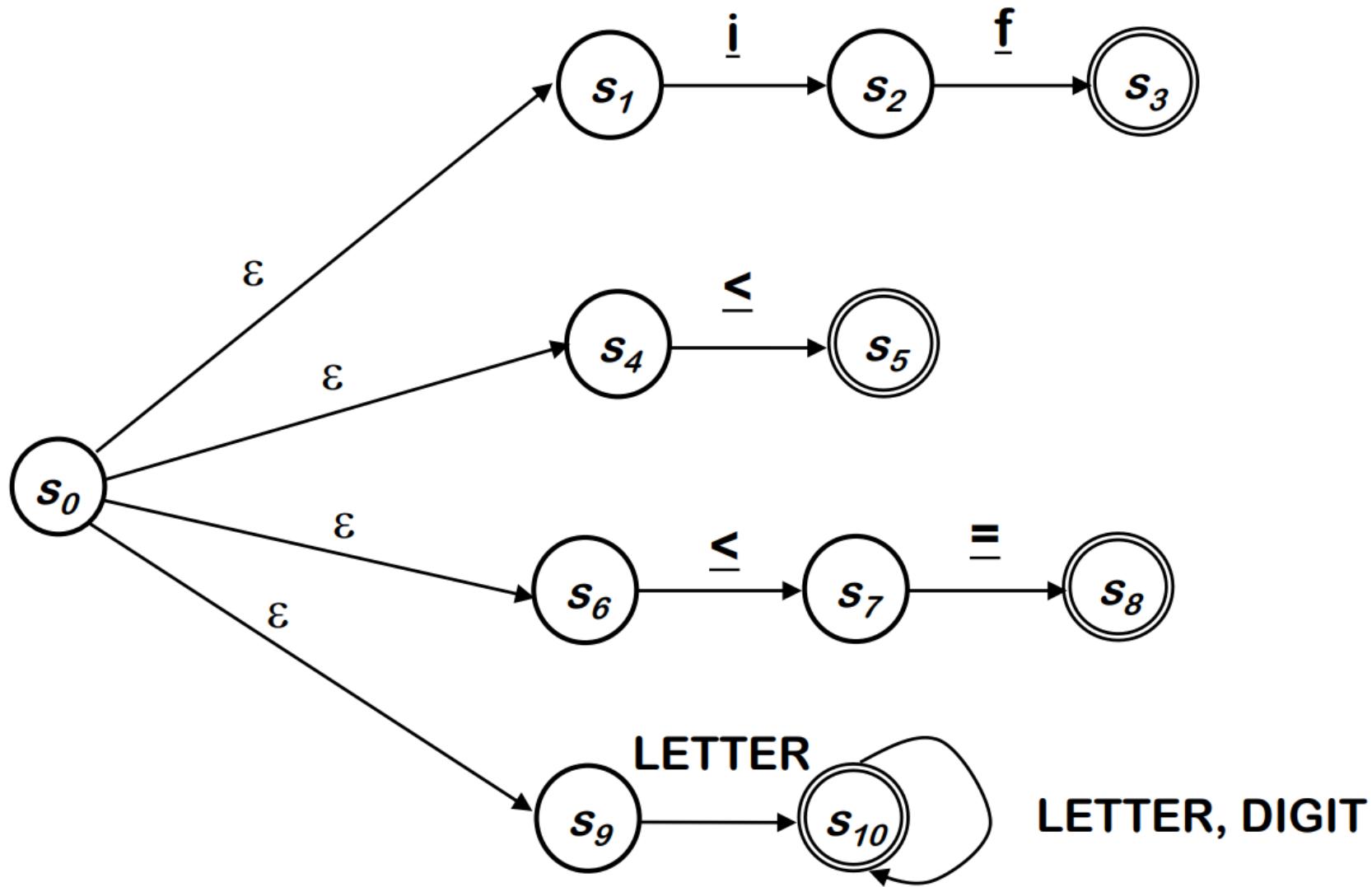
- The first pattern is used if more than one pattern matches.
 - Example: **if** as a keyword (IF) and not as an identifier (ID).
- The longest possible match is taken.
 - Example: “**<=**” as one token.

NFA for JLex Example Spec



- The NFAs for the different REs are combined as above.
- Instead of an NFA, a DFA can also be used for each pattern.

NFA for JLex Example Spec



Running JLex on a Sample Scanner Spec

Scanner.l: the spec passed to the scanner generator JLex

java JLex.Main Scanner.l

```
Processing first section -- user code.  
Processing second section -- JLex declarations.  
Processing third section -- lexical rules.  
Creating NFA machine representation.  
NFA comprised of 315 states.  
Working on character classes.....  
NFA has 46 distinct character classes.  
Creating DFA transition table.  
Working on DFA states.....  
Minimizing DFA transition table.  
109 states after removal of redundant states.  
Outputting lexical analyzer code.
```

Scanner.l.java: the generated scanner

javac Scanner.l.java

Limitations of Regular Languages

Advantages of Regular Expressions

- Simple & powerful notation for specifying patterns
- Automatic construction of fast recognizers (scanners)
- Many patterns can be specified with REs

Example — an expression grammar

Term: $[a-zA-Z] ([a-zA-z] | [0-9])^*$

Op: $\pm | - | * | /$

Expr : (Term Op)^{*} Term

Of course, this would generate a DFA ...

If REs are so useful ...

Why not use them for everything?

Limitations of Regular Languages (cont.)

- Regular expressions are limited in what can be expressed.
 - Example: it is not possible to specify a regular expression for $0^n 1^n$ (the set of strings of N zeroes followed by N ones).

For the same reason, we cannot specify the set of arithmetic expressions for which left and right parentheses match.
e.g., (a), ((a)), (((a))), (((((a)))))
- Regular expressions only usable to specify keywords, identifiers, literals, operators and punctuation characters of a programming language. For everything else (expressions, statements, nested statements, ...) we need a more powerful mechanism

Regular Expressions and Context-Free Grammars

- The “languages” that can be defined by REs and Context-free grammars (CFGs) have been extensively studied by theoretical computer scientists. These are some important conclusions / terminologies:
- RE are a “weaker” formalism than CFGs: Any language expressible by a RE can be expressed by a CFG **but not the other way around!**
- The languages expressible as RE are called regular languages
- Generally: a language that exhibits “**self embedding**” cannot be expressed by a RE.
- Programming languages exhibit self embedding.
Example: an expression can contain another expression

```
expr ::= id | integer | - expr | ( expr ) | expr op expr  
op   ::= + | - | * | /
```

How many expressions can be derived? Infinitely many.

