# Lexical Analysis

CMPT 379: Compilers

Instructor: Anoop Sarkar

anoopsarkar.github.io/compilers-class

# Building a Lexical Analyzer

- Token ⇒ Pattern
- Pattern ⇒ Regular Expression
- Regular Expression ⇒ NFA
- NFA ⇒ DFA
- DFA ⇒ Table-driven implementation of DFA

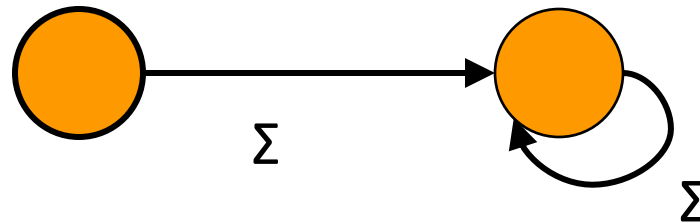# Thompson's construction

- Converts regexps to equivalent NFA
- Six simple rules
  - Empty language
  - Symbols ($\Sigma$)
  - Empty String ($\varepsilon$)
  - Alternation ($r_1$ or $r_2$)
  - Concatenation ($r_1$ followed by $r_2$)
  - Repetition ($r_1*$)

Used by Ken Thompson for pattern-based search in text editor QED (1968)
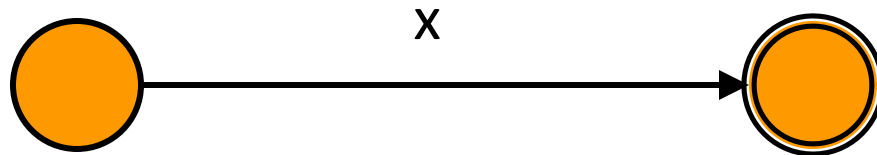
3

# Thompson Rule 0

- For the empty language φ (optionally include a *sinkhole* state)
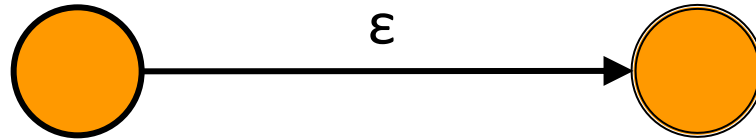
# Thompson Rule 1

- For each symbol *x* of the alphabet, there is a NFA that accepts it
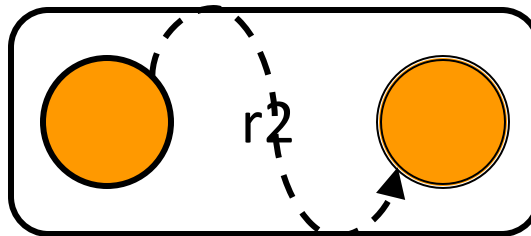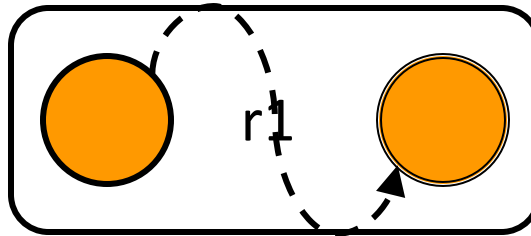
# Thompson Rule 2

- There is an NFA that accepts only ε

# Thompson Rule 3

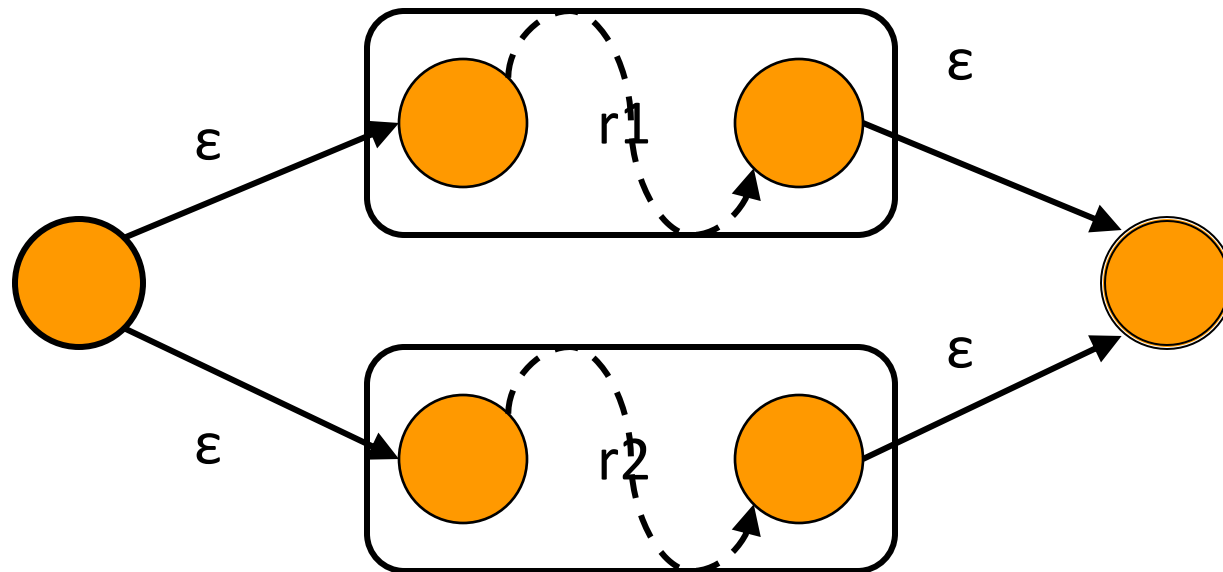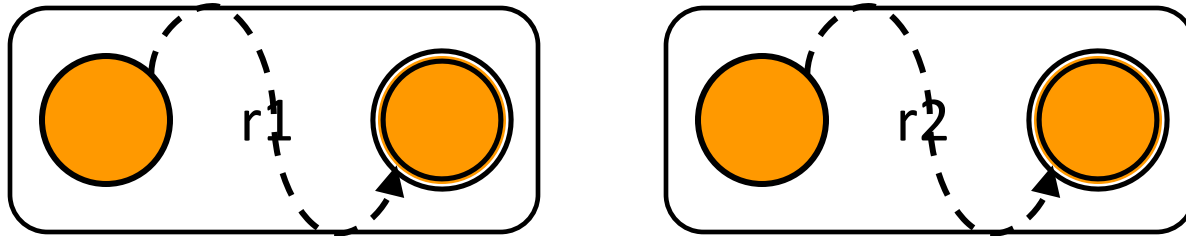- Given two NFAs for $r_1$, $r_2$, there is a NFA that accepts $r_1|r_2$

