# Regular Expressions

COMP 455 – 002, Spring 2019

# **Regular Expressions**

- Regular expressions are simply *algebraic notation* for defining languages.
- A regular expression defines a language.
- In practice, regular expressions can usually define languages in a concise "user-friendly" manner.
  - At least compared to describing a finite automaton...

# Definition of Regular Expressions

- Regular expressions are defined by taking the *union*, *concatenation*, and *closure* of languages.
  - $\bullet \text{Union: } L_1 \cup L_2 \equiv \{x \mid x \in L_1 \text{ or } x \in L_2\}$
  - **Concatenation**:  $L_1L_2 \equiv \{xy \mid x \in L_1 \text{ and } y \in L_2\}$
  - ♦ Closure ("Kleene closure"):  $L^* \equiv \bigcup_{i \ge 0} L^i$ , where:
    - $\Box L^0 = \{\varepsilon\}$
    - $\Box L^1 = L$
    - $\Box L^i = \text{the concatenation of } i \text{ copies of } L$

# Kleene Closure: Example

- Let language  $L = \{00, 11\}$ .
- (From the previous definition) L<sup>i</sup> represents i strings from L concatenated together.
  - $\bigstar L^2 = \{0000, 0011, 1100, 1111\}$
  - $\bigstar L^3 = \{000000, 000011, 001100, \dots\}$
- L\* is any number of strings from L concatenated together.
  - $\bigstar L^* = \{\varepsilon, 00, 11, 0000, 0011, 1111, 000000, \dots\}$

# Definition of Regular Expressions

Regular expressions are formally recursively defined (for an alphabet  $\Sigma$ ):

- 1. Ø is a regular expression denoting the empty set.
- 2.  $\varepsilon$  is a regular expression denoting { $\varepsilon$ }.
- 3. For each  $a \in \Sigma$ , a is a regular expression denoting  $\{a\}$ .
- 4. If *E* and *F* are regular expressions denoting L(E) and L(F) (respectively), then:
  - (E + F) is a regular expression denoting  $L(E) \cup L(F)$
  - (*EF*) is a regular expression denoting L(E)L(F)
  - ( $E^*$ ) is a regular expression denoting  $L(E)^*$ .

Note that we will use **bold** characters for literal symbols in regular expressions.

# Example Regular Expressions

(The following examples assume  $\Sigma = \{0, 1\}$ .)

**▶ 01** 

The language consisting only of the string 01.

► 0 + 1

The language consisting of the strings {0, 1}.

▶ 01 + 10

The language consisting of the strings {01, 10}.

# Example Regular Expressions

(The following examples assume  $\Sigma = \{0, 1\}$ .)

► **10**(**0** + **1**)

The language consisting of {100, 101}.

▶ 10(0 + 1)\*

The language consisting of all strings that start with 10 and are followed by any number of 0s or 1s.

▶  $10(0+1)^* + \varepsilon$ 

The same language as above, but also containing the empty string.

Jim Anderson (modified by Nathan Otterness)

#### Finite Automata and Regular Expressions

Regular expressions and finite automata define the same class of languages.



# From DFAs to Regular Expressions

**Theorem 3.4**: If L = L(M) for some DFA *M*, then a regular expression *R* exists for which L = L(R).

**Proof**:

• Let  $M = \{\{q_1, q_2, ..., q_n\}, \Sigma, \delta, q_1, F\}$ 

We can re-label any set of states to start with  $q_1$ .

Let  $R_{ij}^k \equiv$  all strings that take us from  $q_i$  to  $q_j$  without entering a state higher than  $q_k$ .

Note:  $R_{ij}^n = all$  strings that go from state  $q_i$  to  $q_j$ .

No state "higher" than k

Let R<sup>k</sup><sub>ij</sub> ≡ all strings that take us from q<sub>i</sub> to q<sub>j</sub> without entering a state higher than q<sub>k</sub>.
 Define R<sup>k</sup><sub>ij</sub> as follows:

$$R_{ij}^{k} = R_{ik}^{k-1} (R_{kk}^{k-1})^{*} R_{kj}^{k-1} \cup R_{ij}^{k-1}$$

$$R_{ij}^{0} = \{a \mid \delta(q_{i}, a) = q_{j}\}, \text{ if } i \neq j$$

$$R_{ii}^{0} = \{a \mid \delta(q_{i}, a) = q_{i}\} \cup \{\varepsilon\}$$



Let  $R_{ij}^k \equiv$  all strings that take us from  $q_i$  to  $q_j$  without entering a state higher than  $q_k$ .

#### Define $R_{ij}^k$ as follows:

Strings that go directly from  $q_i$  to  $q_j$ .  $R_{ij}^k = R_{ik}^{k-1} \left( R_{kk}^{k-1} \right)^* R_{kj}^{k-1} \cup R_{ij}^{k-1}$   $R_{ij}^0 = \left\{ a \mid \delta(q_i, a) = q_j \right\}, \text{ if } i \neq j$   $R_{ii}^0 = \left\{ a \mid \delta(q_i, a) = q_i \right\} \cup \{\varepsilon\}$ 

Let  $R_{ij}^k \equiv$  all strings that take us from  $q_i$  to  $q_j$  without entering a state higher than  $q_k$ .