Table-driven vs. Direct-coded Scanners

Jeonggeun Kim
Kyungpook National University



Table-driven vs. Directed-coded Scanners

δ	r	<u>0,1,2,3</u> <u>,4,5,6,</u> <u>7,8,9</u>	All others
s_0	s_1	s_e	s_e
s_1	s_e	s_2	s_e
s_2	s_e	s_2	s_e
s_e	s_e	s_e	s_e

Table encoding RE

Char \leftarrow next character
 State $\leftarrow s_0$
 while (Char \neq EOF)
 State $\leftarrow \delta(\text{State}, \text{Char})$
 Char \leftarrow next character
 if (State is a final state)
 then report success
 else report failure

Table-driven

Skeleton recognizer

```

switch (currentChar) {
  case '+': return token representation for '+'
  case '<':
    takeIt(); // adds '<' to the current lexeme
    if (currentChar == '=')
      return token representation for '<='
    else
      return token representation for '<'
}
  
```

Direct-coded

Ease of Implementation & Performance

Table-driven scanner encodes the automaton as a big table

- Costly lookup (computation of table cell based on state and input character); not cache-friendly if tablesizexceeds size of cache.
- Easy, automated way to implement. Hard to implement manually.

Direct-coded scanner encodes the automaton in the **program counter**:

```
"if '<' then read next char and check for "<="
else ..."
```

- Automated and manual implementations feasible
- No table and thus no table lookup needed
- More complex control flow
 - Ok for architectures with good branch prediction (x86, ...).
- Tends to be faster than table-driven scanner on general-purpose CPU

Training for Lexical Analysis

Jeonggeun Kim
Kyungpook National University



Writing Regular Expressions

Write a regular expression for each of the following sets of tokens:

- 1) Ruby binary literals consisting of “0b” followed by the binary number. Examples: 0b001011, 0b01.
- 2) Ruby binary literals, with an optional underscore (“_”) between a pair of binary digits.
Examples: 0b0_101 , 0b11_01 , but not 0b_1 and not 0b1_ .
- 3) Ada identifiers: a letter followed by any number of letters, digits, and underscores.
An identifier must not end in an underline or have two underscores in a row.
- 4) A floating-point number: one or more digits (whole-number part) followed by a decimal point (“.”) and one or more digits (fractional part).
Examples: 2.24, 0.1234.
- 5) Floating-point number in scientific notation: same as above, but optionally followed by “e” or “E”, and a signed integer exponent.
Examples: 1.2e-2, 2.3E+34, 2.3E34.

1) Ruby binary literals consisting of “0b” followed by the binary number.

Examples: 0b001011, 0b01.

- 2) Ruby binary literals, with an optional underscore (“_”) between a pair of binary digits.

Examples: 0b0_101 , 0b11_01 , but not 0b_1 and not 0b1_ .

- 3) Ada identifiers: a letter followed by any number of letters, digits, and underlines. An identifier must not end in an underline or have two underlines in a row.

- 4) A floating-point number: one or more digits (whole-number part) followed by a decimal point (“.”) and one or more digits (fractional part).

Examples: 2.24, 0.1234.

- 5) Floating-point number in scientific notation: same as above, but optionally followed by “e” or “E”, and a signed integer exponent.

Examples: 1.2e-2, 2.3E+34, 2.3E34.

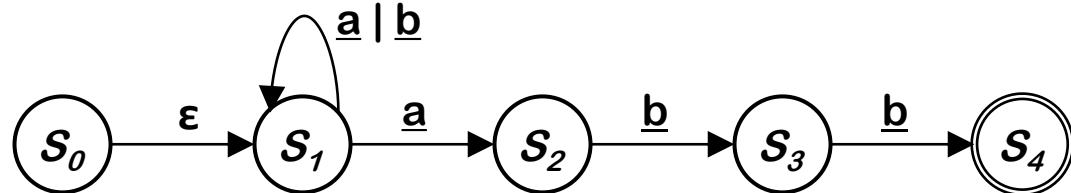
Subset Construction Algorithm (Animation)

Jeonggeun Kim
Kyungpook National University



Example: Simulating an NFA ``on the fly''

Input =



- **On-the-fly simulation** of an NFA is useful in scenarios where a regular expression is used only once (example: in a text editor or IDE).
- In a compiler, the regular expressions for the programming language's tokens are used over and over again → need a more efficient approach → next slide.

Algorithm: NFA → DFA with Subset Construction

Subset construction works on sets of NFA states.

Each set of NFA states becomes a DFA state.

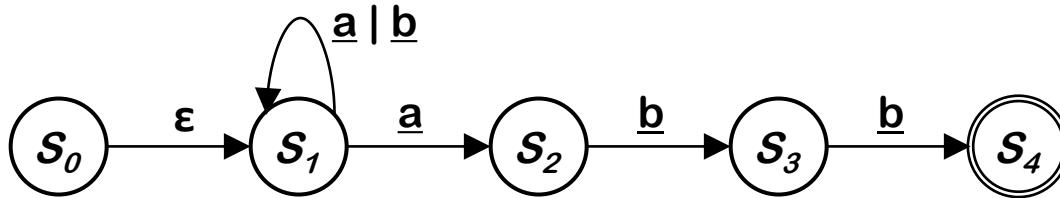
Two key functions

- *Move (S, a)* : set of states reachable from set S by a
- *ϵ -closure (S)* : set of states reachable from set S by ϵ

The algorithm:

- Start state derived from s_0 of the NFA:
 - Take its ϵ -closure $S_0 = \epsilon\text{-closure}(\{s_0\})$
- Compute $\text{Move}(S_0, \alpha)$ for each $\alpha \in \Sigma$, and take its ϵ -closure
- Iterate until no more states are added

Sounds more complex than it is...

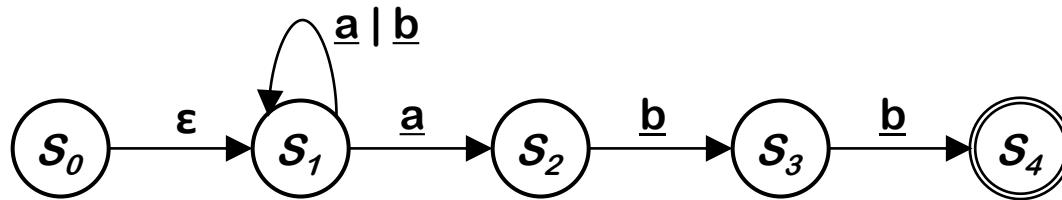


```

Dstates ← {};
add  $\epsilon$ -closure( $S_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $\alpha \in \Sigma$  {
     $U \leftarrow \epsilon$ -closure( $\text{Move}(T, \alpha)$ )
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, \alpha] \leftarrow U$ 
  }
}
  
```

Note:

- here δ refers to the transition function of the **DFA**.
- δ takes a set of NFA-states (= a **DFA state**) and a symbol α as input, and returns a set of NFA-states (= a **DFA state**).



