

Regular Expressions and Automata

Introduction

Regular expressions: convenient to specify a text pattern

Finite automata: efficient to match a text pattern

Both are equivalent

Outline

- Regular Expressions to NFA
 - Illustration
 - Algorithm
- NFA to Regular Expressions

Regular Expressions to NFA

Regex and NFA:

- To say they are equivalent, the **language** produced by the regex $L(R)$ must be equal to the **language** produced by the NFA $L(N)$.

Given a Regex, we need to convert the following parts into NFA

- Basis Regular Expressions:
 - Empty set
 - Empty string
 - Input symbol
- Regular Operations:
 - Concatenation
 - Union
 - Kleene-Closure

Basis Regular Expressions:

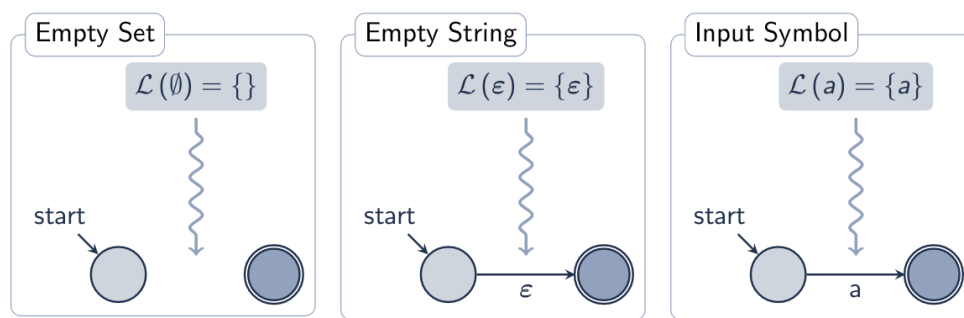


Figure 1: Basis Regular Expressions -> NFA

Regular Operations:

Before handling operations, convert the basis regular expressions into a NFA using the conversions shown in figure 1.

Example:

Given REGEX: $e_1 e_2 e_3^*$, where e_i are expressions

Convert each e_i to NFAs so that they are N_1, N_2, N_3

Then handle the concatenation of N_1 and N_2 and the Kleene Closure of N_3

Concatenation

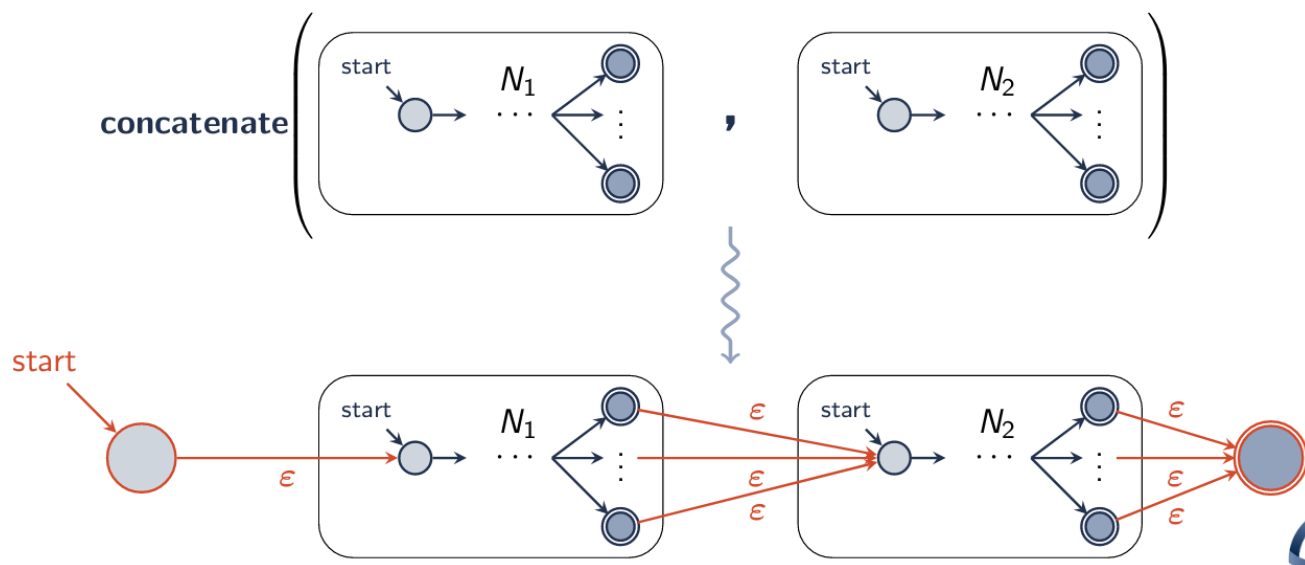


Figure 2: Example Conatention: N_1 and N_2 are basis expressions that were converted to NFAs