# Thompson Rule 3

- Given two NFAs for $r_1, r_2$, there is a NFA that accepts $r_1 | r_2$

# Thompson Rule 4

- Given two NFAs for $r_1, r_2$, there is a NFA that accepts $r_1 r_2$

# Thompson Rule 4

- Given two NFAs for $r_1$, $r_2$, there is a NFA that accepts $r_1 r_2$

# Thompson Rule 4

- Given two NFAs for $r_1$, $r_2$, there is a NFA that accepts $r_1 r_2$
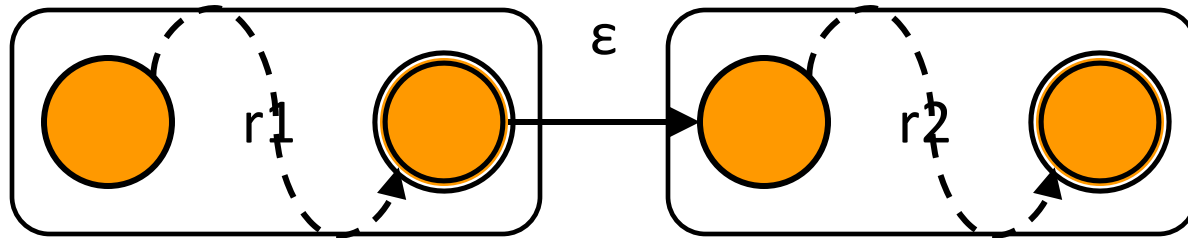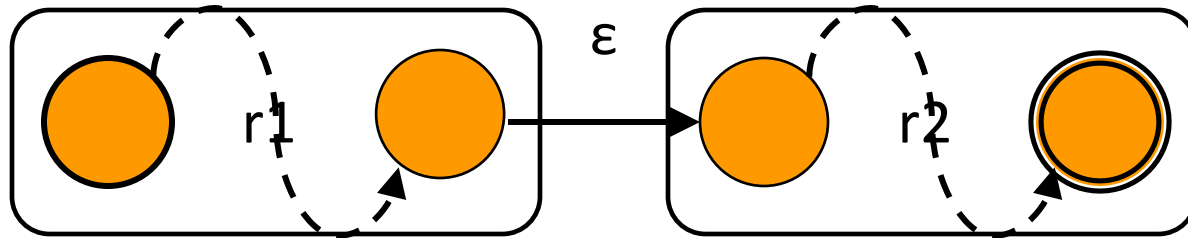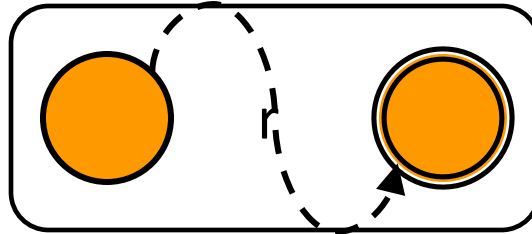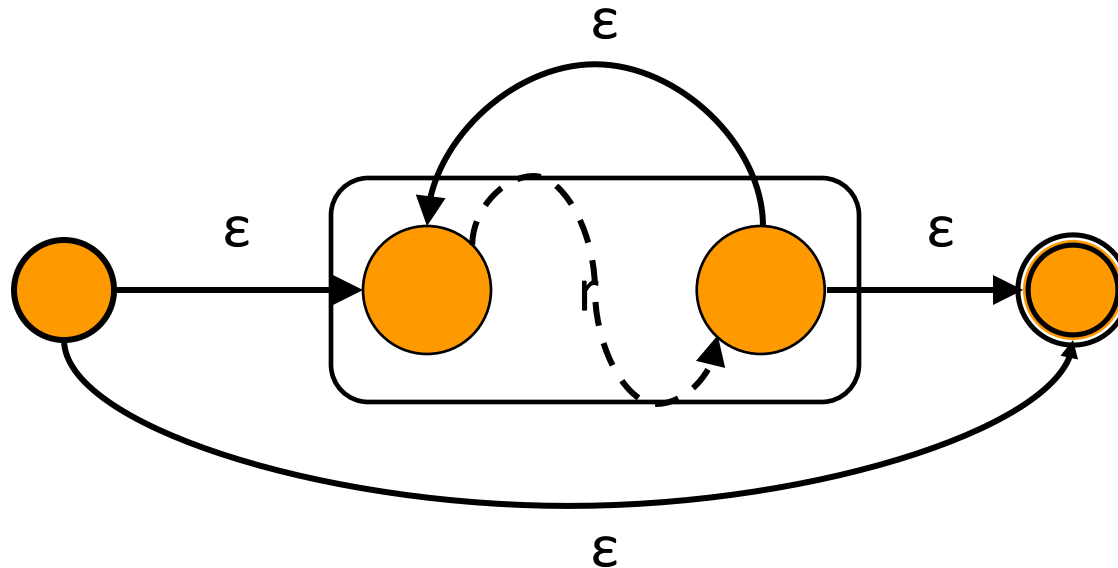
# Thompson Rule 5