$Dstates \leftarrow \{\}$;

add ϵ -closure(s_0) as an unmarked state to $Dstates$;
while (there is an unmarked state T in $Dstates$) {

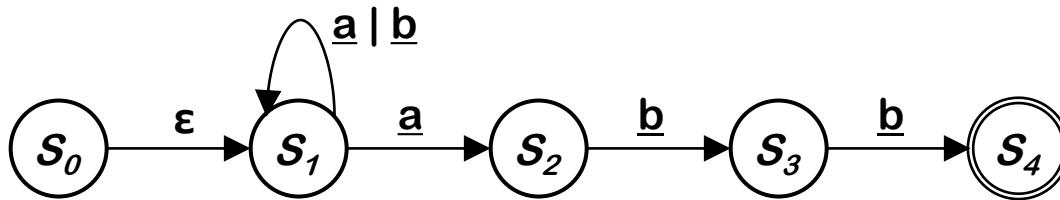
mark T ;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to $Dstates$;

$\delta [T, a] \leftarrow U$

}

$Dstates = \{\}$

ϵ -closure(s_0) = { s_0, s_1 }



$Dstates \leftarrow \{\}$;

add ϵ -closure(S_0) as an
unmarked state to $Dstates$;
while (there is an unmarked
state T in $Dstates$) {

mark T ;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
 if ($U \notin Dstates$) **then**
 add U as an unmarked
 state to $Dstates$;
 $\delta [T, a] \leftarrow U$

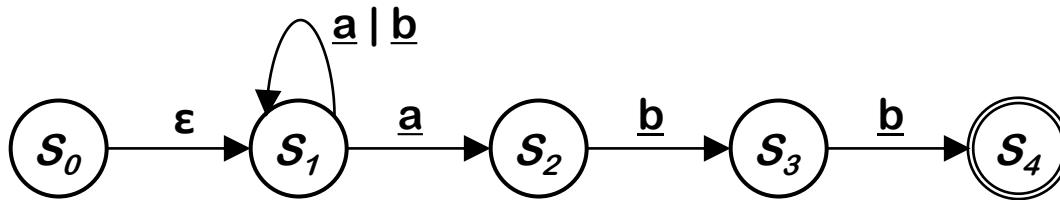
}

$Dstates = \{ \{ s_0, s_1 \} \}$

$Dstates$ contains now one element, i.e.,
the state $\{ s_0, s_1 \}$.

Note1: unmarked states are printed
in red.

Note2: $Dstates$ means ``deterministic
states, the states that will make up the
DFA''.



Dstates $\leftarrow \{\}$;

add ϵ -closure(S_0) as an
unmarked state to *Dstates*;
while (there is an unmarked
state T in *Dstates*) {

mark T ;

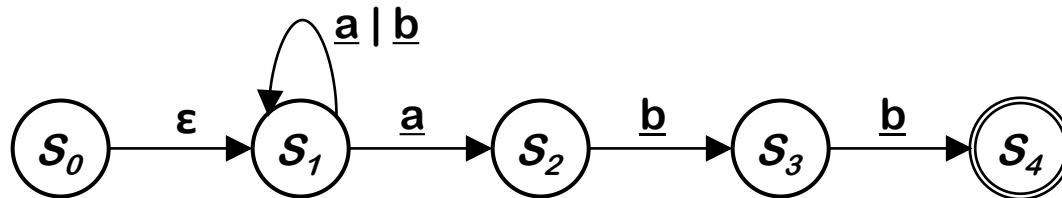
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($\text{Move}(T, a)$)
if ($U \notin \text{Dstates}$) **then**
add U as an unmarked
state to *Dstates*;
 $\delta [T, a] \leftarrow U$

}

Dstates = { { s_0, s_1 } }

Now we pick state { s_0, s_1 } from *Dstates*
and mark it.

Note: marked states are printed in green.



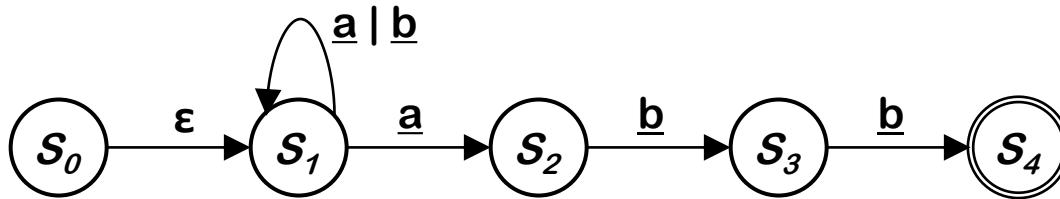
$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
 unmarked state to $Dstates$;
while (there is an unmarked
 state T in $Dstates$) {
mark T ;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
 state to $Dstates$;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \} \}$$

$$T = \{ s_0, s_1 \}$$

$$\Sigma = \{ \underline{a}, \underline{b} \}$$

The for loop is going to iterate over all symbols in Σ .
 First iteration: a
 Second iteration: b.



$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
 unmarked state to $Dstates$;
while (there is an unmarked
 state T in $Dstates$) {
mark T ;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
 state to $Dstates$;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \} \}$$

$$T = \{ s_0, s_1 \}$$

$$\Sigma = \{ \underline{a}, \underline{b} \}$$

The for loop is going to iterate over all symbols in Σ .

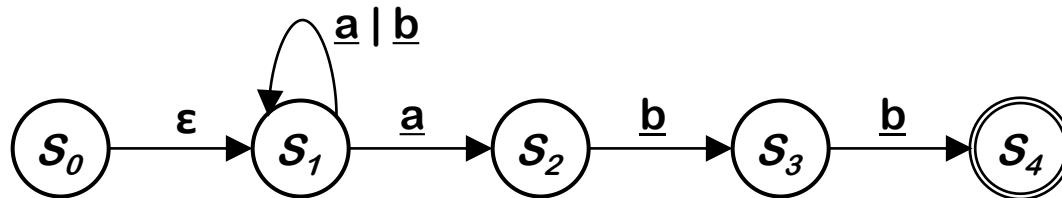
First iteration: $a = \underline{a}$

Second iteration: $a = \underline{b}$.

$$Move (T, \underline{a}) = \{ s_1, s_2 \}$$

$$\epsilon\text{-closure}(\{ s_1, s_2 \}) = \{ s_1, s_2 \}$$

$$U = \{ s_1, s_2 \}$$



$Dstates \leftarrow \{\}$;

add ϵ -closure(S_0) as an unmarked state to $Dstates$;
while (there is an unmarked state T in $Dstates$) {

mark T ;

for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to $Dstates$;
 $\delta [T, a] \leftarrow U$

}

$Dstates = \{ \{ s_0, s_1 \} \}$

$T = \{ s_0, s_1 \}$

$\Sigma = \{ \underline{a}, \underline{b} \}$

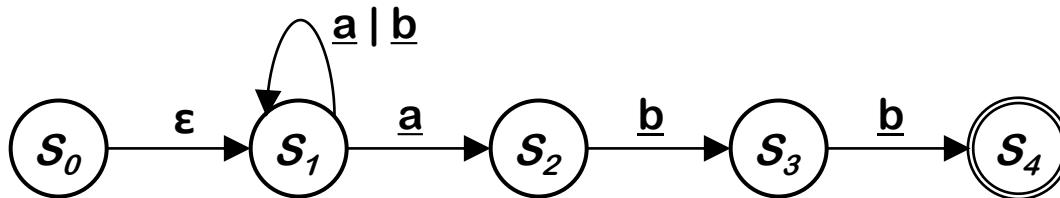
The for loop is going to iterate over all symbols in Σ .

First iteration: $a = \underline{a}$

Second iteration: $a = \underline{b}$.

$U = \{ s_1, s_2 \}$

State U is not contained in $Dstates$ yet, so we add it as an unmarked state.