Union

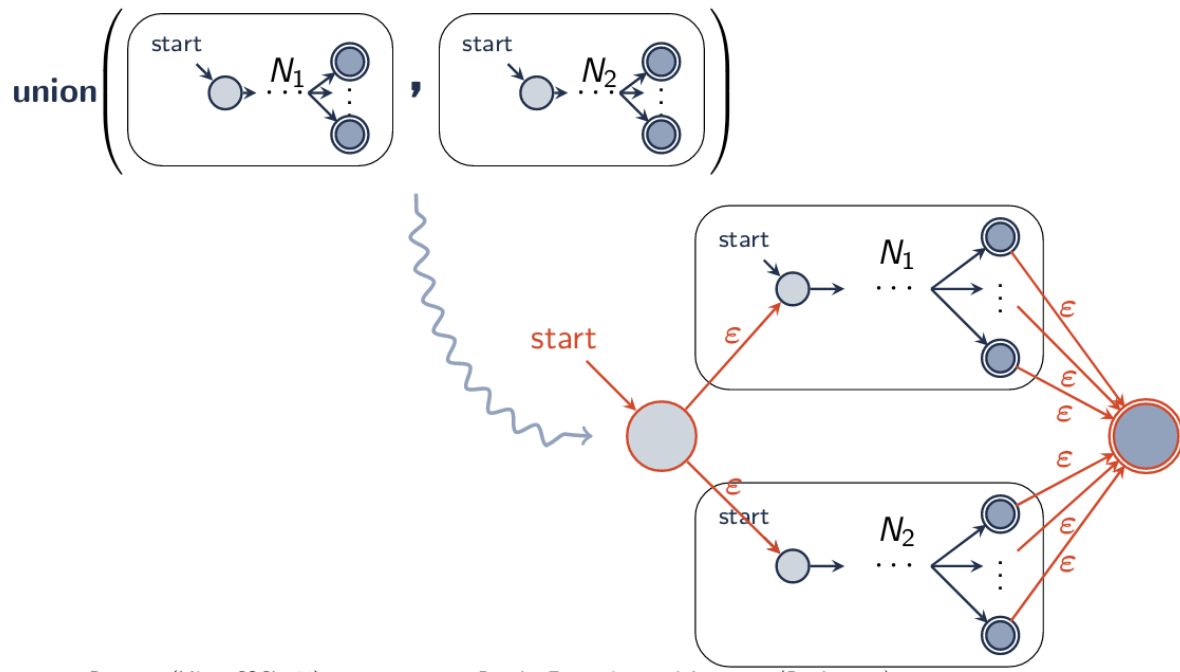


Figure 3: Example Union

Kleene

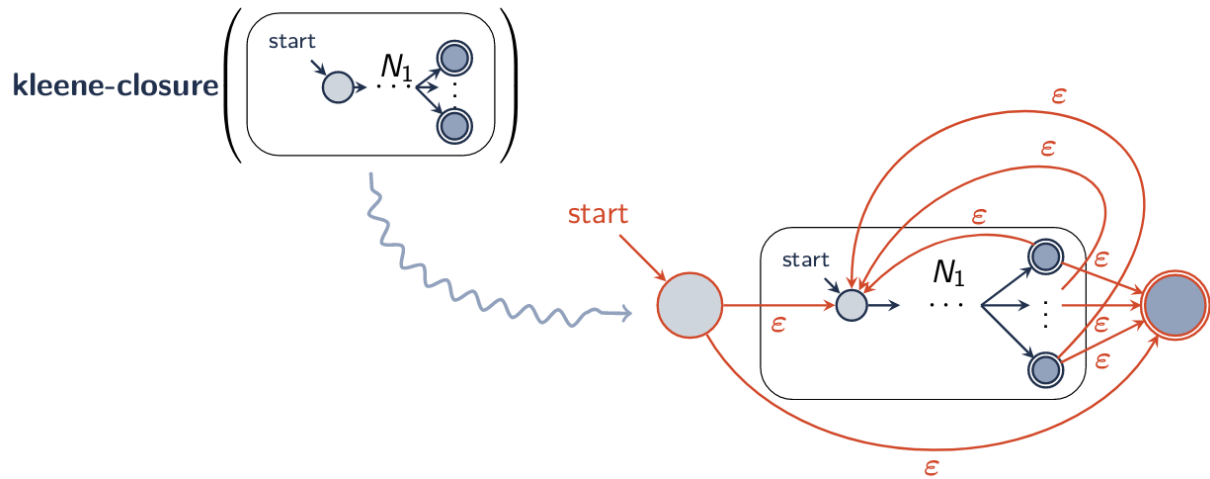
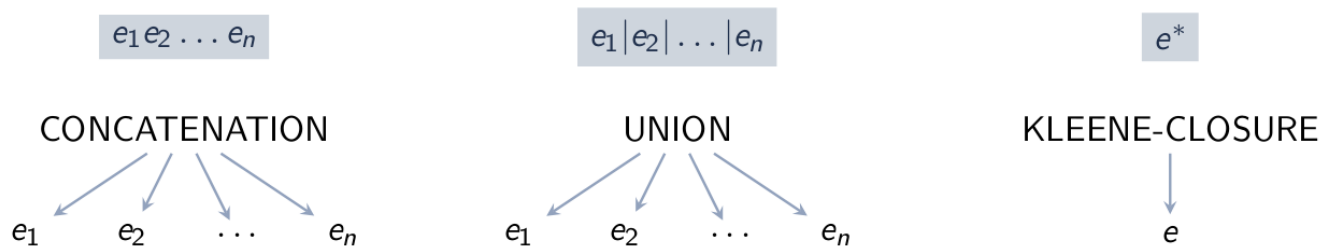


