

Lect
07

CST8152 Compilers

Algonquin College

Computer Engineering
Technology

CST8152 Compilers

Fall, 2023



Lect
07

Based on resources developed
by prof. **Svillen Ranev**.

Prof. Paulo Sousa

Algonquin College

Computer Engineering
Technology

CST8152 Compilers

Fall, 2023



**Lect
07**

Formal Representations

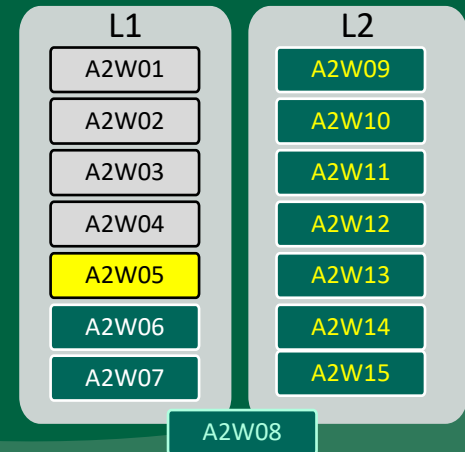


Based on resources developed
by prof. **Svillen Ranev**.

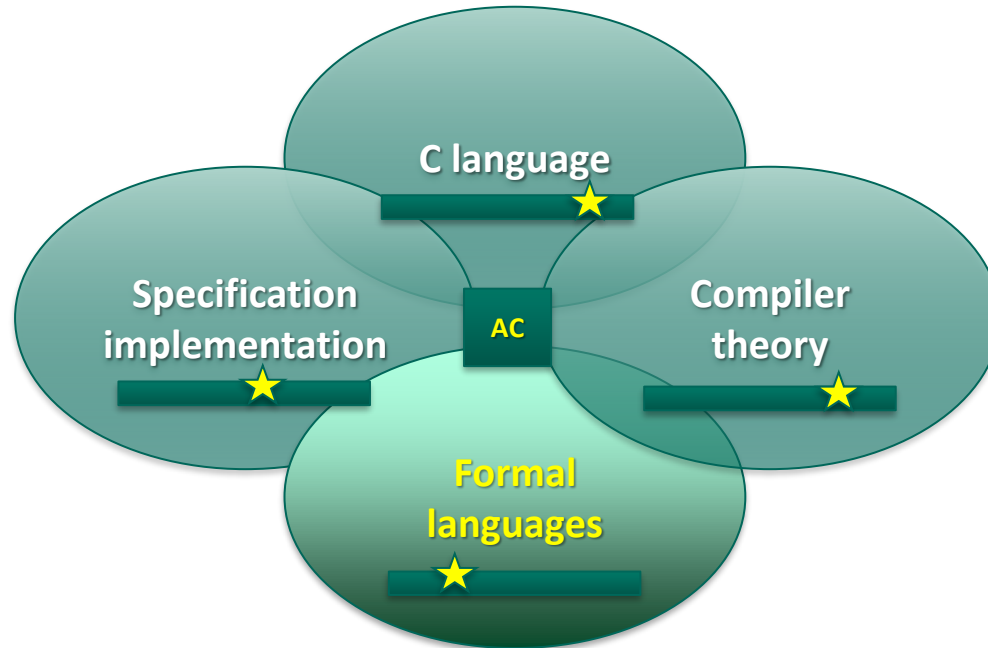
Prof. Paulo Sousa

Art 7: Formal Representations

- *General View*
- *Regular Expressions*
- *Grammar*
- *Automata*



Let's start...





Compilers – Art. 7

Formal Representations

Universal Concepts

Alphabet:

- An **alphabet** Σ is a finite, non-empty, set of symbols. For example:
- The **binary** alphabet is {0, 1}
- The **decimal** alphabet is {0,1,2,3,4,5,6,7,8,9}
- **Note:** The metasympols { , and } used here that are **not** in the alphabet.

IMPLEMENTATION NOTE

- For the scanner, the alphabet may be characters in the **ASCII** character set.
- For the parser the alphabet is the set of tokens produced by the scanner.
- Ex. Important sets: *keywords*

MOLD Language keywords: { "data",
"code", "int", "real", "string", "if", "then",
"else", "while", "do" }

Universal Concepts

String:

- A string is a **finite set** of symbols from an alphabet (not necessarily in a grammar).
- **Example:**
 - For the alphabet $\Sigma = \{a, b, c\}$ some strings are: $\{a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, abc, acb, bac, bca, cab, cba\}$
- **Note:** The order of symbols in a string **matter**.