$Dstates \leftarrow \{\}$;

add ϵ -closure(s_0) as an unmarked state to $Dstates$;
while (there is an unmarked state T in $Dstates$) {

mark T ;

for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to $Dstates$;

$\delta [T, a] \leftarrow U$

}



$Dstates = \{ \{ s_0, s_1 \} \}$

$T = \{ s_0, s_1 \}$

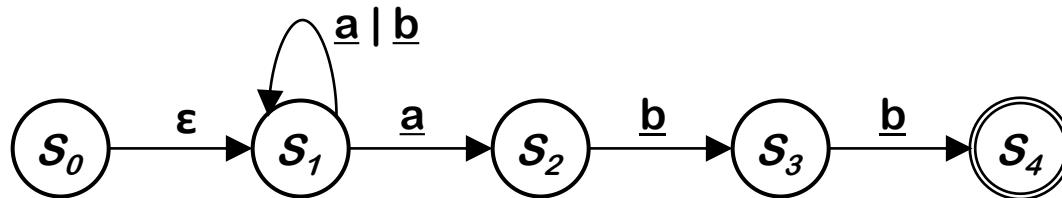
$\Sigma = \{ \underline{a}, \underline{b} \}$

$U = \{ s_1, s_2 \}$

Here we start assembling the transition function δ for our DFA.

If we read an \underline{a} in state T , then the DFA will move to DFA state $\{s_1, s_2\}$:

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	



```

 $Dstates \leftarrow \{\}$ ;
add  $\epsilon$ -closure( $S_0$ ) as an
unmarked state to  $Dstates$ ;
while ( there is an unmarked
state  $T$  in  $Dstates$  ) {
    mark  $T$ ;
    for each  $a \in \Sigma$  {
         $U \leftarrow \epsilon$ -closure( $Move(T, a)$ )
        if (  $U \notin Dstates$  ) then
            add  $U$  as an unmarked
            state to  $Dstates$ ;
         $\delta [T, a] \leftarrow U$ 
    }
}

```

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \} \}$$

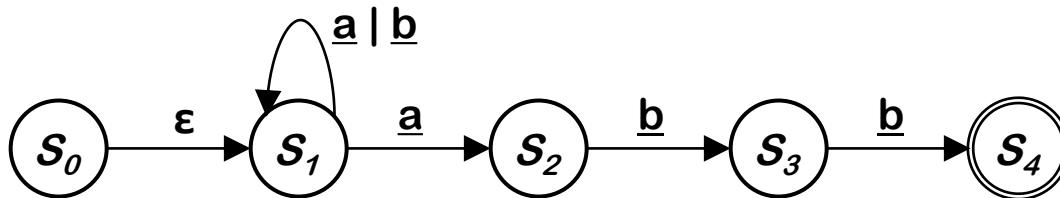
$$T = \{ s_0, s_1 \}$$

$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$$U = \{ s_1, s_2 \}$$

Second iteration of the for loop: $a = \underline{b}$.

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	



$Dstates \leftarrow \{\}$;

add ϵ -closure(S_0) as an
unmarked state to $Dstates$;
while (there is an unmarked
state T in $Dstates$) {

mark T ;

for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to $Dstates$;
 $\delta [T, a] \leftarrow U$

}

$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \} \}$

$T = \{ s_0, s_1 \}$

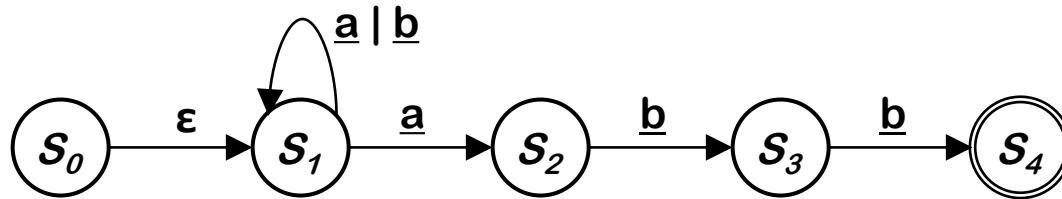
$\Sigma = \{ \underline{a}, \underline{b} \}$

Second iteration of the for loop: $a = \underline{b}$.

$Move(T, \underline{b}) = \{ s_1 \}$

ϵ -closure({ s_1 }) = { s_1 }

$U = \{ s_1 \}$



$Dstates \leftarrow \{\}$;

add ϵ -closure(s_0) as an
unmarked state to $Dstates$;
while (there is an unmarked
state T in $Dstates$) {

mark T ;

for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to $Dstates$;
 $\delta [T, a] \leftarrow U$

}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \} \}$$

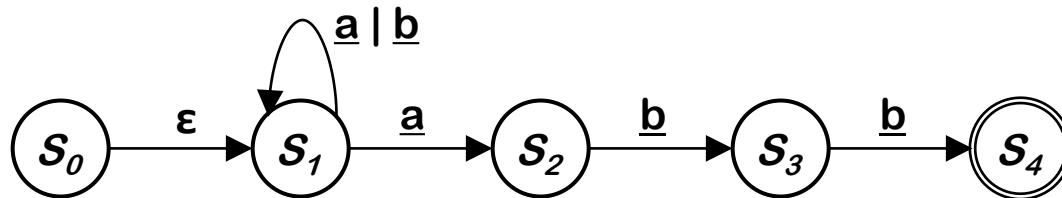
$$T = \{ s_0, s_1 \}$$

$$\Sigma = \{ \underline{a}, \underline{b} \}$$

Second iteration of the **for** loop: $a = \underline{b}$.

$$U = \{ s_1 \}$$

State U is not contained in $Dstates$
yet, so we add it as an unmarked state.



$Dstates \leftarrow \{\}$;

add ϵ -closure(S_0) as an unmarked state to $Dstates$;
while (there is an unmarked state T in $Dstates$) {

mark T ;

for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to $Dstates$;

$\delta [T, a] \leftarrow U$

}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \} \}$$

$$T = \{ s_0, s_1 \}$$

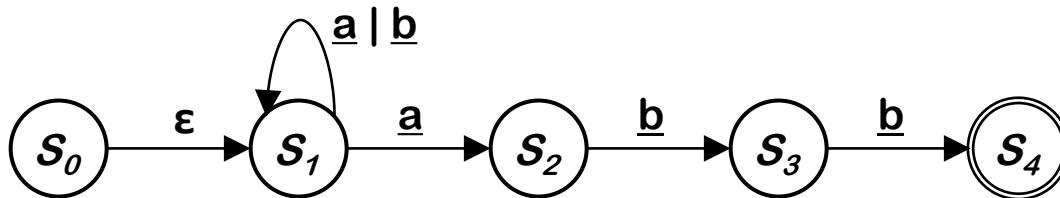
$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$$U = \{ s_1 \}$$

Second iteration of the **for** loop: $a = \underline{b}$.