Figure 4: Example Kleene

Abstract Syntax of Regular Expressions



Concatenation \rightarrow (:concatenation e-1 ... e-n)

Union \rightarrow (:union e-1 ... e-n)

Kleene \rightarrow (:kleene e)

McNaughton-Yamada-Thompson Algorithm

Input: A regular expression R

Output: The equivalent NFA N, where $L(R) = L(N)$

Overview: Recursive construction:

Base case: If R is a basis element (\emptyset , ϵ , a) construct the equivalent NFA directly

Recursive case: Otherwise, recurse on the children of R and combine the results according to the operation of R

Algorithm 1: MYT

```
Input:  R ;                                // Regex
Output: (Q, E, s, a) ;                    // states, edges, start, accept
1 if R is  $\emptyset$ ,  $\epsilon$ , or symbol then // Base Case
2   | return MYT-base(R);
3 else if car(R) = KLEENE-CLOSURE then
4   | return kleene-closure(MYT(child(R)));
5 else if car(R) = CONCATENATION then
6   | return MYT-concatenate(cdr(R));
7 else if car(R) = UNION then
8   | return MYT-union(cdr(R));
9 else // Malformed Regex
10  | ERROR;
```

Procedure MYT-base

Input: R

Output: (Q, E, s, a)

```
1  $s \leftarrow \text{newstate}();$ 
2  $a \leftarrow \text{newstate}();$ 
3  $Q \leftarrow \{s, a\};$ 
4 if  $R = \emptyset$  then // empty set  $\emptyset$ 
5   |  $E \leftarrow \emptyset;$ 
6 else // empty string  $\varepsilon$  or symbol  $\sigma$ 
7   |  $E \leftarrow \{s \xrightarrow{R} a\};$ 
```

Procedure MYT-concatenate

Input: c

Output: (Q, E, s, a)

```
1 function  $f(M, R)$  is
2   |  $\text{concatenate}(M, \text{MYT}(R));$ 
3 if  $\emptyset = c$  then
4   | return  $\text{MYT-base}(\varepsilon);$ 
5 else
6   | return  $\text{fold-left}(f, \text{MYT}(\text{car}(c)), \text{cdr}(c));$ 
```

Procedure MYT-union

Input: c

Output: (Q, E, s, a)

```
1 function  $f(M, R)$  is
2   |  $\text{union}(M, \text{MYT}(R));$ 
3 if  $\emptyset = c$  then
4   | return  $\text{MYT-emptyset}(\emptyset);$ 
5 else
6   | return  $\text{fold-left}(f, \text{MYT}(\text{car}(c)), \text{cdr}(c));$ 
```

NFA to Regular Expressions

Input: NFA N

Output: Equivalent regular expression R ,

$$L(N) = L(R)$$

Approach: Construction generalized NFA: an NFA with regular expressions on its edges 1. Convert N to an initial GNFA 2. Iteratively remove (rip) states from the GNFA 3. When the GNFA has only two states (start and accept), the edge between them is the equivalent regular expression

Generalized NFA (GNFA)

A generalized NFA $\tilde{N} = (Q, \Sigma, \delta, q_{\text{start}}, q_{\text{accept}})$:

- Q is the finite set of states
- Σ is the input alphabet
- $\delta: (Q \setminus \{q_{\text{accept}}\}) \times (Q \setminus \{q_{\text{start}}\}) \rightarrow \text{REGEX}$
- $q_{\text{start}} \in Q$ is the start state
- $q_{\text{accept}} \in Q$ is the accept state

Ripping a state from the GNFA

Ripping is taking a section of the GNFA and converting it into a regex

- Given: $\text{state}_i \rightarrow \text{state}_{\text{rip}} \rightarrow \text{state}_j$
- Convert: $\text{state}_i \rightarrow \text{state}_j$ where the transition is the regex

