

2. Lexical Analysis

Oscar Nierstrasz

Thanks to Jens Palsberg and Tony Hosking for their kind permission to reuse and adapt the CS132 and CS502 lecture notes.

<http://www.cs.ucla.edu/~palsberg/>

<http://www.cs.purdue.edu/homes/hosking/>

Roadmap



- > Introduction
- > Regular languages
- > Finite automata recognizers
- > From RE to DFAs and back again
- > Limits of regular languages

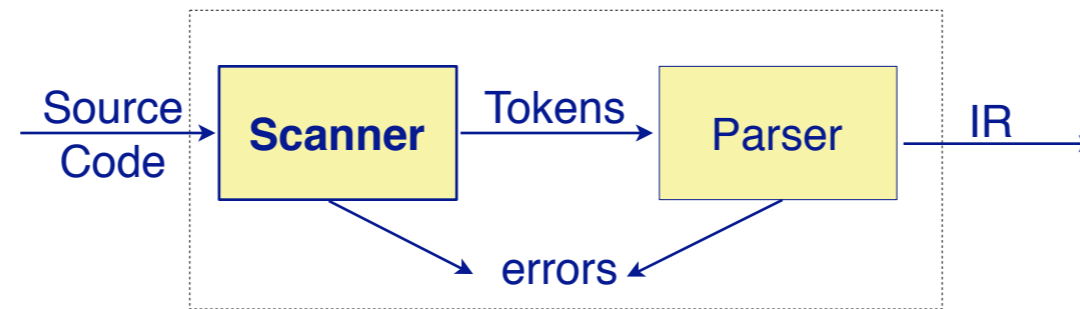
See, *Modern compiler implementation in Java* (Second edition), chapter 2.

Roadmap



- > **Introduction**
- > Regular languages
- > Finite automata recognizers
- > From RE to DFAs and back again
- > Limits of regular languages

Lexical Analysis



1. Maps sequences of characters to *tokens*
2. Eliminates white space (tabs, blanks, comments *etc.*)

`x = x + y` → `<ID,x> <EQ> <ID,x> <PLUS> <ID,y>`

The string value of a token is a *lexeme*.

How to specify rules for token classification?

A scanner must recognize various parts of the language's syntax

Some parts are easy:

White space

```
<WS> ::= <WS> ' '  
      | <WS> '\t'  
      | ' '  
      | '\t'
```

Keywords and operators

specified as literal patterns: do, end

Comments

opening and closing delimiters: /* ... */

Specifying patterns

Other parts are harder:

Identifiers

alphabetic followed by k alphanumerics ($_$, $\$$, $\&$, ...)

Numbers

integers: 0 or digit from 1–9 followed by digits from 0–9

decimals: integer '.' digits from 0–9

reals: (integer or decimal) 'E' (+ or –) digits from 0–9

complex: '(' real ', ' real ')'

We need an expressive notation to specify these patterns!

A key issue is ...



6

why don't we write it by hand?

Roadmap



- > Introduction
- > **Regular languages**
- > Finite automata recognizers
- > From RE to DFAs and back again
- > Limits of regular languages

Languages and Operations

A *language* is a set of strings

| <i>Operation</i> | <i>Definition</i> |
|------------------|---|
| Union | $L \cup M = \{ s \mid s \in L \text{ or } s \in M \}$ |
| Concatenation | $LM = \{ st \mid s \in L \text{ and } t \in M \}$ |
| Kleene closure | $L^* = \bigcup_{i=0, \infty} L^i$ |
| Positive closure | $L^+ = \bigcup_{i=1, \infty} L^i$ |

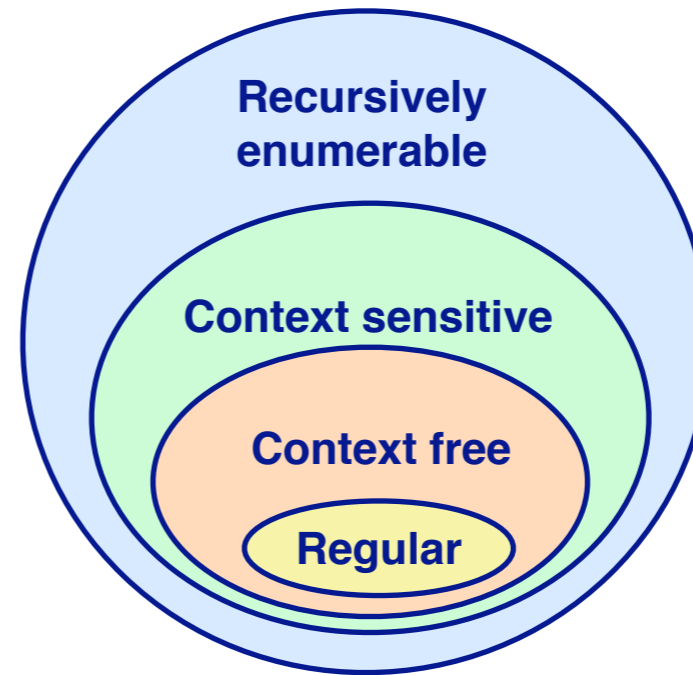
How do you define a language?

Recognizer.

Production grammar.

Production Grammars

- > Powerful formalism for language description
 - Non-terminals (A, B)
 - Terminals (a,b)
 - Production rules ($A \rightarrow abA$)
 - Start symbol (S_0)
- > Rewriting



Detail: The Chomsky Hierarchy

> Type 0: $\alpha \rightarrow \beta$

—Unrestricted grammars generate *recursively enumerable languages*. Minimal requirement for recognizer: Turing machine.

> Type 1: $\alpha A \beta \rightarrow \alpha \gamma \beta$

—Context-sensitive grammars generate *context-sensitive languages*, recognizable by linear bounded automata

> Type 2: $A \rightarrow \gamma$

—Context-free grammars generate *context-free languages*, recognizable by non-deterministic push-down automata

> Type 3: $A \rightarrow a$ and $A \rightarrow aB$

—Regular grammars generate *regular languages*, recognizable by finite state automata

NB: A is a non-terminal; α, β, γ are strings of terminals and non-terminals

10

Individual identifiers in a classical programming language form a regular language.

The language is on the other hand **context free** most of the time.

Grammars for regular languages

Regular grammars generate regular languages

Definition:

In a *regular grammar*, all productions have one of two forms:

1. $A \rightarrow aA$
2. $A \rightarrow a$

where A is any non-terminal and a is any terminal symbol

These are also called type 3 grammars (Chomsky)

Regular languages can be described by *Regular Expressions*

Regular expressions (RE) over an alphabet Σ :

1. ϵ is a RE denoting the set $\{\epsilon\}$
2. If $a \in \Sigma$, then a is a RE denoting $\{a\}$
3. If r and s are REs denoting $L(r)$ and $L(s)$, then:
 - > $(r) | (s)$ is a RE denoting $L(r) \cup L(s)$
 - > $(r)(s)$ is a RE denoting $L(r)L(s)$
 - > $(r)^*$ is a RE denoting $L(r)^*$

We adopt a *precedence* for operators: *Kleene closure*, then *concatenation*, then *alternation* as the order of precedence.

For any RE r , there exists a grammar g such that $L(r) = L(g)$

12

Epsilon (the set with the “empty” string)

As you can see, we don't define a^+ or $[a]$

Patterns are often specified as *regular languages*.

Notations used to describe a regular language (or a regular set) include both *regular expressions* and *regular grammars*

Examples

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

> $a \mid b$ denotes $\{a,b\}$

> $(a \mid b)(a \mid b)$ denotes $\{aa,ab,ba,bb\}$

> a^* denotes $\{\varepsilon,a,aa,aaa,\dots\}$

> $(a \mid b)^*$ denotes the set of all strings of a's and b's
(including ε)

> Universit(ä | ae)t Bern(e |) ...

Algebraic properties of REs

| | |
|--------------------------------------|--|
| $r s = s r$ | $ $ is commutative |
| $r (s t) = (r s) t$ | $ $ is associative |
| $r(st) = (rs)t$ | concatenation is associative |
| $r(s t) = rs rt$ $(s t)r = sr tr$ | concatenation distributes over $ $ |
| $\epsilon r = r$ $r\epsilon = r$ | ϵ is the identity for concatenation |
| $r^* = (r \epsilon)^*$ | ϵ is contained in * |
| $r^{**} = r^*$ | * is idempotent |

Examples of using REs to specify lexical patterns

identifiers

$letter \rightarrow (a | b | c | \dots | z | A | B | C | \dots | Z)$

$digit \rightarrow (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)$

$id \rightarrow letter (letter | digit)^*$

numbers

$integer \rightarrow (+ | - | \epsilon) (0 | (1 | 2 | 3 | \dots | 9) digit^*)$

$decimal \rightarrow integer . (digit)^*$

$real \rightarrow (integer | decimal) \mathbb{E} (+ | -) digit^*$

$complex \rightarrow '(real , real)'$

Numbers can get much more complicated.