If we read a \underline{b} in state T , then the DFA will move to state $\{s_1\}$:

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{ s_1 \}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

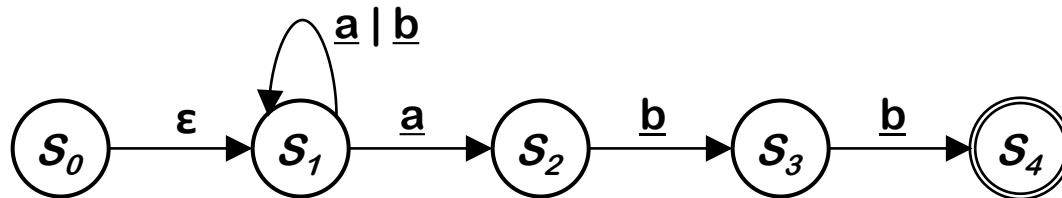
$Dstates \leftarrow \{ \};$
add ϵ -closure(S_0) as an
unmarked state to Dstates;
while (*there is an unmarked state T in Dstates*) {
mark T;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to Dstates;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

Dstates contains now two unmarked states, $\{ s_1, s_2 \}$ and $\{ s_1 \}$.

We pick state $\{ s_1, s_2 \}$ next.

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{ s_1 \}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

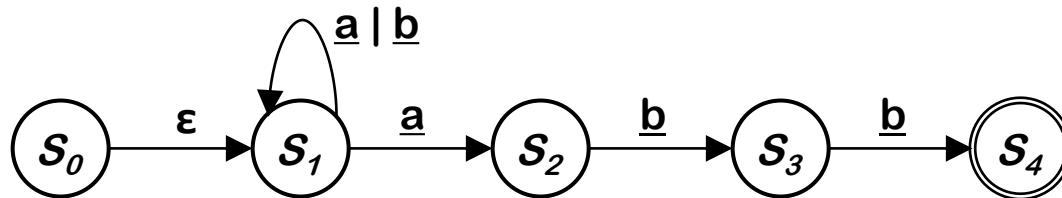
Dstates ← {};
add ε-closure( $s_0$ ) as an
    unmarked state to Dstates;
while ( there is an unmarked
    state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $\alpha \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
                state to Dstates;
         $\delta [T, \alpha] \leftarrow U$ 
    }
}

```

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

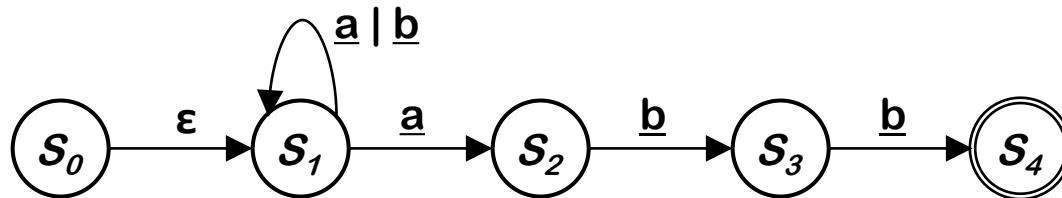
$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{a}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{ s_1 \}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add ε-closure( $s_0$ ) as an
    unmarked state to Dstates;
while ( there is an unmarked
    state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $\alpha \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
                state to Dstates;
         $\delta [T, \alpha] \leftarrow U$ 
    }
}

```

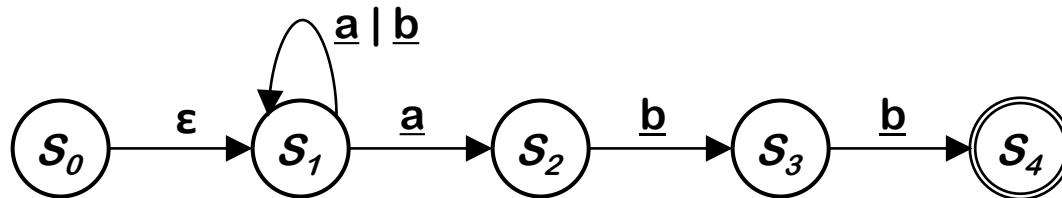
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{a}$$

$$U = \{ s_1, s_2 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{ s_1 \}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

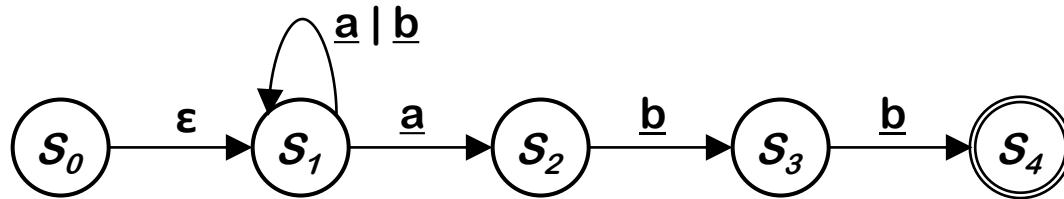
$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{a}$$

$$U = \{ s_1, s_2 \}$$

State U is already contained in Dstates
(we are currently working on it!). Nothing
to be done here...

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{ s_1 \}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $a \in \Sigma$  {
     $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, a] \leftarrow U$ 
  }
}
  
```

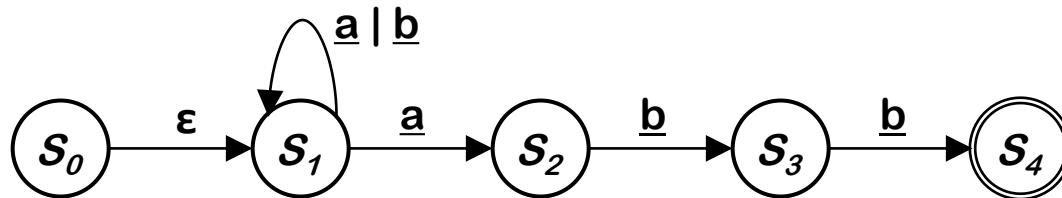
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{a}$$

$$U = \{ s_1, s_2 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

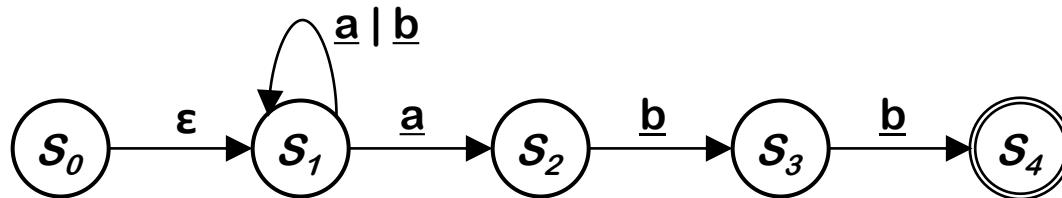
$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
 unmarked state to $Dstates$;
while (there is an unmarked
 state T in $Dstates$) {
mark T ;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
 state to $Dstates$;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{b}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add ε-closure( $s_0$ ) as an
    unmarked state to Dstates;
while ( there is an unmarked
    state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $\alpha \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
                state to Dstates;
         $\delta [T, \alpha] \leftarrow U$ 
    }
}