Empty string:

- ϵ is the **empty** string.
- It is the string consisting of no symbols.
- The length of a string s is $|s|$ and is equal to the number of symbols in the string.
- So: $|\epsilon| = 0, |abc| = 3$.



Universal Concepts

Equipotency of ϵ :

- ϵ can be considered as a neutral element in operations:

$$\epsilon = \epsilon\epsilon = \epsilon^k$$

$$\epsilon X = X\epsilon = X$$

Think about this [2]:

- What is the best model for PL?



Empty Language:

Useless

- The regular expression that matches nothing: Φ .
 - Defined by Φ , it is the pattern for nothing; it generates the set containing nothing;
 - Representation: $L(\Phi) = \{\}$
- Note:** This is not the same as the empty string.
 - By contrast, ϵ is the pattern for the set that contains the **string** that contains no characters.

Operations in Languages

Concatenation of sets

- The concatenation of two sets A and B is defined by:

$$AB = \{ xy \mid x \text{ in } A \text{ and } y \text{ in } B \}$$

which reads “*the set of strings xy such that x is in A and y is in B*”.

- For example.
- If $A = \{a,b\}$ and $B = \{c,d\}$ then $AB = \{ac, ad, bc, bd\}$

Powers of sets

- The power of a set A: The repetition of A several times.

$$A^4 = \{ x \mid \text{4-symbol string} \}$$

which reads “*the set of strings with four symbols*”.

- This is just repeated:
 $A^0 = \{ \varepsilon \}, A^1 = A, A^2 = AA, A^3 = AAA, \dots$
- Note that $A^0 = \{ \varepsilon \}$ (for any set)

Operations in Languages

Union of sets

- The union of two sets A and B is defined by:

$$A \cup B = \{ x \mid x \text{ in } A \text{ or } x \text{ in } B \}$$

which reads “the set of strings x such that x is in A or x is in B ”.

- For example.
- If $A = \{a,b\}$ and $B = \{c,d\}$ then $A \cup B = \{a, b, c, d\}$

Kleene closure

- The Kleene closure of a set A is the *** operator** defined as the set of all strings including the empty string:

$$A^* = \bigcup_{i=0}^{\infty} A^i$$

- It is the union of all powers of A.

$$A^* = A^0 + A^1 + A^2 + A^3 + \dots$$

Operations in Languages

Positive closure

- The positive closure is the **+** operator defined as the set of all strings excluding the empty string:

$$A^+ = \bigcup_{i=1}^{\infty} A^i$$

- It means:

$$A^+ = A^* - \{\epsilon\} = A^1 + A^2 + A^3 + \dots$$

Examples

Let L be the set {A, B, ... Z, a, b, ... z} Let D be the set {0, 1,9}

- $L \cup D$ is the set of letters and digits: {A, B, ... Z, a, b, ... z, 0, 1,9}
- LD is the set of strings consisting of a letter followed by a digit: {A0, A1, A2 ..., B0, B1,... Z9, a0, b0, a1, b1 ... }
- L^4 is the set of all four-letter strings
- L^* is the set of all strings of letters, including ϵ .
- $L(L \cup D)^*$ is the set of all strings of letters and digits beginning with a letter
- D^+ is the set of all strings of one or more digits



Compilers – Art. 7

RE (Regular Expressions)

Remember Kleene Theorem

- **Main Idea:**

The language that can be defined by:

1. **Regular Expressions** (compact language);
2. **Regular Grammar** (syntax production rules);
3. **Finite Automaton** (DFA);
4. **Transition graph** (transition / state diagrams).



Prof. Kleene

Source: Wikipedia

Model 1: Regular Expressions:

1. Regular expressions are a **convenient notation** (or means or tools) for specifying certain simple (though possibly infinite) **set of strings** over some alphabet.
2. A regular expression is a shorthand equivalent to a regular **grammar**.

$$L(RE) = L(G)$$

Remember Kleene Theorem

- **TIP:** A regular expression can be used to construct a **Deterministic Finite Automaton (DFA)** which therefore can recognize strings (words) of the grammar, which is the purpose of the Scanner.
- The sets of strings defined by regular expression are termed **regular sets**.

To define the RE (as any expression notation) use **operands** and **operations**.

- The **operands** are **alphabet symbols** or **strings** defined by regular expressions (regular definitions).
- The standard operations are **catenation** (concatenation) , union or **alternation** (|), and **recursion** or Kleene closure (*)
- Regular expressions use the **metasymbols** |, (,), {, } , [,], * , + (and others ?, ^) to define its operations.