Most programming language tokens can be described with REs.

Roadmap



- > Introduction
- > Regular languages
- > **Finite automata recognizers**
- > From RE to DFAs and back again
- > Limits of regular languages

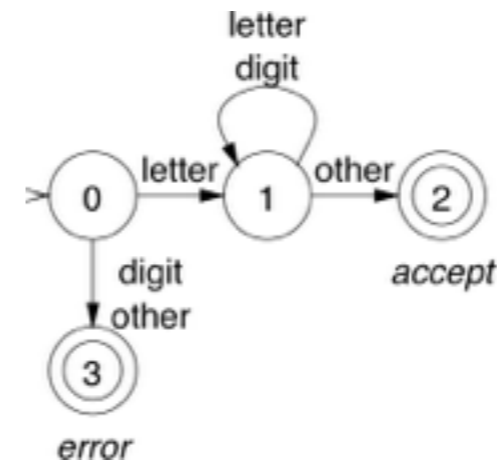
REs are cool for specifying.

FAs are good for implementing REs.

Recognizers

$letter \rightarrow (a | b | c | \dots | z | A | B | C | \dots | Z)$
 $digit \rightarrow (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)$
 $id \rightarrow letter (letter | digit)^*$

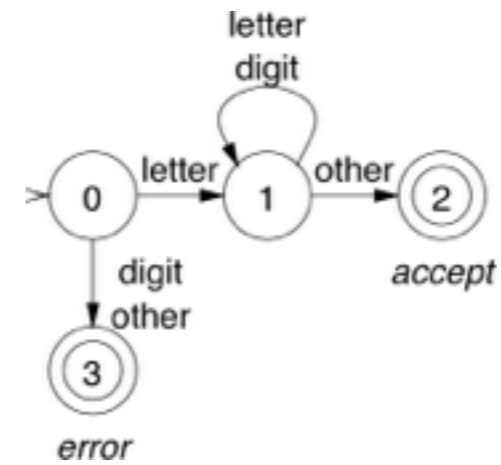
From a regular expression we can construct a *deterministic finite automaton* (DFA)



DFA for recognizing an identifier.
why D? why F? why A?

Code for the recognizer

```
char ← next_char();
state ← 0;      /* code for state 0 */
done ← false;
token_value ← "" /* empty string */
while( not done ) {
  class ← char_class[char];
  state ← next_state[class,state];
  switch(state) {
    case 1: /* building an id */
      token_value ← token_value + char;
      char ← next_char();
      break;
    case 2: /* accept state */
      token_type = identifier;
      done = true;
      break;
    case 3: /* error */
      token_type = error;
      done = true;
      break;
  }
}
return token_type;
```



I.e., encode the transitions in the next_state matrix

Tables for the recognizer

Two tables control the recognizer

| | | | | | |
|------------|--------------|--------|--------|-------|-------|
| char_class | <i>char</i> | a-z | A-Z | 0-9 | other |
| | <i>value</i> | letter | letter | digit | other |

| | | | | | |
|------------|--------|---|---|---|---|
| next_state | | 0 | 1 | 2 | 3 |
| | letter | 1 | 1 | — | — |
| | digit | 3 | 1 | — | — |
| | other | 3 | 2 | — | — |

To change languages, we can just change tables

Automatic construction

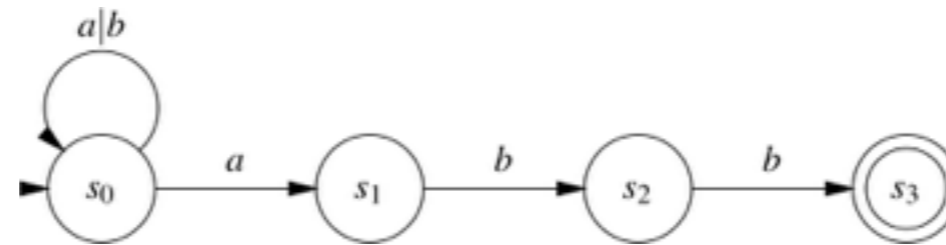
- > *Scanner generators* automatically construct code from regular expression-like descriptions
 - construct a DFA
 - use *state minimization* techniques
 - emit code for the scanner (table driven or direct code)

- > A key issue in automation is an interface to the parser

- > *lex* is a scanner generator supplied with UNIX
 - emits C code for scanner
 - provides macro definitions for each token (used in the parser)
 - nowadays JavaCC is more popular

NFA example

What about the RE $(a|b)^*abb$?



State s_0 has multiple transitions on a !

This is a non-deterministic finite automaton

Review: Finite Automata

A non-deterministic finite automaton (**NFA**) consists of:

1. a set of *states* $S = \{ s_0, \dots, s_n \}$
2. a set of *input symbols* Σ (the alphabet)
3. a transition function *move* (δ) mapping state-symbol pairs to sets of states
4. a distinguished *start state* s_0
5. a set of distinguished *accepting (final) states* F

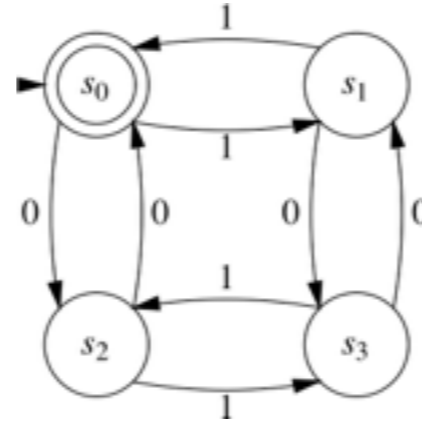
A Deterministic Finite Automaton (**DFA**) is a special case of an NFA:

1. no state has a ϵ -transition, and
2. for each state s and input symbol a , there is at most one edge labeled a leaving s .

A DFA accepts x iff there exists a *unique* path through the transition graph from the s_0 to an accepting state such that the labels along the edges spell x .

DFA example

Example: the set of strings containing an even number of zeros and an even number of ones



The RE is $(00 \mid 11)^*((01 \mid 10)(00 \mid 11)^*(01 \mid 10)(00 \mid 11)^*)^*$

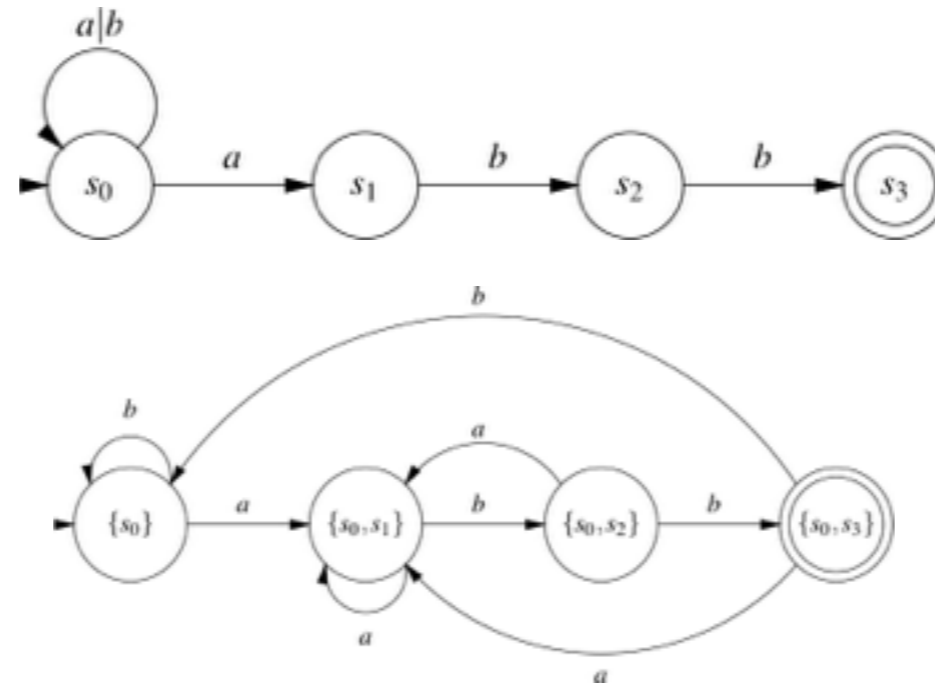
Note how the RE walks through the DFA.