```

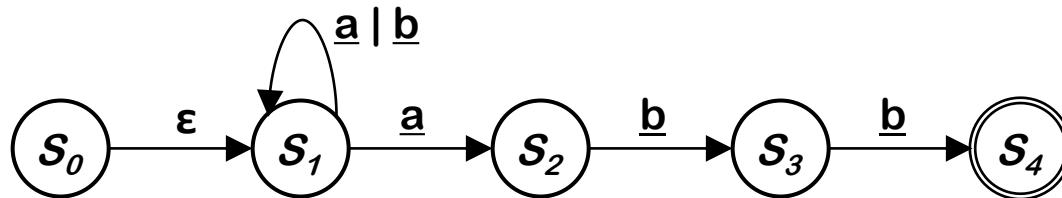
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{b}$$

$$U = \{ s_1, s_3 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

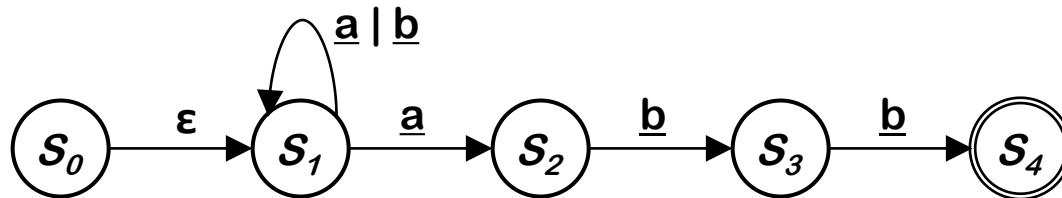
$$Dstates = \{ \{ S_0, S_1 \}, \{ S_1, S_2 \}, \{ S_1 \} \}$$

$$T = \{ S_1, S_2 \}$$

$$\alpha = \underline{b}$$

$$U = \{ S_1, S_3 \}$$

δ	\underline{a}	\underline{b}
$\{ S_0, S_1 \}$	$\{ S_1, S_2 \}$	$\{ S_1 \}$
$\{ S_1, S_2 \}$	$\{ S_1, S_2 \}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $a \in \Sigma$  {
     $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, a] \leftarrow U$ 
  }
}
  
```

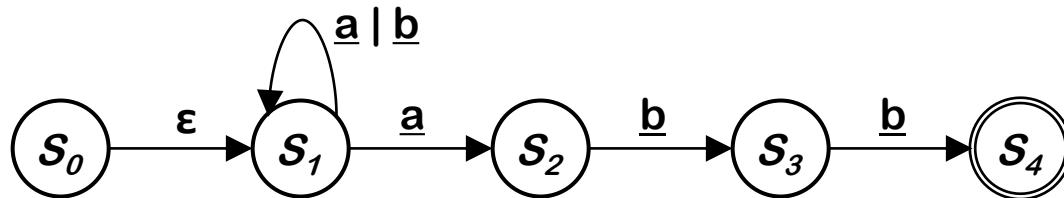
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$a = \underline{b}$$

$$U = \{ s_1, s_3 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
→ while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $a \in \Sigma$  {
     $U \leftarrow \epsilon$ -closure( $\text{Move}(T, a)$ )
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, a] \leftarrow U$ 
  }
}
  
```

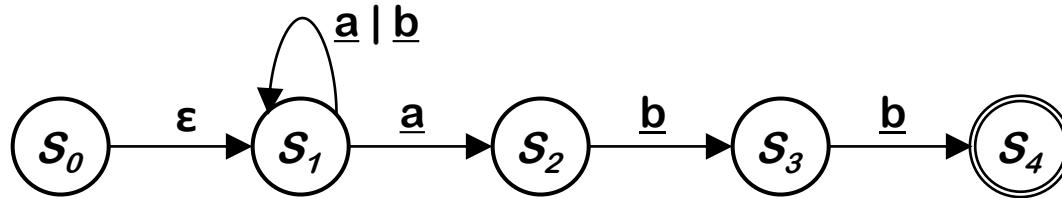
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1, s_2 \}$$

$$\alpha = \underline{b}$$

{ s_1 } is the next unmarked state in Dstates.

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

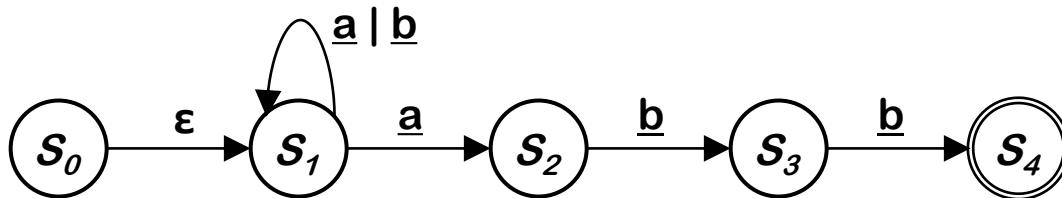
$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{a}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

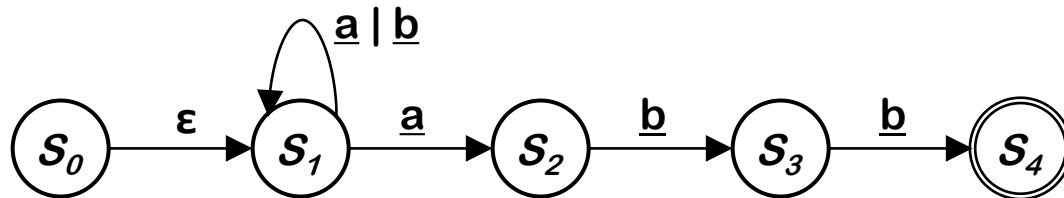
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{a}$$

$$U = \{ s_1, s_2 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $a \in \Sigma$  {
     $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, a] \leftarrow U$ 
  }
}
  
```

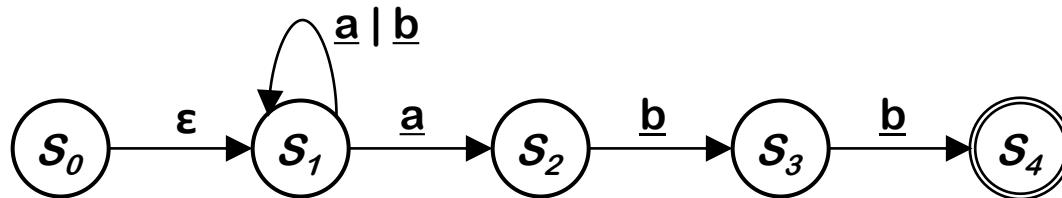
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{a}$$

$$U = \{ s_1, s_2 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

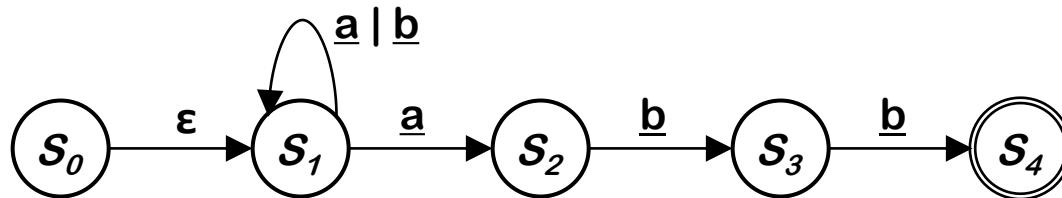
$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{b}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $a \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
            state to Dstates;
         $\delta [T, a] \leftarrow U$ 
    }
}