RE Operations:

- **Catenation**

1. RE: For the alphabet $\Sigma = \{a, b, c\}$, if $x = ab$ and $y = b\&c$, then $x.y = abbc$.
2. Note: Since ε is a valid word (empty string), remember that:
 $\varepsilon X = X\varepsilon = X$

Alternation:

1. RE: For example if $\Sigma = \{a, b, c, d\}$ then $L(a|b) = \{a, b\}$ and $L(c|d) = \{c, d\}$ so
 - $a|b$ is either a or b
 - $c|d$ is either c or d

RE Operations:

- **Recursion**

1. RE: For a regular expression **r**, Kleene Closure is defined by:

- $L(r^*) = L(r)^*$
- which means concatenation with all powers of L .

Example:

$$L((a|bb)^*) = L(a|bb)^* = L(a|bb)^0 + L(a|bb)^1 + L(a|bb)^2 + L(a|bb)^3 + \dots$$

- $L(a|bb)^0 = \{\epsilon\}$
- $L(a|bb)^1 = \{a, bb\}$
- $L(a|bb)^2 = \{a, bb\}\{a, bb\} = \{aa, abb, bba, bbbb\}$ or:
- $L((a|bb)^*) = \{\epsilon, a, bb, aa, abb, bba, bbbb, aaa, abba, bbaa, bbbba, aabb, \dots\}$
- $\{dc, ddc, dddc, ddddc, \dots\} = d^+c(a|c|dd)^*$

More about RE

Special operations

- Positive Closure (One or more +)
 - $a^+ = aa^*$ (also $a^* = a^+ | \epsilon$)
- Exponentiation or Power operation (exactly k)
 - $a^k = aaa...a$ (exactly k times)
- Optional Inclusion (Zero or one ?)
 - $a? = a | \epsilon$

$(a|c|dd+)^*? = \epsilon | (a|c|dd+)^* = \epsilon | a | c | dd | aa$

Character classes:

- Specify a **range** of characters or numbers that follow a sequence.
- Examples:
 - $[a-z]$ means any character in the range a to z.
 - $[A-Z]$ means any character in the range A to Z.
- A regular expression for the **pattern** for an identifier that begins with a letter or an underscore and is followed by any number of numbers and letters is:

$[a-zA-Z_][A-Za-z0-9]^*$

$\{ a, b, ..., A, B, ..., _, aA, aa0, ... \}$

More about RE

Complements:

- A **complement character class** is specified using the \wedge or (\sim) symbols, or **Not** operator.
- Examples:
 - `[^a-z]` matches any character except a to z.
 - `Comment = // (^CR)*CR`

TIP

Sometimes, the complement is the better (shortest) to represent collections.

Precedence of operators:

- Decrease precedence order:
 - Grouping: ()
 - Character classes: []
 - Power / Recursion (Kleene star): *
 - Concatenation: , , , ,
 - Alternation: |
 - Positive closure: +
 - Optional: ?
 - Iteration: k

More about RE (1)

- Remember:

Example 3.4: Let $\Sigma = \{a, b\}$.

1. The regular expression **$\mathbf{a|b}$** denotes the language $\{a, b\}$.
2. **$\mathbf{(a|b)(a|b)}$** denotes $\{aa, ab, ba, bb\}$, the language of all strings of length two over the alphabet Σ . Another regular expression for the same language is **$\mathbf{aa|ab|ba|bb}$** .
3. **$\mathbf{a^*}$** denotes the language consisting of all strings of zero or more a 's, that is, $\{\epsilon, a, aa, aaa, \dots\}$.
4. **$\mathbf{(a|b)^*}$** denotes the set of all strings consisting of zero or more instances of a or b , that is, all strings of a 's and b 's: $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$. Another regular expression for the same language is **$\mathbf{(a^*b^*)^*}$** .
5. **$\mathbf{a|a^*b}$** denotes the language $\{a, b, ab, aab, aaab, \dots\}$, that is, the string a and all strings consisting of zero or more a 's and ending in b .



More about RE (2)

- RE Properties:

LAW	DESCRIPTION
$r s = s r$	$ $ 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 $ $
$\epsilon r = r\epsilon = r$	ϵ is the identity for concatenation
$r^* = (r \epsilon)^*$	ϵ is guaranteed in a closure
$r^{**} = r^*$	$*$ is idempotent

Figure 3.7: Algebraic laws for regular expressions

More about RE (3)

- **Common Syntax for RE (Lex):**
- Conventions are used in order to express RE in a specific notation.
- For instance, the Lex file (used for automated generated parsers), has the following notation...