- Given a NFA for r, there is an NFA that accepts *r**

# Thompson Rule 5

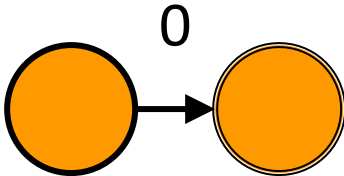- Given a NFA for r, there is an NFA that accepts *r\**

# Example

- Set of all binary strings that are divisible by four (include 0 in this set)

- Defined by the regexp: ((0|1)*00) | 0
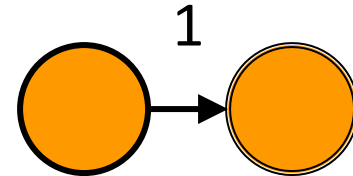
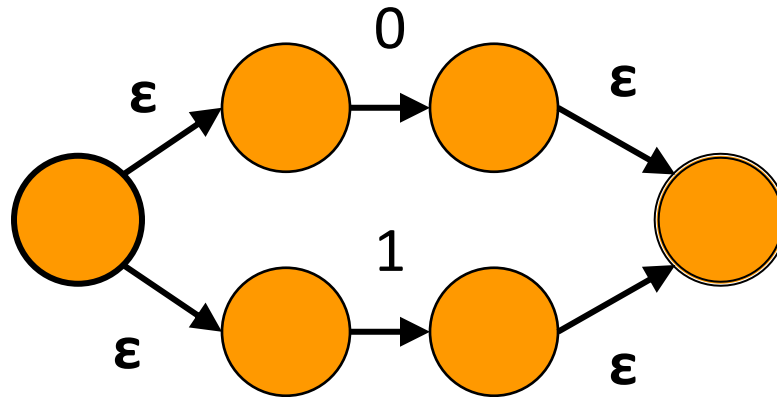- Apply Thompson's Rules to create an NFA