```

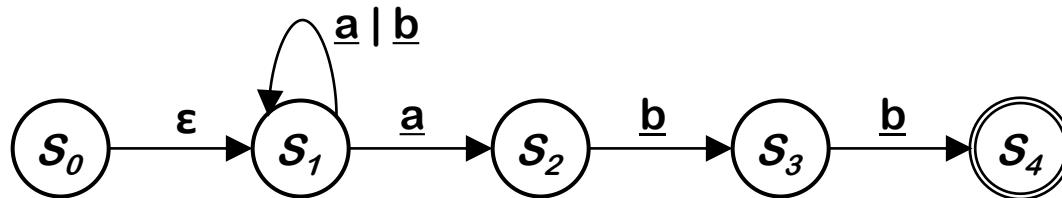
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{b}$$

$$U = \{ s_1 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add ε-closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $\alpha \in \Sigma$  {
     $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, \alpha] \leftarrow U$ 
  }
}
  
```

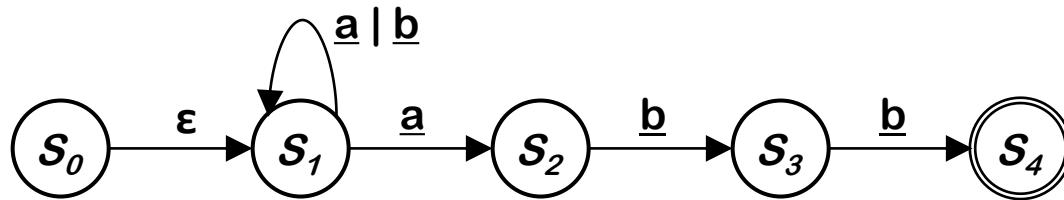
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1 \}$$

$$\alpha = \underline{b}$$

$$U = \{ s_1 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

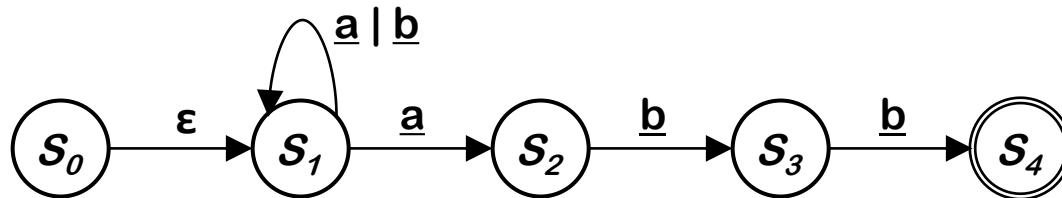
Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $a \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
            state to Dstates;
         $\delta [T, a] \leftarrow U$ 
    }
}

```

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$\alpha =$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {

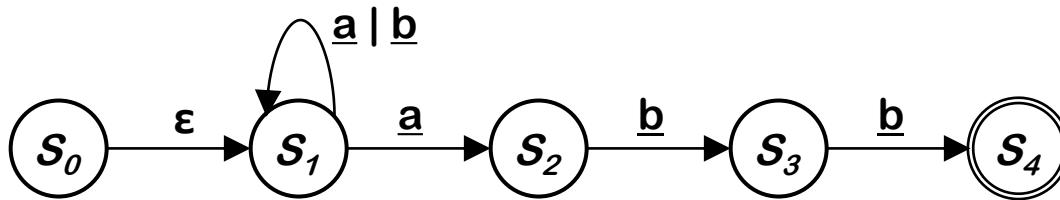
→
mark T;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, a] \leftarrow U$
}
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s1, s3 \}$$

$$\alpha =$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(S_0) as an
unmarked state to Dstates;
while (*there is an unmarked state T in Dstates*) {
mark T;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked state to Dstates;
 $\delta [T, a] \leftarrow U$
}

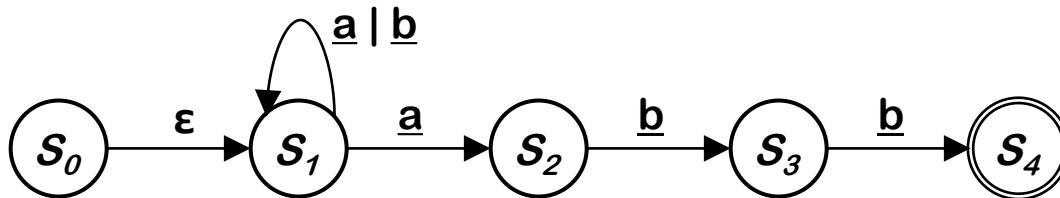
$$Dstates = \{ \{ S_0, S_1 \}, \{ S_1, S_2 \}, \{ S_1 \}, \{ S_1, S_3 \} \}$$

$$T = \{ S_1, S_3 \}$$

$$a = \underline{a}$$

If we are in state T, reading an \underline{a} , we get to $\{ S_1, S_2 \}$. This transition is added to δ .

δ	\underline{a}	\underline{b}
$\{ S_0, S_1 \}$	$\{ S_1, S_2 \}$	$\{ S_1 \}$
$\{ S_1, S_2 \}$	$\{ S_1, S_2 \}$	$\{ S_1, S_3 \}$
$\{ S_1 \}$	$\{ S_1, S_2 \}$	$\{ S_1 \}$
$\{ S_1, S_3 \}$	$\{ S_1, S_2 \}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

```

Dstates ← {};
add ε-closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $\alpha \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
        if (  $U \notin \text{Dstates}$  ) then
            add  $U$  as an unmarked
            state to Dstates;
         $\delta [T, \alpha] \leftarrow U$ 
    }
}

```

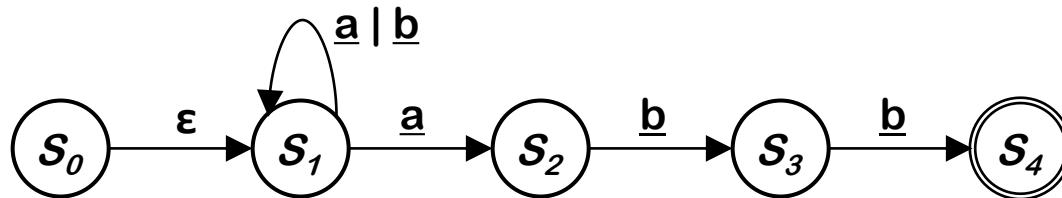
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1, s_3 \}$$

$$\alpha = \underline{b}$$

If we are in state T , reading a \underline{b} , we get to $\{ s_1, s_4 \} \dots$

δ	<u>a</u>	<u>b</u>
$\{ s_0, s_1 \}$	$\{ s_1, s_2 \}$	$\{ s_1 \}$
$\{ s_1, s_2 \}$	$\{ s_1, s_2 \}$	$\{ s_1, s_3 \}$
$\{ s_1 \}$	$\{ s_1, s_2 \}$	$\{ s_1 \}$
$\{ s_1, s_3 \}$	$\{ s_1, s_2 \}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $\alpha \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, \alpha)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, \alpha] \leftarrow U$
}