NB: The states capture whether there are an even number of 0s or 1s => 4 possible states.

DFAs and NFAs are equivalent

1. DFAs are a subset of NFAs
2. Any NFA can be converted into a DFA, by *simulating sets of simultaneous states*:
 - each DFA state corresponds to a set of NFA states
 - NB: possible exponential blowup

NFA to DFA using the subset construction



Roadmap

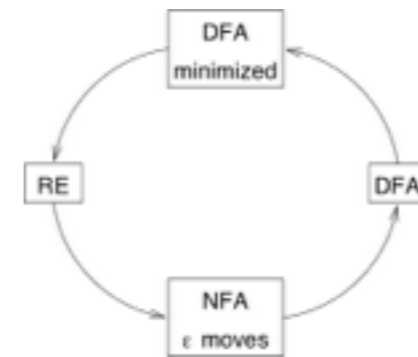


- > Introduction
- > Regular languages
- > Finite automata recognizers
- > **From RE to DFAs and back again**
- > Limits of regular languages

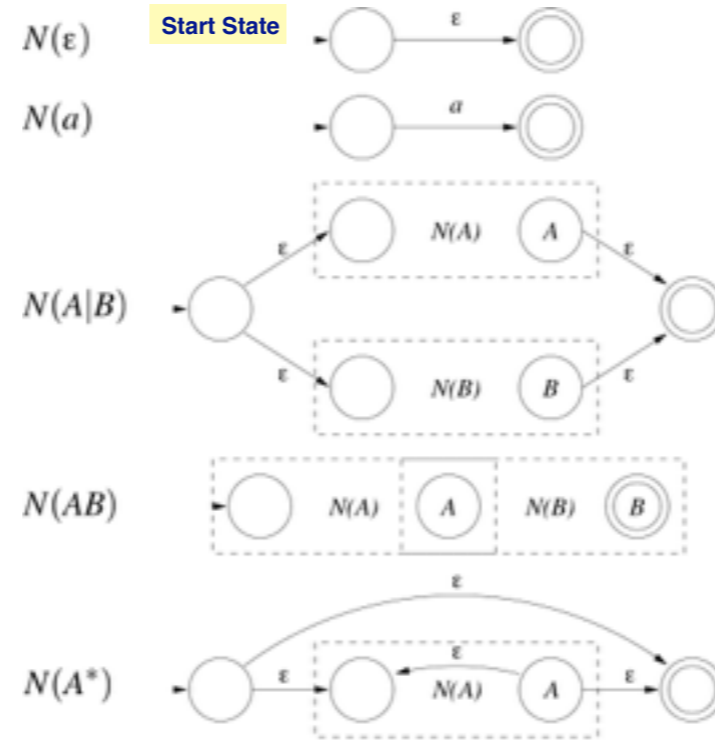
Constructing a DFA from a RE

- > RE → NFA
 - Build NFA for each term; connect with ϵ moves
- > NFA → DFA
 - Simulate the NFA using the subset construction
- > DFA → minimized DFA
 - Merge equivalent states

- > DFA → RE
 - Construct $R_{ij}^k = R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \cup R_{ij}^{k-1}$
 - Or convert via Generalized NFA (GNFA)

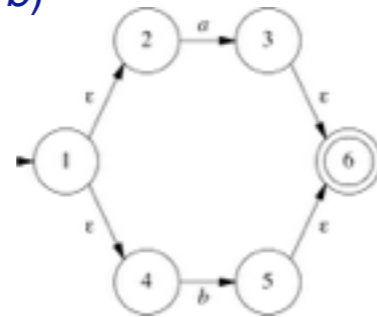


RE to NFA

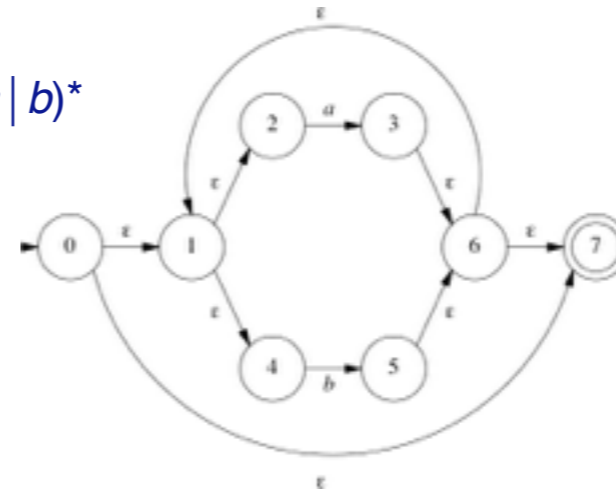


RE to NFA example: $(a | b)^*abb$

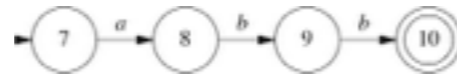
$(a | b)$



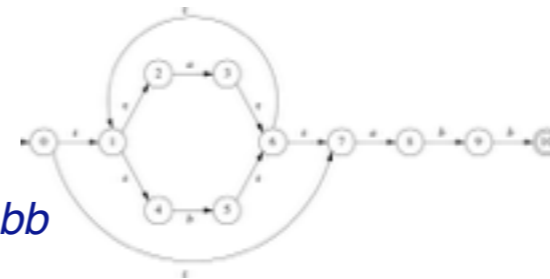
$(a | b)^*$



abb



$(a | b)^*abb$



NFA to DFA: the subset construction

Input: NFA N

Output: DFA D with states S_D and transitions T_D such that $L(D) = L(N)$

Method: Let s be a state in N and P be a set of states. Use the following operations:

- > ϵ -closure(s) — set of states of N reachable from s by ϵ transitions alone
- > ϵ -closure(P) — set of states of N reachable from some s in P by ϵ transitions alone
- > $\text{move}(T, a)$ — set of states of N to which there is a transition on input a from some s in P

add state $P = \epsilon$ -closure(s_0)
unmarked to S_D

while \exists unmarked state P in S_D

mark P

for each input symbol a

$U = \epsilon$ -closure($\text{move}(P, a)$)