EXPRESSION	MATCHES	EXAMPLE
<i>c</i>	the one non-operator character <i>c</i>	<i>a</i>
<i>\c</i>	character <i>c</i> literally	<i>*</i>
<i>"s"</i>	string <i>s</i> literally	<i>"**"</i>
<i>.</i>	any character but newline	<i>a.*b</i>
<i>^</i>	beginning of a line	<i>^abc</i>
<i>\$</i>	end of a line	<i>abc\$</i>
<i>[s]</i>	any one of the characters in string <i>s</i>	<i>[abc]</i>
<i>[^s]</i>	any one character not in string <i>s</i>	<i>[^abc]</i>
<i>r*</i>	zero or more strings matching <i>r</i>	<i>a*</i>
<i>r+</i>	one or more strings matching <i>r</i>	<i>a+</i>
<i>r?</i>	zero or one <i>r</i>	<i>a?</i>
<i>r{m,n}</i>	between <i>m</i> and <i>n</i> occurrences of <i>r</i>	<i>a{1,5}</i>
<i>r₁r₂</i>	an <i>r₁</i> followed by an <i>r₂</i>	<i>ab</i>
<i>r₁ r₂</i>	an <i>r₁</i> or an <i>r₂</i>	<i>a b</i>
<i>(r)</i>	same as <i>r</i>	<i>(a b)</i>
<i>r₁/r₂</i>	<i>r₁</i> when followed by <i>r₂</i>	<i>abc/123</i>

Figure 3.8: Lex regular expressions

Grammar Examples

$VID = AVID \mid SVID$

- Examples: {a, b, c, a1, b12, c123, ..., a\$, b\$, c\$, a1\$, b12\$, c123\$}

$AVID = L (L \mid D)^*$

- Examples: {a, b, c,abc, abcdf, abc123...}

$L = a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z$

- Examples: {a, b, c, ..., z, A, ..., Z}

$D = 0 \mid \dots \mid 9$

- Examples: {0, 1, 2, 3, ..., 9}

$SVID = AVID\$$

- Examples: {a\$, b\$, c\$,abc\$, abcdf\$, abc123\$...}

Suppose a grammar in which string variables must finish with "\$"

$IL = DIL$

- Examples: {1, 2, ..., 0, 00, 111, 123, ..., 777, ...}

$DIL = ZDS \mid NzDD^*$

- Examples: {1, 2, ..., 0, 00, 111, 123, ..., 777, ...}

$ZDS = 00^*$

- Examples: {0, 00, 000, ..., 000000000}

$NzD = 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

- Examples: {1, 2, ..., 9}

$D = 0 \mid NzD$

- Examples: {0, 1, 2, ..., 9}

Suppose a grammar in which numbers starting with "0" cannot include other digits.



Compilers – Art. 7

BNF – Grammar (review)

Remember Kleene Theorem

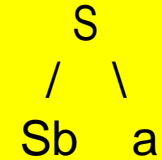
- Model 2: Grammar Formalization**

$$G = \{V_T, V_N, P, S\}$$

- Where:

1. V_T = Terminals;
2. V_N = Non-Terminals ($V_N \rightarrow V_N \mid V_T$)
3. P = Production rules: $\{A \rightarrow X_1 X_2 \dots X_N\}$;
4. S = Start symbol.

Example:



1. $G1: \{\{a,b\}, \{S\}, P = \{S \rightarrow Sb \mid a\}, S\}$

The infinite set of words can be given by:

$a, ab, abb, abbb, \dots$

The corresponding RE1: ab^*

Note: The grammar $G1$ and RE1 are equivalent!

Grammar Examples

$\langle \text{variable identifier} \rangle \rightarrow \langle \text{arithmetic variable identifier} \rangle \mid \langle \text{string variable identifier} \rangle$

- Example: {a, b, c, a1, b12, c123, ..., a\$, b\$, c\$, a1\$, b12\$, c123\$, ...}

$\langle \text{arithmetic variable identifier} \rangle \rightarrow \langle \text{letter} \rangle \langle \text{opt_letters or digits} \rangle$

- Example: {a, b, c, ..., abc, abcdf, abc123...}

$\langle \text{opt_letters or digits} \rangle \rightarrow \langle \text{letters or digits} \rangle \mid \epsilon$

- Example: { ϵ , a, b, c, ..., abc, abcdf, abc123...}

$\langle \text{letters or digits} \rangle \rightarrow \langle \text{letter or digit} \rangle \mid \langle \text{letters or digits} \rangle \langle \text{letter or digit} \rangle$