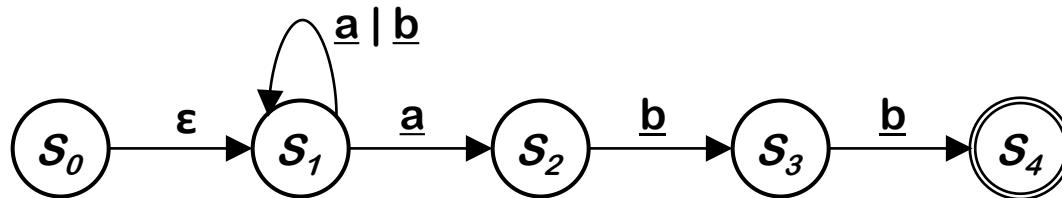
$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

$$T = \{ s_1, s_3 \}$$

$$\alpha = \underline{b}$$

$$U = \{ s_1, s_4 \}$$

δ	a	b
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$Dstates \leftarrow \{ \}$;
add ϵ -closure(s_0) as an
unmarked state to Dstates;
while (*there is an unmarked*
state T in Dstates) {
mark T;
for each $a \in \Sigma$ {
 $U \leftarrow \epsilon$ -closure($Move(T, a)$)
if ($U \notin Dstates$) **then**
add U as an unmarked
state to Dstates;
 $\delta [T, a] \leftarrow U$
}

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \} \}$$

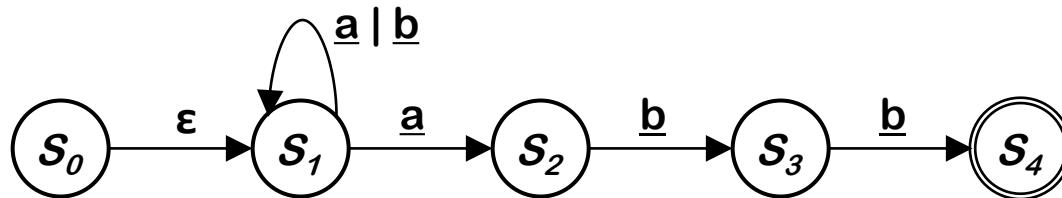
$$T = \{ s_1, s_3 \}$$

$$a = \underline{b}$$

$$U = \{ s_1, s_4 \}$$

U not yet in Dstates → add.

δ	<u>a</u>	<u>b</u>
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \}, \{ s_1, s_4 \} \}$$

```

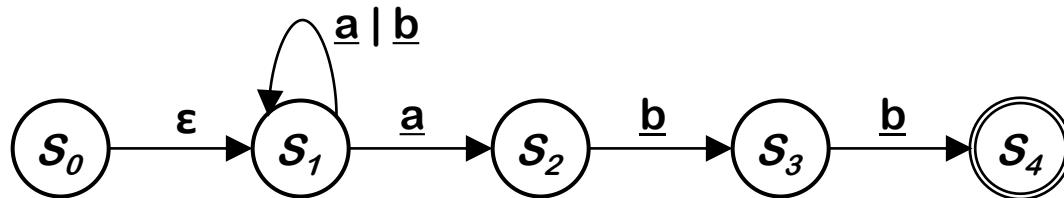
Dstates ← {};
add ε-closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $\alpha \in \Sigma$  {
     $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, \alpha))$ 
    if (  $U \notin Dstates$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, \alpha] \leftarrow U$ 
  }
}
  
```

$$T = \{ s_1, s_3 \}$$

$$\alpha = \underline{b}$$

$$U = \{ s_1, s_4 \}$$

δ	\underline{a}	\underline{b}
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	$\{s_1, s_4\}$



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

Dstates = { { s_0, s_1 }, { s_1, s_2 }, { s_1 }, { s_1, s_3 }, { s_1, s_4 } }

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
  mark  $T$ ;
  for each  $a \in \Sigma$  {
     $U \leftarrow \epsilon$ -closure( $\text{Move}(T, a)$ )
    if (  $U \notin \text{Dstates}$  ) then
      add  $U$  as an unmarked
      state to Dstates;
     $\delta [T, a] \leftarrow U$ 
  }
}
  
```

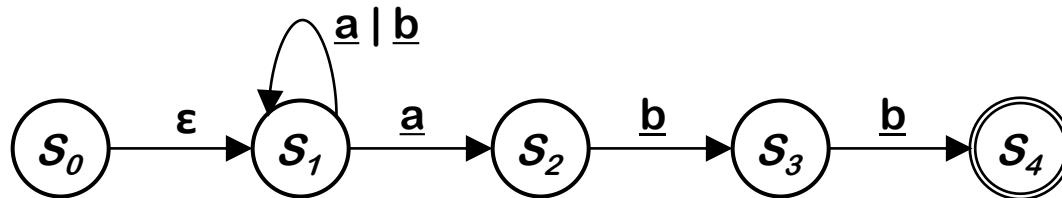
$T = \{ s_1, s_4 \}$ is the last unmarked state.

Reading an \underline{a} takes us to { s_1, s_2 }.

Reading a \underline{b} takes us to { s_1 }.

Those two transitions are added to δ .

δ	<u>a</u>	<u>b</u>
{ s_0, s_1 }	{ s_1, s_2 }	{ s_1 }
{ s_1, s_2 }	{ s_1, s_2 }	{ s_1, s_3 }
{ s_1 }	{ s_1, s_2 }	{ s_1 }
{ s_1, s_3 }	{ s_1, s_2 }	{ s_1, s_4 }
{ s_1, s_4 }	{ s_1, s_2 }	{ s_1 }



$$\Sigma = \{ \underline{a}, \underline{b} \}$$

$$Dstates = \{ \{ s_0, s_1 \}, \{ s_1, s_2 \}, \{ s_1 \}, \{ s_1, s_3 \}, \{ s_1, s_4 \} \}$$

```

Dstates ← {};
add  $\epsilon$ -closure( $s_0$ ) as an
unmarked state to Dstates;
while ( there is an unmarked
state  $T$  in Dstates ) {
    mark  $T$ ;
    for each  $a \in \Sigma$  {
         $U \leftarrow \epsilon\text{-closure}(\text{Move}(T, a))$ 
        if (  $U \notin Dstates$  ) then
            add  $U$  as an unmarked
            state to Dstates;
         $\delta [T, a] \leftarrow U$ 
    }
}

```

No more unmarked states in Dstates.

The algorithm stops.

δ	<u>a</u>	<u>b</u>
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	$\{s_1, s_4\}$
$\{s_1, s_4\}$	$\{s_1, s_2\}$	$\{s_1\}$

δ	a	b
$\{s_0, s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_2\}$	$\{s_1, s_2\}$	$\{s_1, s_3\}$
$\{s_1\}$	$\{s_1, s_2\}$	$\{s_1\}$
$\{s_1, s_3\}$	$\{s_1, s_2\}$	$\{s_1, s_4\}$
$\{s_1, s_4\}$	$\{s_1, s_2\}$	$\{s_1\}$

All DFA states that contain an accepting NFA state become accepting states in the DFA!

The DFA for the above transition function δ :