if $U \notin S_D$

then add U unmarked to S_D

$T_D[P, a] = U$

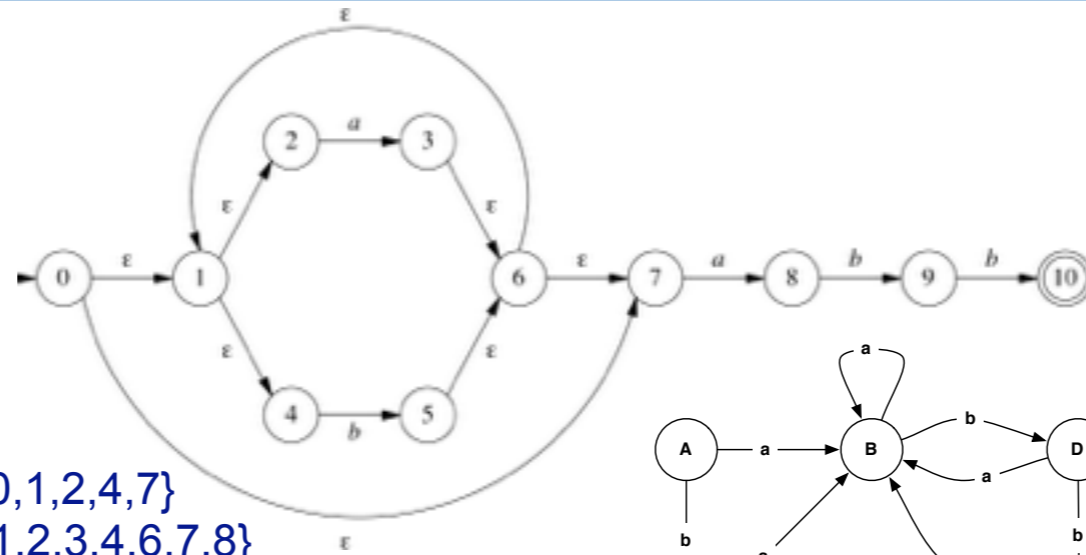
end for

end while

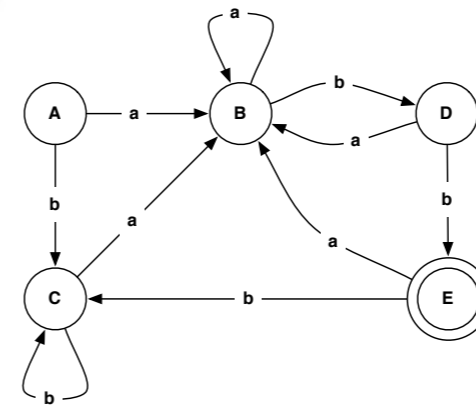
ϵ -closure(s_0) is the start state of D

A state of D is accepting if it contains an accepting state of N

NFA to DFA using subset construction: example



A = {0,1,2,4,7}
B = {1,2,3,4,6,7,8}
C = {1,2,4,5,6,7}
D = {1,2,4,5,6,7,9}
E = {1,2,4,5,6,7,10}

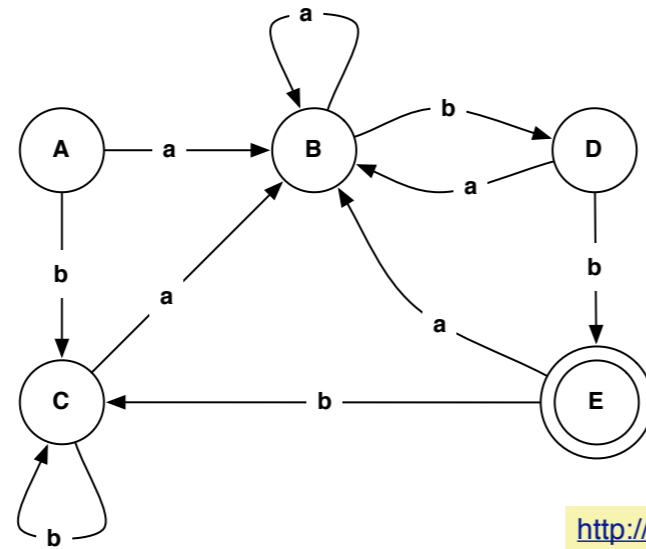


31

Are NFAs more powerful than DFAs?
Fewer states and easier to construct!
But the transformation is not minimal.

DFA Minimization

Theorem: For each regular language that can be accepted by a DFA, there exists a DFA with a minimum number of states.



Minimization approach:
merge *equivalent* states.

States A and C are indistinguishable, so they can be merged!

http://en.wikipedia.org/wiki/DFA_minimization

32

After b^*a we always end up in state B.

DFA Minimization algorithm

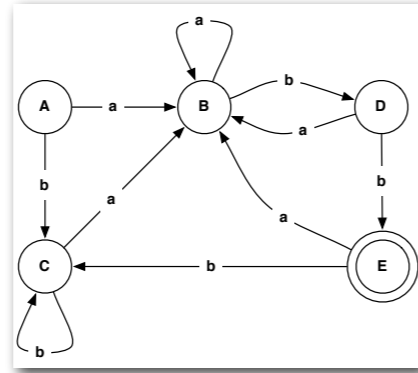
- > Create lower-triangular table DISTINCT, initially blank
- > For every pair of states (p, q) :
 - If p is final and q is not, or vice versa
 - $DISTINCT(p, q) = \varepsilon$
- > Loop until no change for an iteration:
 - For every pair of states (p, q) and each symbol α
 - If $DISTINCT(p, q)$ is blank and $DISTINCT(\delta(p, \alpha), \delta(q, \alpha))$ is not blank
 - $DISTINCT(p, q) = \alpha$
- > Combine all states that are not distinct

33

Distinguish final state from all others. Then take single steps to check what is distinguishable. The intuition:

- if one state is final and the other not, then they are clearly distinct
- otherwise, for every (state, state, symbol) tuple we see whether the δ is in DISTINCT

Minimization in action



C and A are *indistinguishable*
so can be merged

| | | | | | |
|---|---|---|---|---|---|
| A | | | | | |
| B | | | | | |
| C | | | | | |
| D | | | | | |
| E | | | | | |
| | A | B | C | D | E |

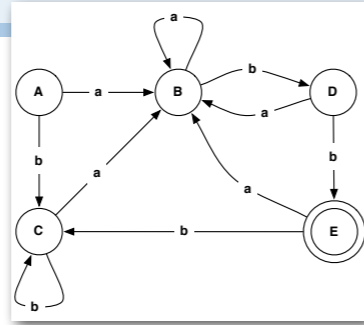
| | | | | | |
|---|------------|------------|------------|------------|---|
| A | | | | | |
| B | | | | | |
| C | | | | | |
| D | | | | | |
| E | ϵ | ϵ | ϵ | ϵ | |
| | A | B | C | D | E |

| | | | | | |
|---|------------|------------|------------|------------|---|
| A | | | | | |
| B | | | | | |
| C | | | | | |
| D | b | b | b | | |
| E | ϵ | ϵ | ϵ | ϵ | |
| | A | B | C | D | E |

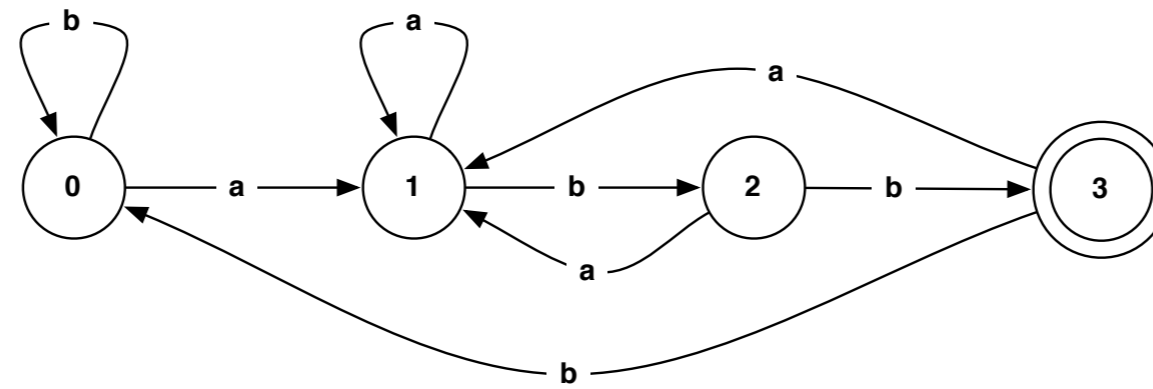
| | | | | | |
|---|------------|------------|------------|------------|---|
| A | | | | | |
| B | b | | | | |
| C | | b | | | |
| D | b | b | b | | |
| E | ϵ | ϵ | ϵ | ϵ | |
| | A | B | C | D | E |

0. initial state. 1. E is final, so different from others.
2. Only a "b" step from D leads to non-blank space.
3. B can make a "b" step to D, so differs from A and C.
4. A and C are indistinguishable. (An "a" takes both to B and "b" takes both to C.)

DFA Minimization example



It is easy to see that this is in fact the minimal DFA for $(a|b)^*abb\dots$



35

Actually it is easy to see that this is the minimal DFA:

Start with the path abb . This gives us 4 states. Now add the missing arrows.

Any a transition brings us to state 1, since we must follow with bb .

Any b not in the path brings us back to state 0, since we must follow with abb .