# Basic Blocks 0 and 1
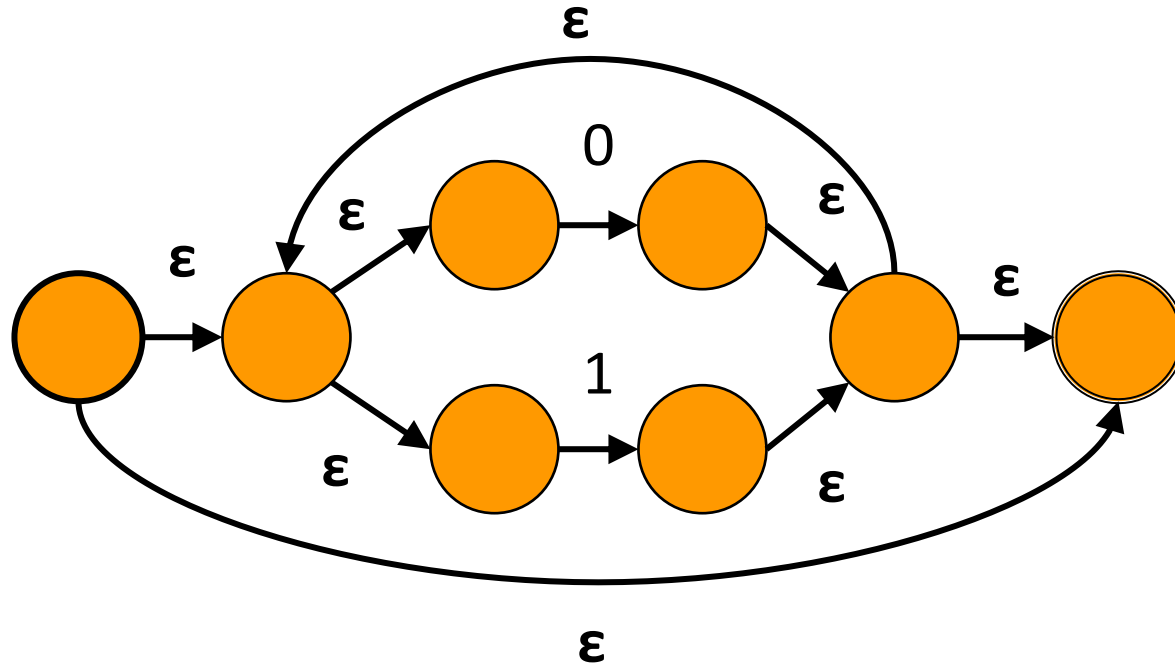
$$((0|1)*00) | 0$$

NFA for 0

NFA for 1

0

1

0|1

((0|1)*00) | 0

(0|1)*

((0|1)*00) | 0

(0|1)*00
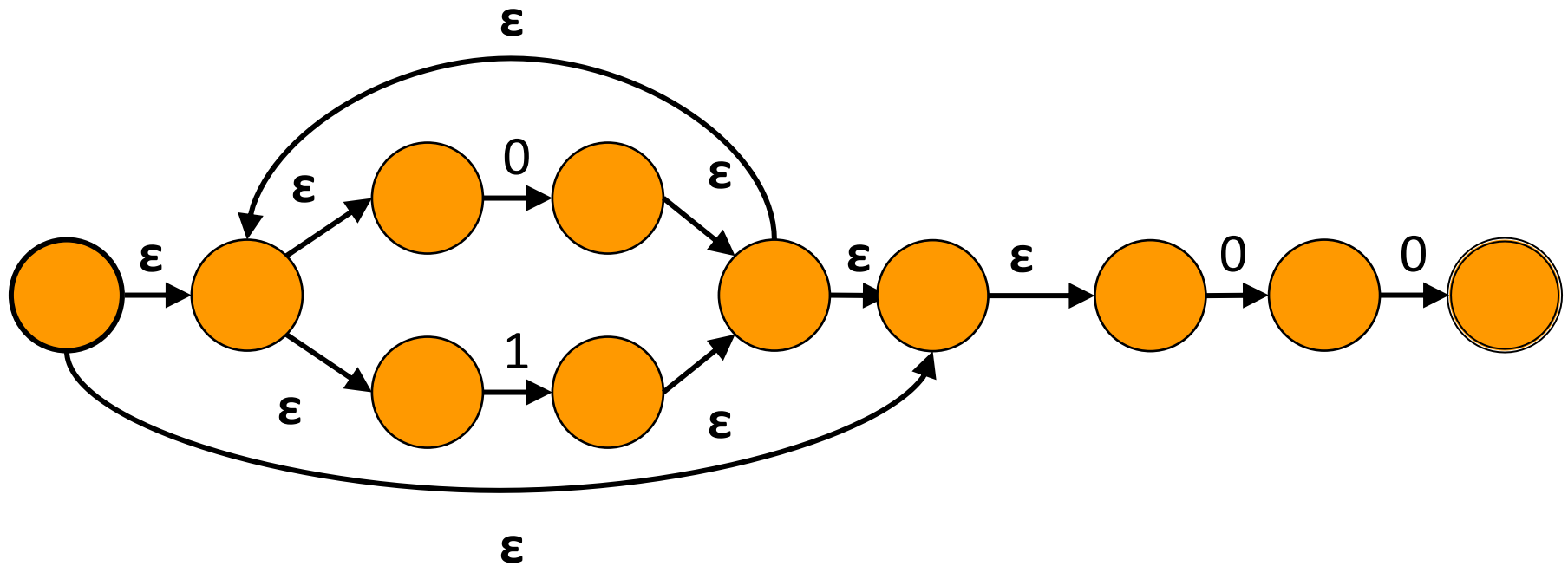
((0|1)*00) | 0

(0|1)*00

((0|1)*00) | 0

$((0|1)*00)|0$

$$((0|1)*00)|0$$

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree |
| --- |

$(a(a|b))c$

$aab|.c.$

| Post-order traversal of regexp tree |
| --- |



n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

(a(a**|**b))c

aab**|**.c.

| Post-order traversal of regexp tree |
|---|



stack

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

23

# Thompson's construction

## Converts regexps to NFA

Build NFA recursively from regexp tree

$(a(a|b))c$

Post-order traversal of regexp tree

$aab|.c.$



push n1

stack
| n1 |

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )
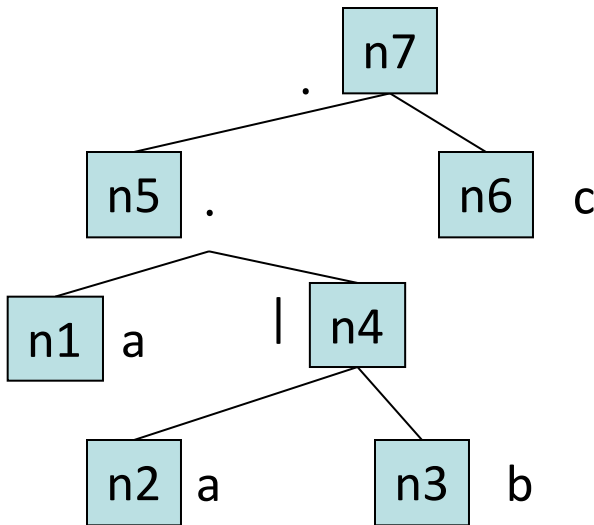
n6= nfa(c)

n7= nfa(n5, n6, **.** )

24

# Thompson's construction
## Converts regexps to NFA

**Build NFA recursively from regexp tree**

(a(a**|**b))c

aab**|**.c.