#### Define $R_{ij}^k$ as follows:

Strings that loop from  $q_i$  back to  $q_i$ .

$$R_{ij}^{k} = R_{ik}^{k-1} (R_{kk}^{k-1})^{*} R_{kj}^{k-1} \cup R_{ij}^{k-1}$$

$$R_{ij}^{0} = \{a \mid \delta(q_{i}, a) = q_{j}\}, \text{ if } i \neq j$$

$$R_{ii}^{0} = \{a \mid \delta(q_{i}, a) = q_{i}\} \cup \{\varepsilon\}$$

**Claim**: There exists a regular expression  $r_{ij}^k$  denoting  $R_{ij}^k$ . We will prove this by induction on k. **Basis**: k = 0

Let  $A \equiv \{a_1, a_2, ..., a_p\}$  denote the symbols that take  $q_i$  to  $q_j$ . Then,

$$r_{ij}^{0} = \begin{cases} a_{1} + a_{2} + \dots + a_{p} & , \text{ if } A \neq \emptyset \text{ and } i \neq j \\ 0 & , \text{ if } A = \emptyset \text{ and } i \neq j \\ a_{1} + a_{2} + \dots + a_{p} + \varepsilon & , \text{ if } A \neq \emptyset \text{ and } i = j \\ \varepsilon & , \text{ if } A = \emptyset \text{ and } i = j \end{cases}$$

**Claim**: There exists a regular expression  $r_{ij}^k$  denoting  $R_{ij}^k$ . We will prove this by induction on k. **Inductive step**: k > 0.

In this case,

$$r_{ij}^{k} = r_{ik}^{k-1} (r_{kk}^{k-1})^{*} r_{kj}^{k-1} + r_{ij}^{k-1}$$

**Claim**: There exists a regular expression  $r_{ij}^k$  denoting  $R_{ij}^k$ . We will prove this by induction on k. **Inductive step**: k > 0.

In this case,



The inductive hypothesis lets us assume that we already have these regular expressions!

#### Finishing up:

We note that  $L(M) = \bigcup_{q_i \in F} R_{1j}^n$ 

We still need to show that we can construct a regular expression that matches L(M).

Any string that takes us from state  $q_1$  to  $q_j$ , where  $q_j$  is a final state.

Therefore, L(M) is denoted by the regular expression:

$$r_{1j_1}^n + r_{1j_2}^n + \dots + r_{1j_{p'}}^n$$
 where  $F = \{q_{j_1}, q_{j_2}, \dots, q_{j_p}\}$ .

	k = 0	k = 1	<i>k</i> = 2	<i>k</i> = 3
$r_{11}^k$	Е			
$r_{12}^{k}$	0			
$r_{13}^{k}$	1			
$r_{21}^k$	Ø			
$r^k_{22}$	$0 + 1 + \varepsilon$			
$r^k_{23}$	Ø			
$r^k_{31}$	Ø			
$r^k_{32}$	Ø			
$r_{33}^{k}$	$0 + 1 + \varepsilon$			



If *A* is the set of symbols going from  $q_i$  to  $q_j$ :  $r_{ij}^0 = \begin{cases} a_1 + a_2 + \dots + a_p &, \text{ if } A \neq \emptyset \text{ and } i \neq j \\ \emptyset &, \text{ if } A = \emptyset \text{ and } i \neq j \\ a_1 + a_2 + \dots + a_p + \varepsilon &, \text{ if } A \neq \emptyset \text{ and } i = j \\ \varepsilon &, \text{ if } A = \emptyset \text{ and } i = j \end{cases}$ 

	k = 0	<i>k</i> = 1	<i>k</i> = 2	<i>k</i> = 3
$r_{11}^{k}$	Е	Е		
$r_{12}^{k}$	0	0		
$r_{13}^{k}$	1	1		
$r_{21}^{k}$	Ø	Ø		
$r_{22}^{k}$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$		
$r_{23}^{k}$	Ø	Ø		
$r_{31}^{k}$	Ø	Ø		
$r_{32}^{k}$	Ø	Ø		
$r_{33}^k$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$		



$$r_{ij}^{k} = r_{ik}^{k-1} (r_{kk}^{k-1})^{*} r_{kj}^{k-1} + r_{ij}^{k-1}$$

Jim Anderson (modified by Nathan Otterness)

	k = 0	k = 1	<i>k</i> = 2	<i>k</i> = 3
$r_{11}^{k}$	Е	Е	Е	
$r_{12}^{k}$	0	0	$0(0+1)^{*}$	
$r_{13}^{k}$	1	1	1	
$r_{21}^{k}$	Ø	Ø	Ø	
$r_{22}^{k}$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	<b>(0 + 1</b> )*	
$r_{23}^{k}$	Ø	Ø	Ø	
$r_{31}^{k}$	Ø	Ø	Ø	
$r_{32}^{k}$	Ø	Ø	Ø	
$r_{33}^k$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	