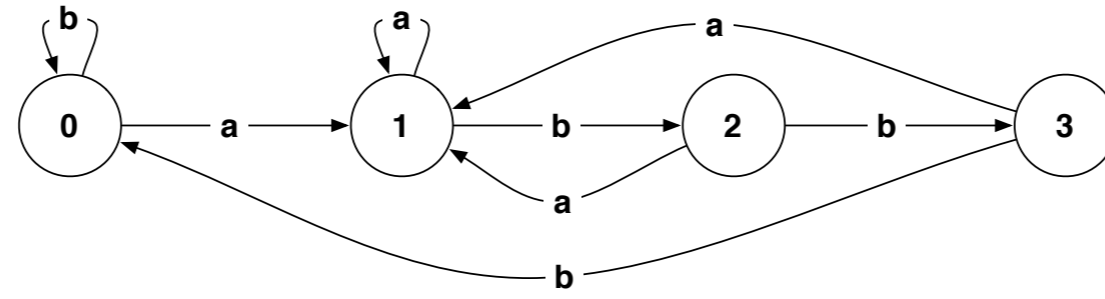
DFA to RE via GNFA

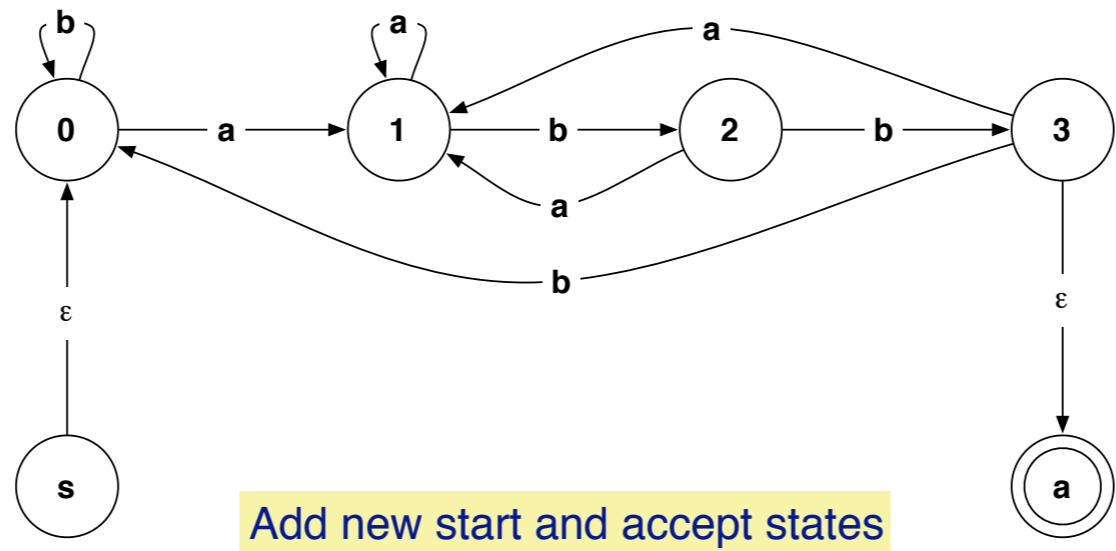
- > A *Generalized NFA* is an NFA where transitions may have any RE as labels
- > Conversion algorithm:
 1. *Add a new start state and accept state* with ϵ -transitions to/from the old start/end states
 2. *Merge multiple transitions* between two states to a single RE choice transition
 3. *Add empty \emptyset -transitions* between states where missing
 4. *Iteratively “rip out” old states* and replace “dangling transitions” with appropriately labeled transitions between remaining states
 5. *STOP when all old states are gone* and only the new start and accept states remain

GNFA conversion algorithm

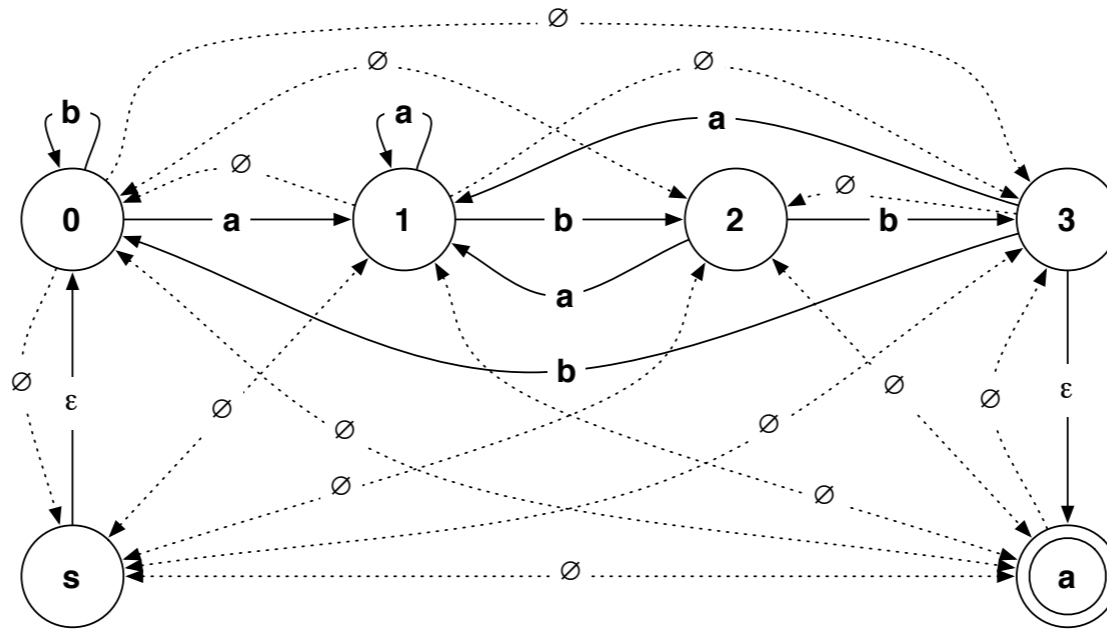
1. Let k be the number of states of G , $k \geq 2$
2. If $k=2$, then RE is the label found between q_s and q_a (start and accept states of G)
3. While $k > 2$, select $q_{rip} \neq q_s$ or q_a
 - $Q' = Q - \{q_{rip}\}$
 - For any $q_i \in Q' - \{q_a\}$ let $\delta'(q_i, q_j) = R_1 R_2^* R_3 \cup R_4$ where:
 $R_1 = \delta'(q_i, q_{rip})$, $R_2 = \delta'(q_{rip}, q_{rip})$, $R_3 = \delta'(q_{rip}, q_j)$, $R_4 = \delta'(q_i, q_j)$
 - Replace G by G'

The initial DFA



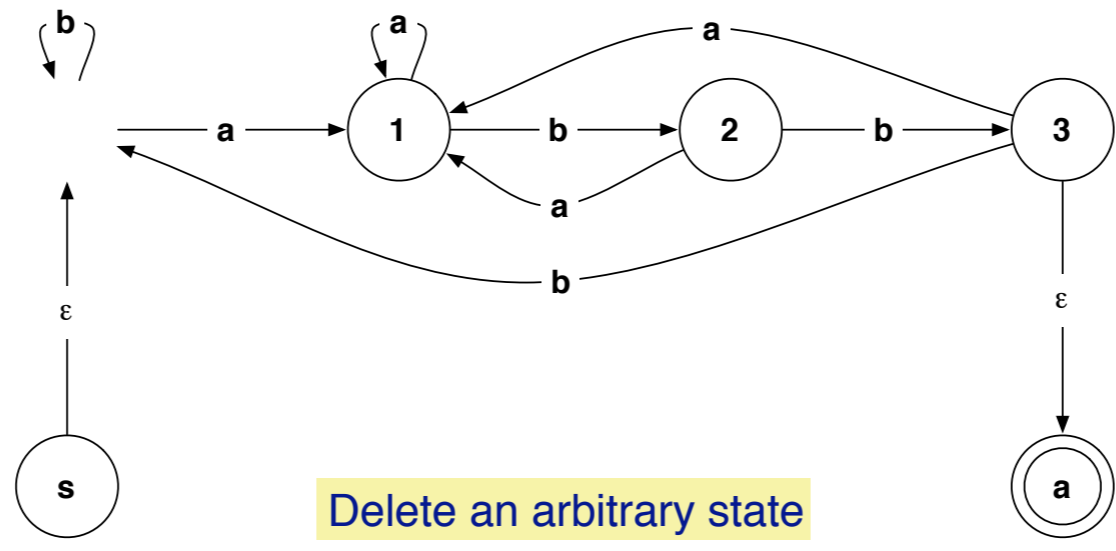


Add new start and accept states

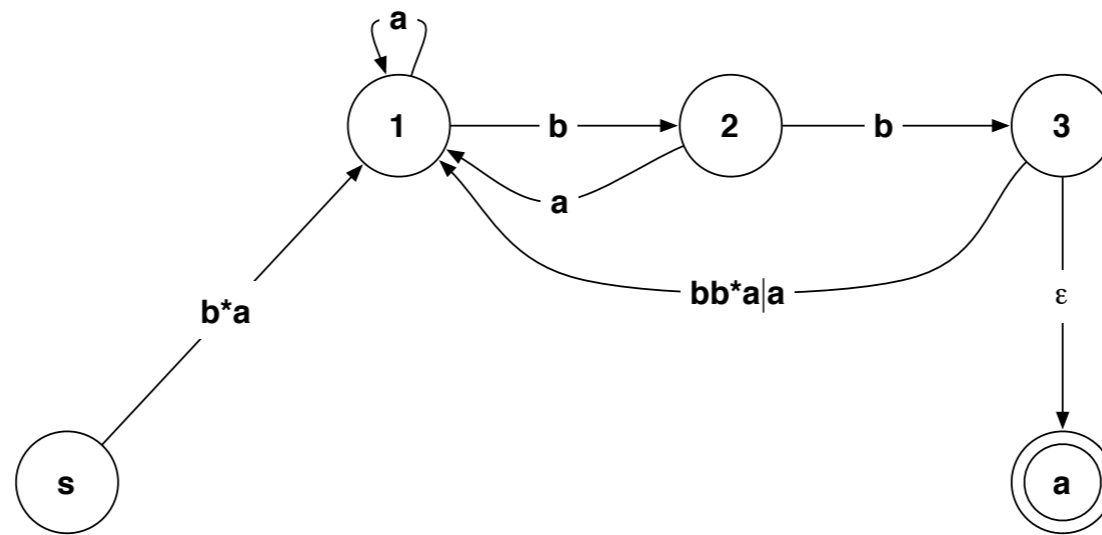


Add missing empty transitions
(we'll just pretend they're there)