Ripping state \tilde{q}

- Predecessor state q_i
- Successor state q_j

Four types of edges:

- Predecessor $E(q_i, \tilde{q})$
- Successor $E(\tilde{q}, q_j)$
- Loop $E(\tilde{q}, \tilde{q})$
- Bypass $E(q_i, q_j)$

All predecessor/successor pairs

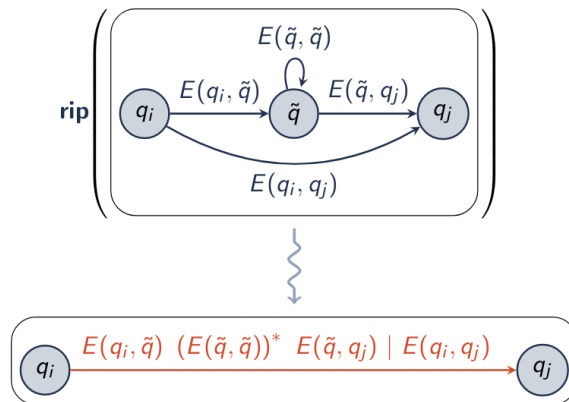


Figure 5: Example rip of a GNFA

NFA to Regular Expression Algorithm

Algorithm 2: NFA to Regex

Input: $N = (Q, \Sigma, E, q_0, F)$; // states, alphabet, edges, start, accept

Output: R ; // Regular Expression

/* Construct the initial GNFA */

1 $N' \leftarrow \text{NFA-to-GNFA}(N)$;

/* Call Convert() subroutine on the GNFA */

2 $R \leftarrow \text{Convert}(N')$;

Algorithm 3: NFA-to-GNFA

Input: $N = (Q, \Sigma, E, q_0, F)$; // states, alphabet, edges, start, accept

Output: $N' = (Q', \Sigma, E', q_{\text{start}}, q_{\text{accept}})$; // states, alphabet, edges, start, accept

1 $Q' \leftarrow Q \cup \{q_{\text{start}}, q_{\text{accept}}\}$; // add new start, accept states

2 $E' \leftarrow \underbrace{\{q_{\text{start}} \xrightarrow{\epsilon} q_0\}}_{\text{edge from new start}} \cup \underbrace{\bigcup_{q \in F} \{q \xrightarrow{\epsilon} q_{\text{accept}}\}}_{\text{edges to new accept}};$

3 **forall** $q_i \in Q$ **do** // Merge multiple edges between nodes into union edges

4 **forall** $q_j \in Q$ **do**

5 $e \leftarrow \{a \xrightarrow{\sigma} b \in E \mid a = q_i \wedge b = q_j\}$; // set of edges from q_i to q_j

6 **if** $|e| \leq 1$ **then**

7 $E' \leftarrow E' \cup e$;

8 **else if** $|e| > 1$ **then**

9 $\ell \leftarrow \left(\bigcup_{a \xrightarrow{\sigma} b \in e} \{\sigma\} \right)$; // set of edge labels from q_i to q_j

10 $r \leftarrow \text{regex}(\ell_0 \cup \dots \cup \ell_n)$;

11 $E' \leftarrow E' \cup \{q_i \xrightarrow{r} q_j\}$;

Function Convert(Q, E)

```
1 if  $|Q| = 2$  then
2    $R \leftarrow E(q_{\text{start}}, q_{\text{accept}})$ ; // Extract label of edge from GNFA start to accept
3   return  $R$ ;
4 else
5    $\tilde{q} \leftarrow$  any state in  $(Q \setminus \{q_{\text{start}}, q_{\text{accept}}\})$ ;
6    $Q' \leftarrow Q \setminus \{\tilde{q}\}$ ;
7    $E' \leftarrow E \setminus \{\tilde{q} \rightarrow \tilde{q}\}$ ;
8   forall  $q_i$  where  $E(q_i, \tilde{q}) \neq \emptyset$  do // predecessors of  $\tilde{q}$ 
9     forall  $q_j$  where  $E(\tilde{q}, q_j) \neq \emptyset$  do // successors of  $\tilde{q}$ 
10       $r \leftarrow \text{regex}(E(q_i, \tilde{q}) (E(\tilde{q}, \tilde{q}))^* E(\tilde{q}, q_j) \cup E(q_i, q_j))$ ;
11       $E' \leftarrow E' \setminus \{(q_i \rightarrow \tilde{q}), (\tilde{q} \rightarrow q_j), (q_i \rightarrow q_j)\}$ ;
12       $E' \leftarrow E' \cup \{q_i \xrightarrow{r} q_j\}$ ;
13 return Convert( $Q', E'$ );
```

Equivalence