**Post-order traversal of regexp tree**

stack

| n1 |
|----|

n1= nfa(a)

n2= nfa(a)
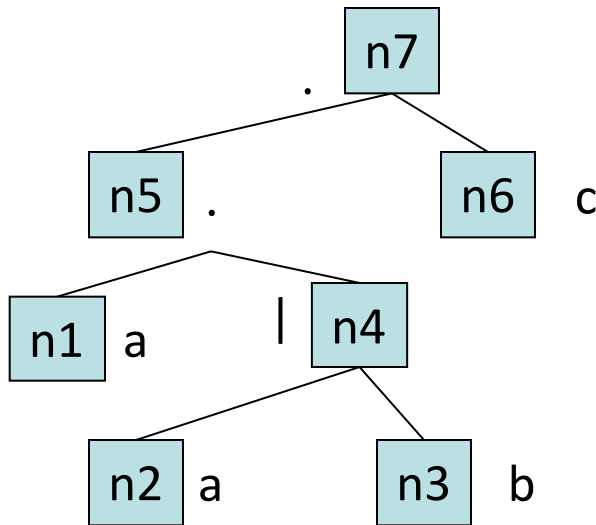
n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

# Thompson's construction

## Converts regexps to NFA

| Build NFA recursively from regexp tree | (a(a❙b))c | Post-order traversal of regexp tree |
|---|---|---|

aab❙.c.

↑



stack

push n2 | n2, n1

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, ❙ )

n5= nfa(n1, n4, . )

n6= nfa(c)
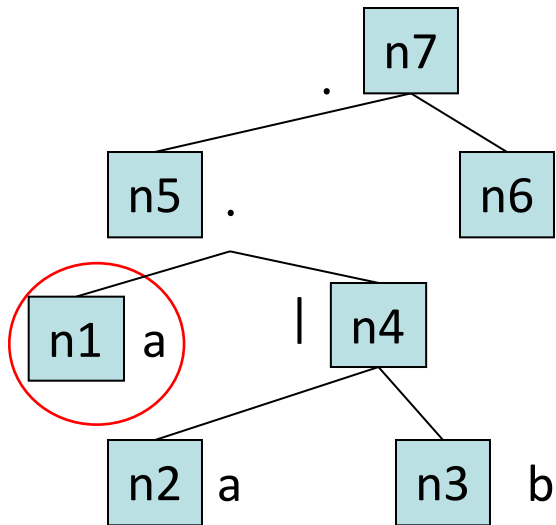
n7= nfa(n5, n6, . )

26

# Thompson's construction
## Converts regexps to NFA

Build NFA recursively from regexp tree

(a(a**|**b))c

aab**|**.c.

Post-order traversal of regexp tree



stack
n2, n1

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

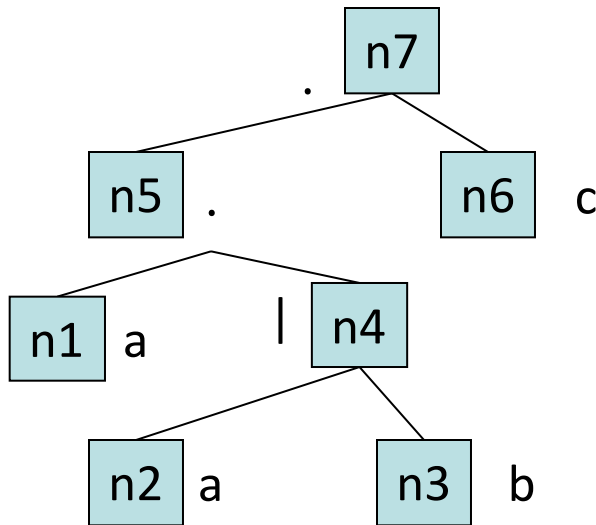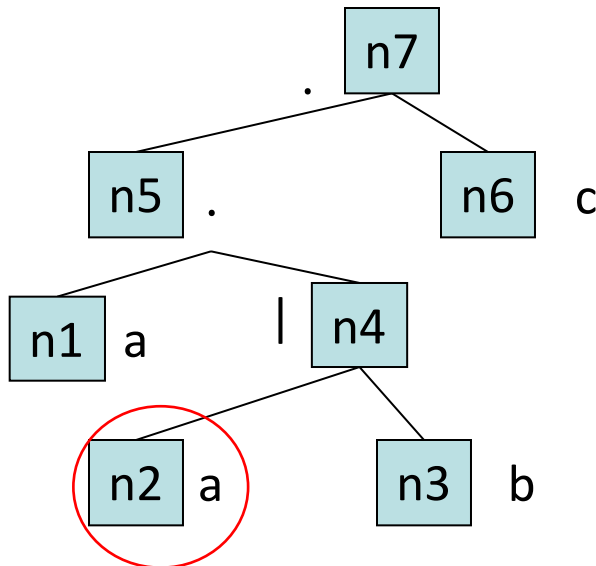n6= nfa(c)

n7= nfa(n5, n6, **.** )