This means “you can't get there from here”



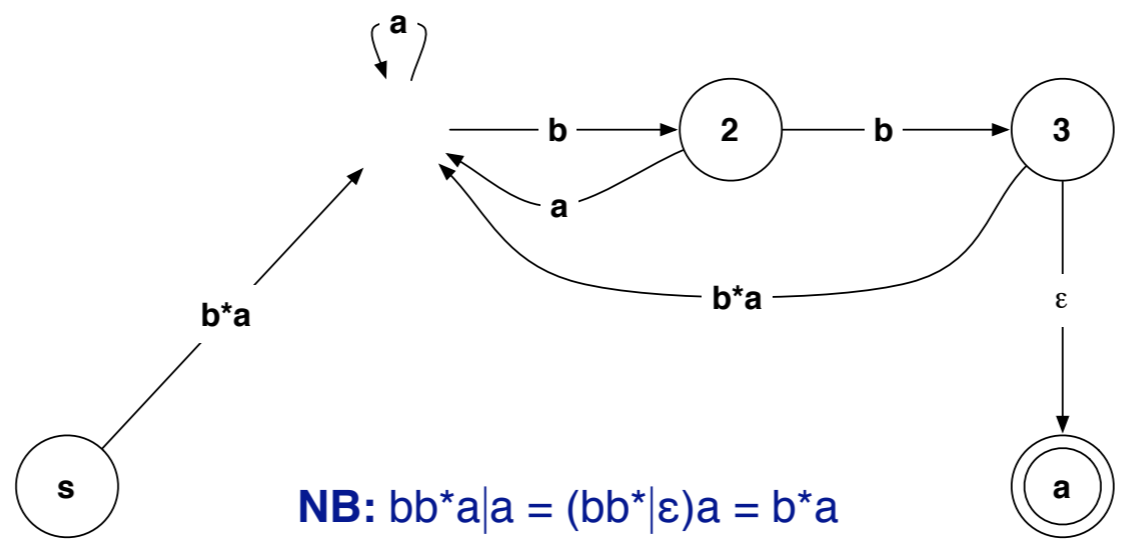
Delete an arbitrary state



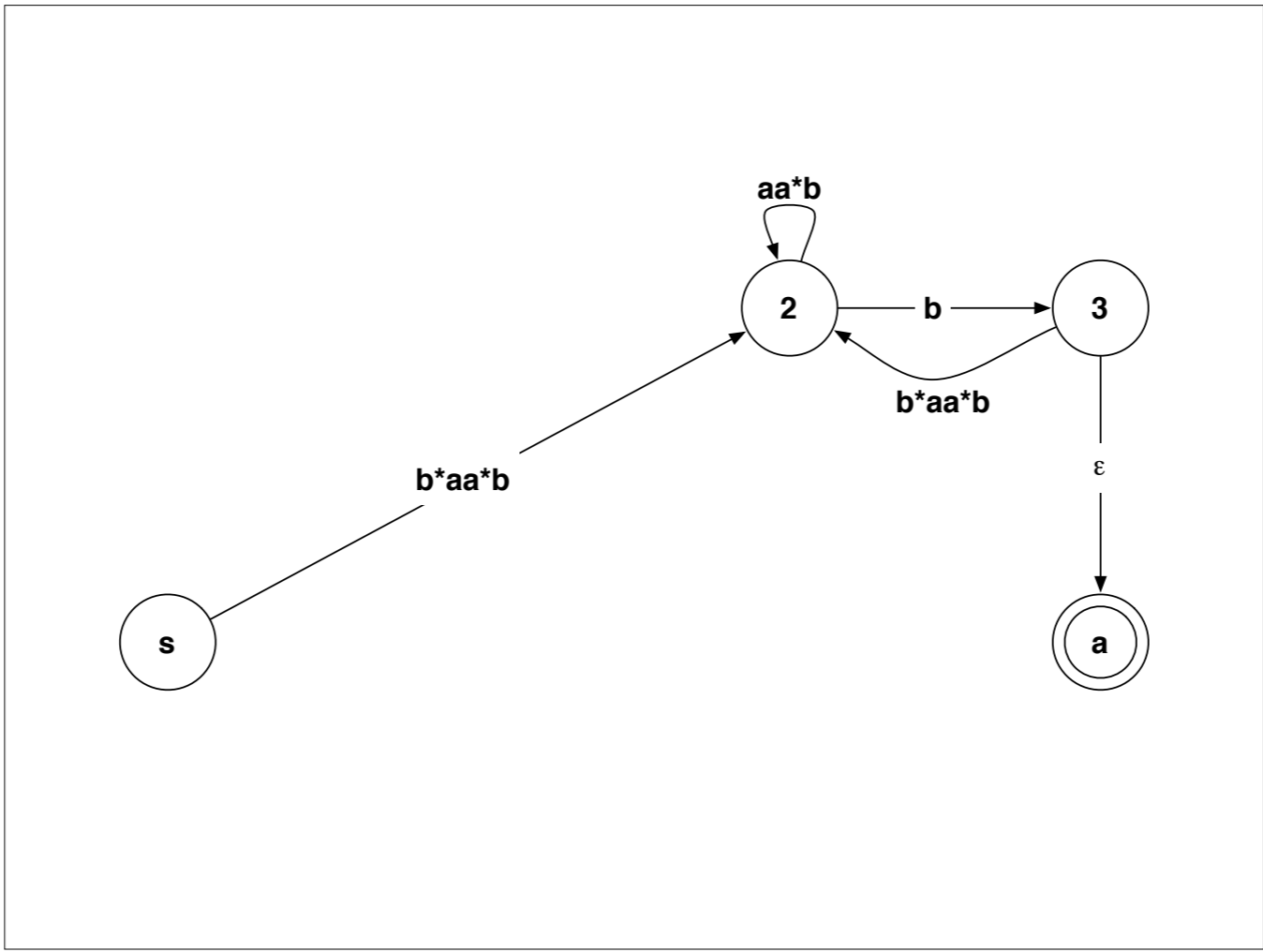
Fix dangling transitions $s \rightarrow 1$ and $3 \rightarrow 1$
 Don't forget to merge the existing transitions!