$$r_{ij}^{k} = r_{ik}^{k-1} (r_{kk}^{k-1})^{*} r_{kj}^{k-1} + r_{ij}^{k-1}$$

	k = 0	k = 1	<i>k</i> = 2	<i>k</i> = 3
$r_{11}^{k}$	Е	Е	Е	Е
$r_{12}^{k}$	0	0	$0(0+1)^{*}$	$0(0+1)^{*}$
$r_{13}^{k}$	1	1	1	$1(0+1)^{*}$
$r_{21}^{k}$	Ø	Ø	Ø	Ø
$r_{22}^{k}$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	<b>(0 + 1</b> )*	<b>(0 + 1</b> )*
$r_{23}^{k}$	Ø	Ø	Ø	Ø
$r_{31}^{k}$	Ø	Ø	Ø	Ø
$r_{32}^{k}$	Ø	Ø	Ø	Ø
$r_{33}^k$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	$(0 + 1)^*$



$$r_{ij}^{k} = r_{ik}^{k-1} (r_{kk}^{k-1})^{*} r_{kj}^{k-1} + r_{ij}^{k-1}$$

Jim Anderson (modified by Nathan Otterness)

	k = 0	k = 1	<i>k</i> = 2	k = 3
$r_{11}^k$	Е	Е	Е	Е
$r_{12}^{k}$	0	0	$0(0+1)^{*}$	$0(0+1)^{*}$
$r_{13}^{k}$	1	1	1	$1(0+1)^*$
$r_{21}^k$	Ø	Ø	Ø	Ø
$r^k_{22}$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	<b>(0 + 1</b> )*	<b>(0 + 1</b> )*
$r^k_{23}$	Ø	Ø	Ø	Ø
$r^k_{31}$	Ø	Ø	Ø	Ø
$r_{32}^{k}$	Ø	Ø	Ø	Ø
$r_{33}^{k}$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	$0 + 1 + \varepsilon$	<b>(0 + 1</b> )*



The expression taking us from  $q_1$  to  $q_2$  through states up to  $q_3$  is  $r_{12}^3$ :  $\mathbf{0}(\mathbf{0} + \mathbf{1})^*$ 

# Comments on this Approach

- This approach is an example of *dynamic programming*, a COMP 550 topic.
- For an *n*-state DFA, this approach may produce a regular expression with 4<sup>n</sup> symbols!
- Section 3.2.2 of the book presents a more efficient method.

### From Regular Expressions to *ε*-NFAs

**Theorem 3.7** (reworded): If *R* is a regular expression, L(R) = L(E) for some NFA with  $\varepsilon$ -transitions *E*.

- We will construct *E* with one final state and *no transitions out of that state*.
- ▶ We will do this by induction on the number of operators in *R*.

#### Proof of Theorem 3.7

**Base case**: *R* has no operators. There are three possibilities here:



**Inductive step**:

We now need to handle the three regular-expression operators:

- ▶ Union  $(R_1 + R_2)$
- Concatenation  $(R_1R_2)$

Closure  $(R_1^*)$ 

#### **Inductive step**: Union: $R = R_1 + R_2$

We can assume by the inductive hypothesis that we already have  $\varepsilon$ -NFAs  $E_1$  and  $E_2$  accepting the same languages as  $R_1$  and  $R_2$ . Construct the  $\varepsilon$ -NFA E to accept the same language as R like this:



**Inductive step**: Concatenation:  $R = R_1 R_2$ 

Once again, assume  $E_1$  and  $E_2$  are  $\varepsilon$ -NFAs accepting the same languages as  $R_1$  and  $R_2$ . Construct the  $\varepsilon$ -NFA E to accept the same language as R like this:



**Inductive step**: Closure:  $R = R_1^*$ 

Assume  $E_1$  is an  $\varepsilon$ -NFA accepting the same language as  $R_1$ . Construct the  $\varepsilon$ -NFA E to accept the same language as R like this:  $\varepsilon$ 



We'll convert the regular expression  $0(0 + 1)^*$  to an NFA with  $\varepsilon$ -transitions.



$$0+1: \quad \text{Start} \rightarrow \underbrace{e}_{\epsilon} \underbrace{a}_{c} \underbrace{d}_{c} \underbrace{b}_{\epsilon} \underbrace{e}_{\epsilon} \underbrace{f}_{c} \underbrace{f}_{d} \underbrace{f}_{\epsilon} \underbrace{f}_{c} \underbrace{f}_{c} \underbrace{f}_{d} \underbrace{f}_{c} \underbrace$$



Jim Anderson (modified by Nathan Otterness)







Note: A *much* simpler machine exists.

Jim Anderson (modified by Nathan Otterness)