27

# Thompson's construction

## Converts regexps to NFA

Build NFA recursively from regexp tree

$(a(a|b))c$

$aab|.c.$
↑

stack

push n3 | n3, n2, n1 |

Post-order traversal of regexp tree

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, | )

n5= nfa(n1, n4, . )

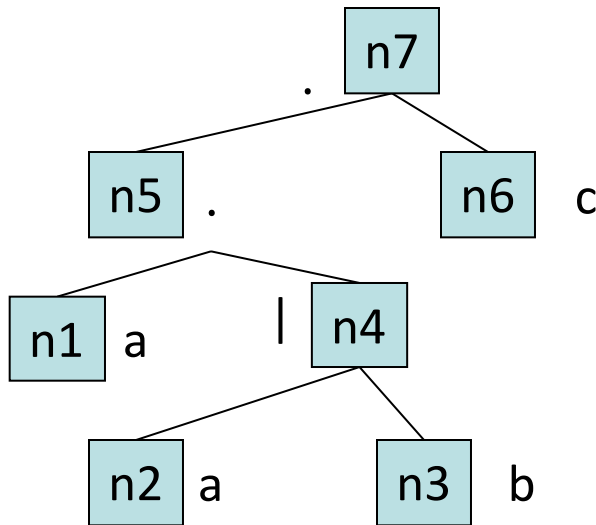n6= nfa(c)

n7= nfa(n5, n6, . )



28

# Thompson's construction
## Converts regexps to NFA

Build NFA recursively from regexp tree

$(a(a\mathbf{|}b))c$

Post-order traversal of regexp tree

$aab\mathbf{|}.c.$



stack

n3, n2, n1

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

# Thompson's construction

## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

(a(a**|**b))c

aab**|**.c.

↑

stack

pop n3,n2 | n1 |

| Post-order traversal of regexp tree |
|---|

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

```
        n7
    .   /  \
  n5 .      n6   c
   /   \
 n1  a   l  n4
        /    \
      n2  a    n3   b
```

30

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

(a(a**|**b))c

aab**|**.c.



↑

push n4

stack

| n4, n1 |
|---|

| Post-order traversal of regexp tree |
|---|

n1= nfa(a)

n2= nfa(a)

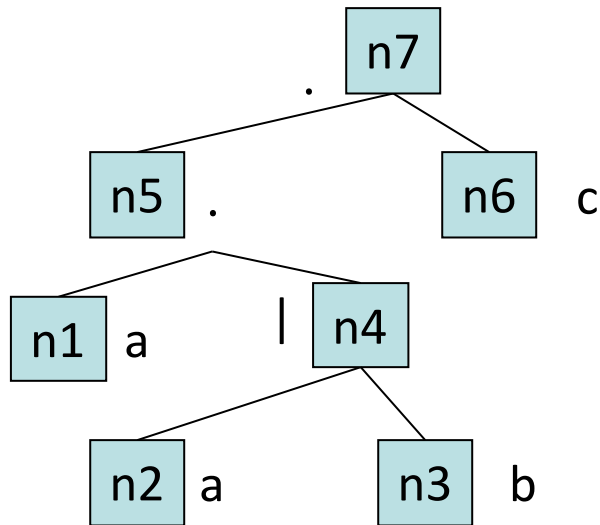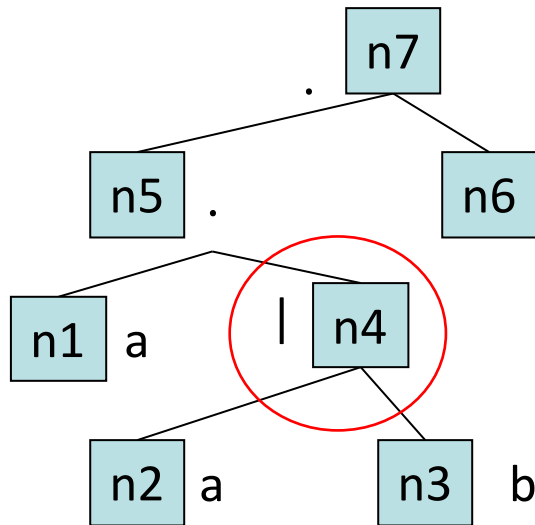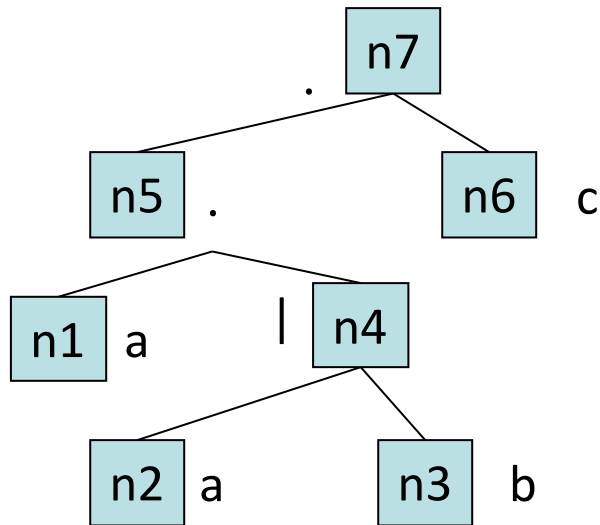n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )
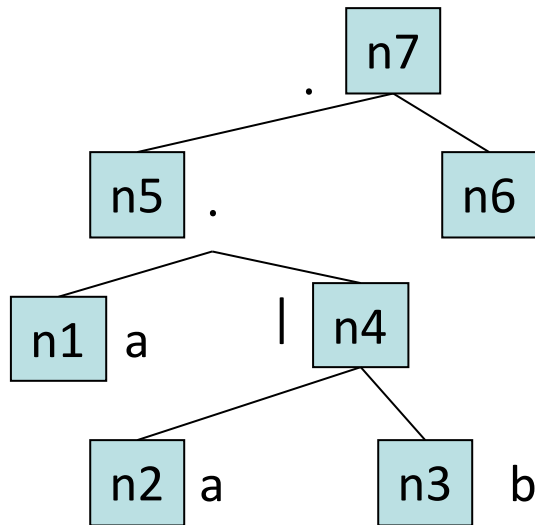
n6= nfa(c)