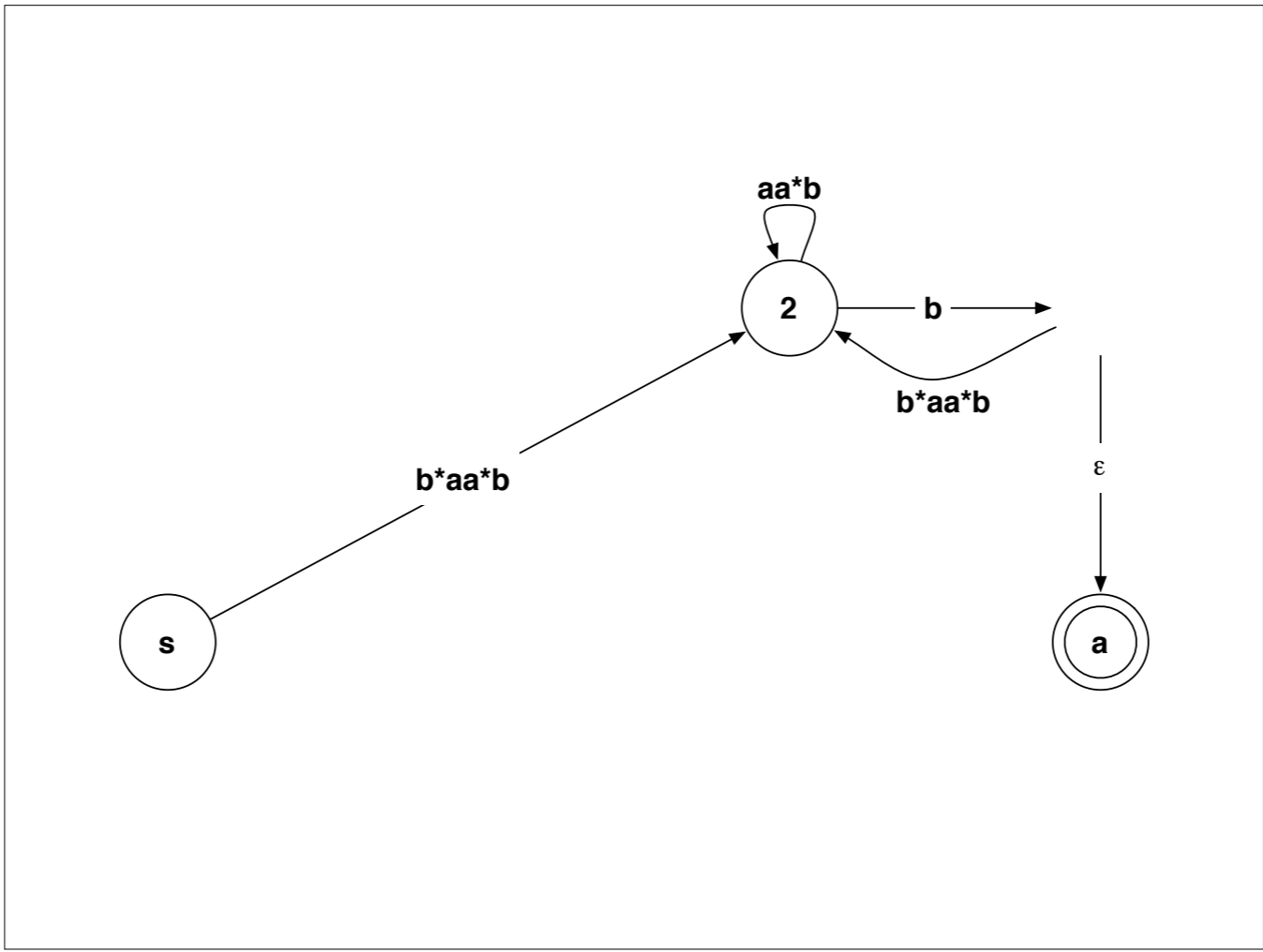
NB: The path from (3) to (1) merges the old path bb^*a from $(3) \rightarrow (0) \rightarrow (1)$ and the path a from $(3) \rightarrow (1)$.