n7= nfa(n5, n6, **.** )

# Thompson's construction
## Converts regexps to NFA

Build NFA recursively
from regexp tree

(a(a**|**b))c

Post-order traversal of
regexp tree

aab**|**.c.

↑

stack

| n4, n1 |
|---|

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

n7

.

n5  .

n6  c

n1  a

l  n4

n2  a

n3  b

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

$$(a(a\mathbf{|}b))c$$

$$aab\mathbf{|}.c.$$

↑

stack

pop n4,n1

| Post-order traversal of regexp tree |
|---|

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

n7

.

n5  .

n6  c

n1  a

|

n4

n2  a

n3  b

33

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree | $(a(a|b))c$ | Post-order traversal of regexp tree |
|---|---|---|

$aab|.c.$

↑



stack

| n5 |
|----|

push n5

# Thompson's construction

## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

(a(a**|**b))c

aab**|**.c.

↑

stack

| n5 |
|---|

| Post-order traversal of regexp tree |
|---|

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

# Thompson's construction
## Converts regexps to NFA

Build NFA recursively from regexp tree

Post-order traversal of regexp tree

$(a(a|b))c$

$aab|.c.$

↑

push n6

stack

n6, n5



n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, | )

n5= nfa(n1, n4, . )
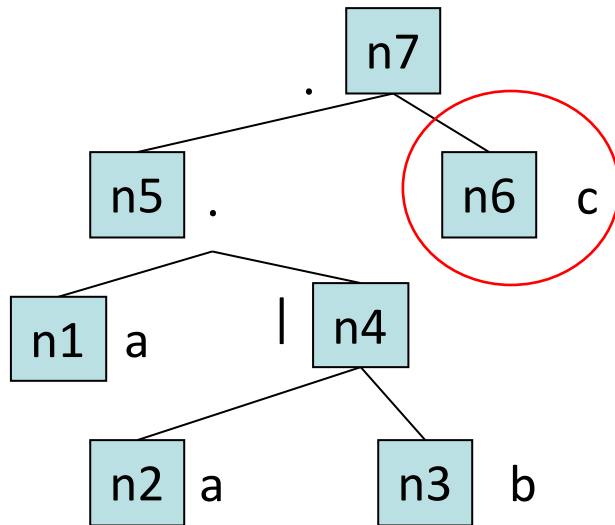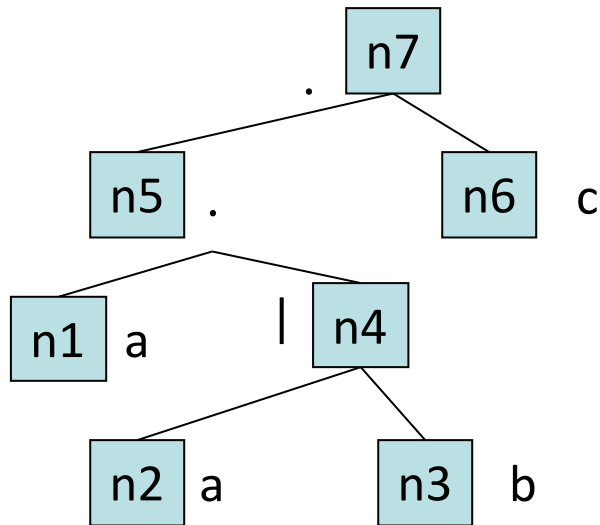
n6= nfa(c)

n7= nfa(n5, n6, . )

# Thompson's construction
## Converts regexps to NFA

Build NFA recursively from regexp tree

$(a(a|b))c$

aab|.c.