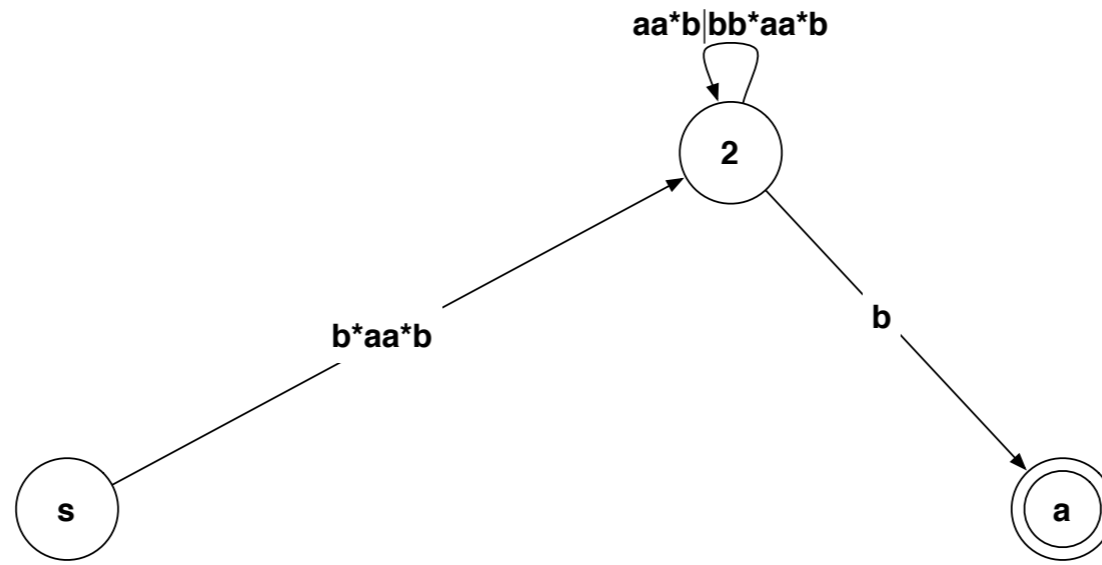
Simplify the RE
Delete another state



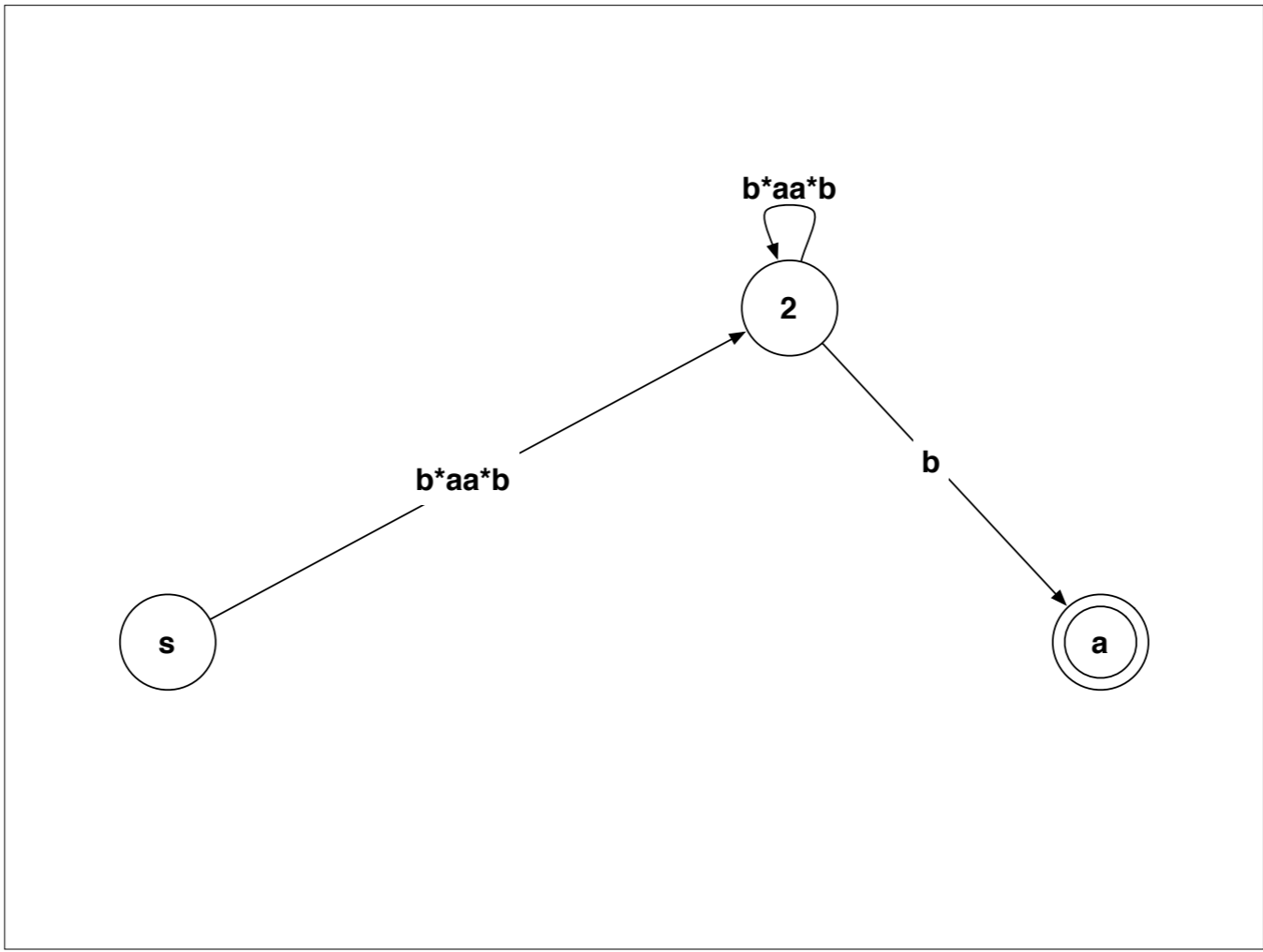
NB: $bb^*a|a = (bb^*|\epsilon)a = b^*a$

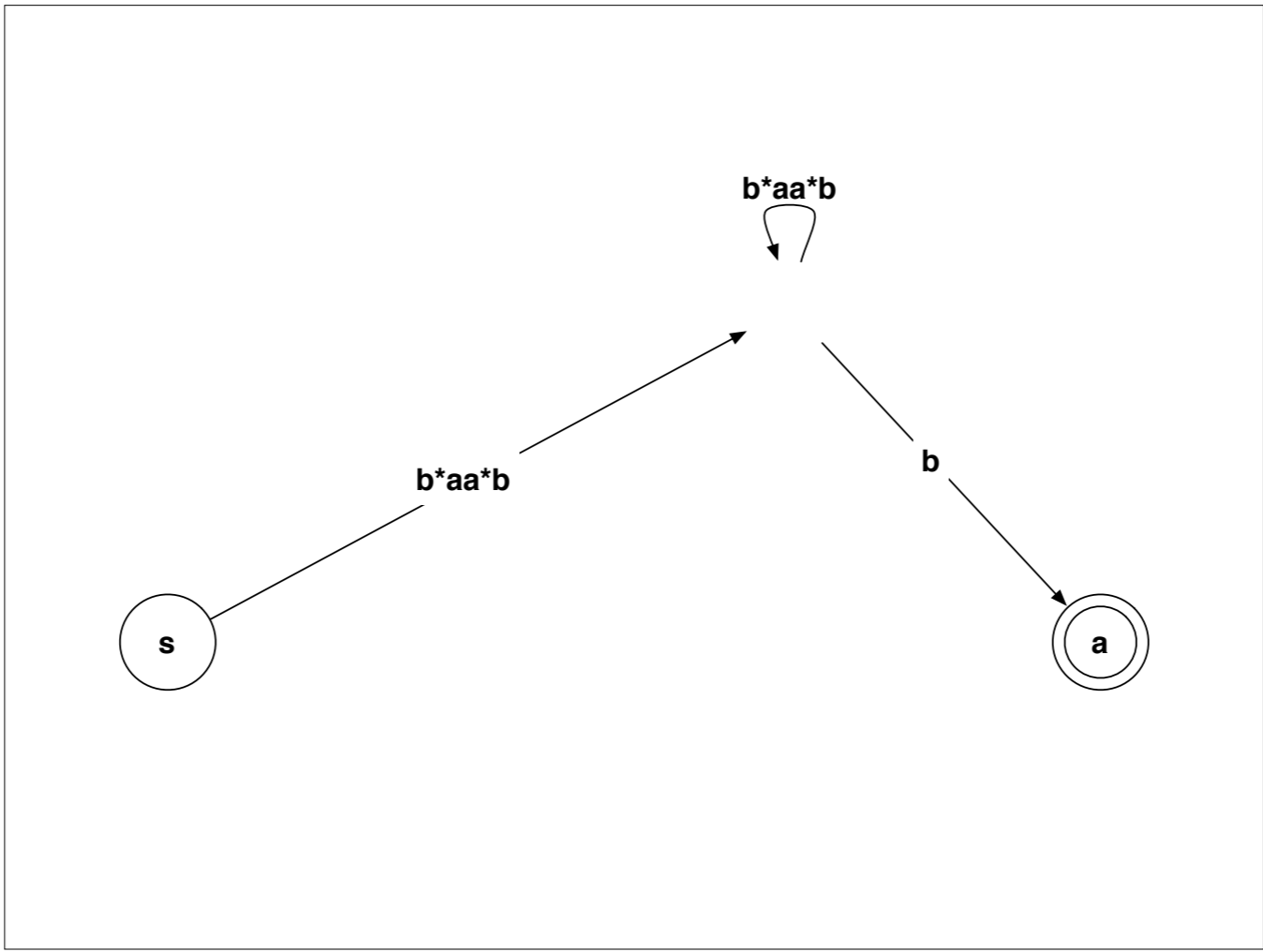




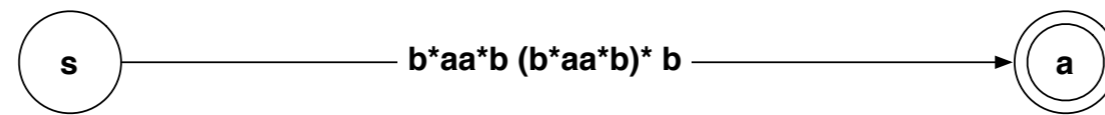


NB: $aa^*b|bb^*aa^*b = (\epsilon|bb^*)aa^*b = b^*aa^*b$





Hm ... not what we expected



Note that $b^*aa^*b = b^*a^*ab$

And so $b^*aa^*b (b^*aa^*b)^* b = (b^*a^*ab)^* b^*a^*abb$

It remains to be shown that $(b^*a^*ab)^* b^*a^* = (alb)^* \dots$

$b^*aa^*b (b^*aa^*b)^* b = (a|b)^*abb$?

> *We can rewrite:*

— $b^*aa^*b (b^*aa^*b)^* b$

— $b^*a^*ab (b^*a^*ab)^* b$

— $(b^*a^*ab)^* b^*a^* abb$

> *But does this hold?*

— $(b^*a^*ab)^* b^*a^* = (a|b)^*$

We can show that the minimal DFAs for these REs are isomorphic ...

Proof: Split any string in $(ab)^*$ by occurrences of ab . This will match $(Xab)^*X$, where X does not contain ab . X is clearly b^*a^* . QED

Proof #2 by @grammarware: $(b^*a^*ab)^*b^*a^* = (b^*a^*b)^*b^*a^* = b^*(a^*b^+)^*a^* = b^*(b^*|(a^*b^+)^*)a^* = b^*(b^*|(a^*b^+)^*|a^*)a^* = b^*(ab)^*a^* = (ab)^*a^* = (ab)^*$

Roadmap



- > Introduction
- > Regular languages
- > Finite automata recognizers
- > From RE to DFAs and back again
- > **Limits of regular languages**

Limits of regular languages

Not all languages are regular!

One cannot construct DFAs to recognize these languages:

$$L = \{ p^k q^k \}$$
$$L = \{ w c w^r \mid w \in \Sigma^*, w^r \text{ is } w \text{ reversed} \}$$

In general, DFAs cannot count!

However, one *can* construct DFAs for:

- Alternating 0's and 1's:

$$(\varepsilon \mid 1)(01)^*(\varepsilon \mid 0)$$

- Sets of pairs of 0's and 1's

$$(01 \mid 10)^+$$

So, what is hard?

Certain language features can cause problems:

- > **Reserved words**
 - PL/I had no reserved words
 - `if then then then = else; else else = then`
- > **Significant blanks**
 - FORTRAN and Algol68 ignore blanks
 - `do 10 i = 1,25`
 - `do 10 i = 1.25`
- > **String constants**
 - Special characters in strings
 - Newline, tab, quote, comment delimiter
- > **Finite limits**
 - Some languages limit identifier lengths
 - Add state to count length
 - FORTRAN 66 — 6 characters(!)

How bad can it get?

```
1      INTEGERFUNCTIONA
2      PARAMETER(A=6,B=2)
3      IMPLICIT CHARACTER*(A-B)(A-B)
4      INTEGER FORMAT(10),IF(10),DO9E1
5      100  FORMAT(4H)=(3)
6      200  FORMAT(4 )=(3)
7      DO9E1=1
8      DO9E1=1,2
9          IF(X)=1
10         IF(X)H=1
11         IF(X)300,200
12  300    CONTINUE
13      END
14      C    this is a comment
        $ FILE(1)
14      END
```

Example due to Dr. F.K. Zadeck of IBM Corporation

*Compiler needs context
to distinguish variables
from control constructs!*

What you should know!

- ✎ What are the key responsibilities of a scanner?*
- ✎ What is a formal language? What are operators over languages?*
- ✎ What is a regular language?*
- ✎ Why are regular languages interesting for defining scanners?*
- ✎ What is the difference between a deterministic and a non-deterministic finite automaton?*
- ✎ How can you generate a DFA recognizer from a regular expression?*
- ✎ Why aren't regular languages expressive enough for parsing?*

Can you answer these questions?

- ✎ Why do compilers separate scanning from parsing?*
- ✎ Why doesn't NFA → DFA translation normally result in an exponential increase in the number of states?*
- ✎ Why is it necessary to minimize states after translation a NFA to a DFA?*
- ✎ How would you program a scanner for a language like FORTRAN?*



Attribution-ShareAlike 4.0 International (CC BY-SA 4.0)

You are free to:

Share — copy and redistribute the material in any medium or format

Adapt — remix, transform, and build upon the material for any purpose, even commercially.

The licensor cannot revoke these freedoms as long as you follow the license terms.

Under the following terms:



Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



ShareAlike — If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.

No additional restrictions — You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

<http://creativecommons.org/licenses/by-sa/4.0/>