↑
stack

| n6, n5 |
|--------|

Post-order traversal of regexp tree

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

n7

.

n5 .

n6 c

n1 a

|

n4

n2 a

n3 b

# Thompson's construction
## Converts regexps to NFA

| Build NFA recursively from regexp tree |
|---|

(a(a|b))c

aab|.c.

↑
stack
pop n6, n5

| |
|---|

.

n7

n5 .

n6   c

n1  a

|   n4

n2  a

n3   b

| Post-order traversal of regexp tree |
|---|

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, | )

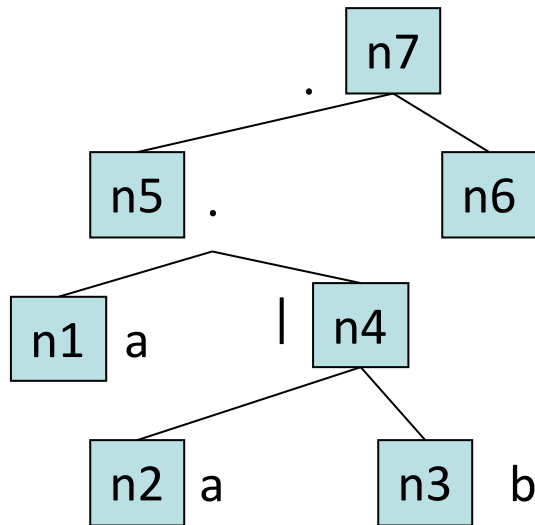n5= nfa(n1, n4, . )

n6= nfa(c)

n7= nfa(n5, n6, . )

# Thompson's construction
## Converts regexps to NFA

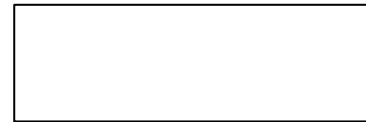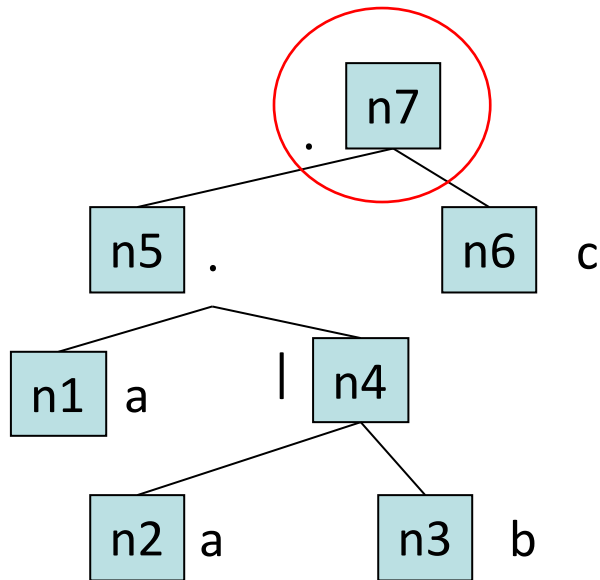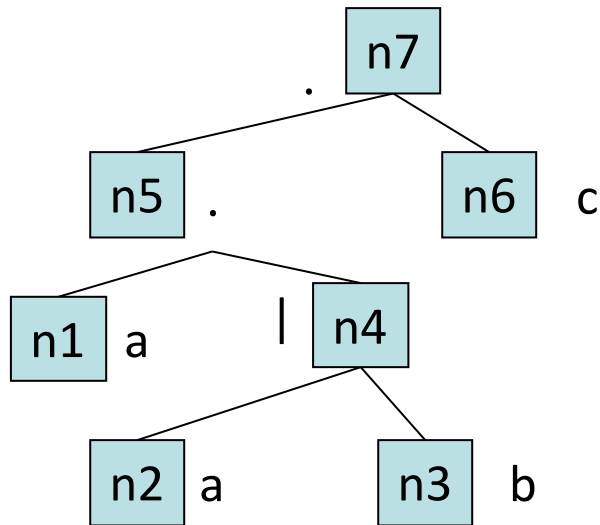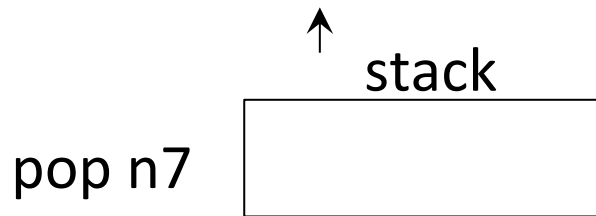Build NFA recursively from regexp tree

$(a(a|b))c$

aab|.c.

Post-order traversal of regexp tree



n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, | )

n5= nfa(n1, n4, . )

n6= nfa(c)

n7= nfa(n5, n6, . )

push n7   stack  n7

# Thompson's construction

## Converts regexps to NFA

Build NFA recursively from regexp tree

$(a(a|b))c$

Post-order traversal of regexp tree

$aab|.c.$

↑

pop n7

stack

n1= nfa(a)

n2= nfa(a)

n3= nfa(b)

n4= nfa(n2, n3, **|** )

n5= nfa(n1, n4, **.** )

n6= nfa(c)

n7= nfa(n5, n6, **.** )

n7 .

n5 . n6 c

n1 a | n4

n2 a n3 b