- Example: {a, b, c, ..., 1, 2, 3, 1s, e2 ...1111, aaaaa,...}

$\langle \text{letter or digit} \rangle \rightarrow \langle \text{letter} \rangle \mid \langle \text{digit} \rangle$

- Example: {a, b, c, ..., 1, 2, 3}

$\langle \text{letter} \rangle \rightarrow a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z$

- Example: {a, b, c, ..., z, A,..., Z}

$\langle \text{digit} \rangle \rightarrow 0 \mid \dots \mid 9$

- Example: {0, 1, 2, 3, ..., 9}

Grammar Examples

`<string variable identifier> -> <arithmetic variable identifier>$`

- Example: {a\$, b\$, c\$, ...abc\$, abcdf\$, abc123\$...}

`<integer literal> -> <decimal integer literal>`

- Example: {1, 2, , 0, 00, 111, 123, ... , 777, ...}

`<decimal integer literal> -> <zeros> | <non zero digit> <opt_digits>`

- Example: {1, 2, , 0, 00, 111, 123, ... , 777, ...}

`<zeros> -> 0 | <zeros>0`

- Example: {0, 00, 000, ... , 00000000}

`<opt_digits> -> <digits> | ϵ`

- Example: { ϵ , 1, 2, , 0, 00, 111, 123, ... , 777, ...}

`<digits> -> <digit> | <digits> <digit>`

- Example: { 1, 2, , 0, 00, 111, 123, ... , 777, ...}

`<digit> -> 0 | <non zero digit>`

- Example: { 0, 1, 2, ... , 9}

`<non zero digit> -> 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9`

- Example: {1, 2, ... , 9}



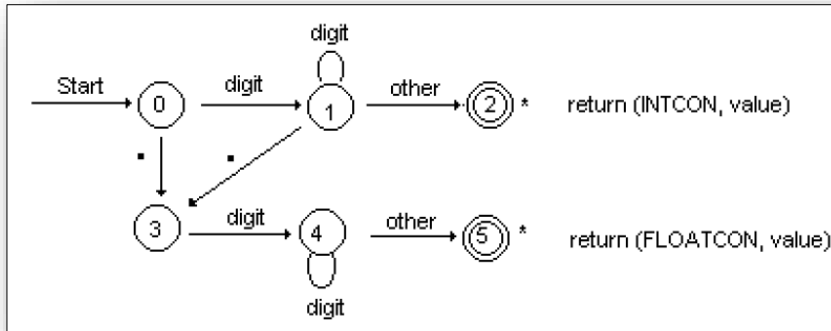
Compilers – Art. 7

Starting Automata

Remember Kleene Theorem

- **Model 3: Automaton:**

1. Functional way to define the evolution of an acceptable string in a language.
2. Can be visual such as:



Finite Automaton:

1. **Mathematical representation** of transitions that transform inputs in outputs that describes a language.

$$L(G) = A(\Sigma, Q, q_0, Q_f, \Delta)$$

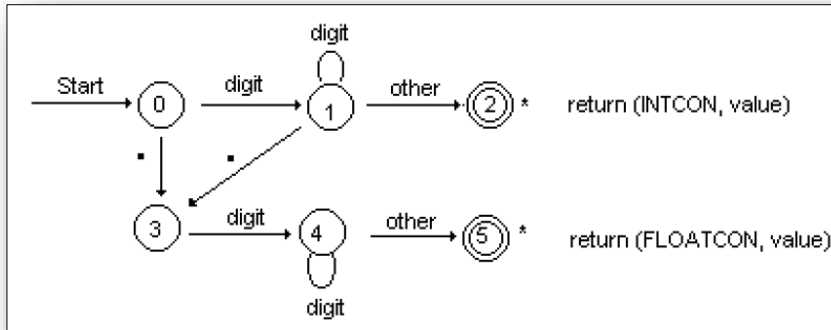
2. This notation includes the **alphabet** (Σ) that, **starting** in a **state** (q_0) can perform words in the **end** (Q_f), by **productions** (Δ) between **states** (Q).

Note: **Empty string** (ϵ) is also acceptable.

Remember Kleene Theorem

- Model 3: Automaton:**

1. Modeling: $L(G) = A(\Sigma, Q, q_0, Q_f, \Delta)$



Delta transitions:

$\Delta = \{$

$\delta(q_0, \text{digit}) = q_1;$

$\delta(q_0, \text{period}) = q_3;$

$\delta(q_1, \text{digit}) = \delta(\delta(q_0, \text{digit}), \text{digit}) = q_1;$

$\delta(q_1, \text{period}) = \delta(\delta(q_0, \text{digit}), \text{period}) = q_3;$

//...

$\}$

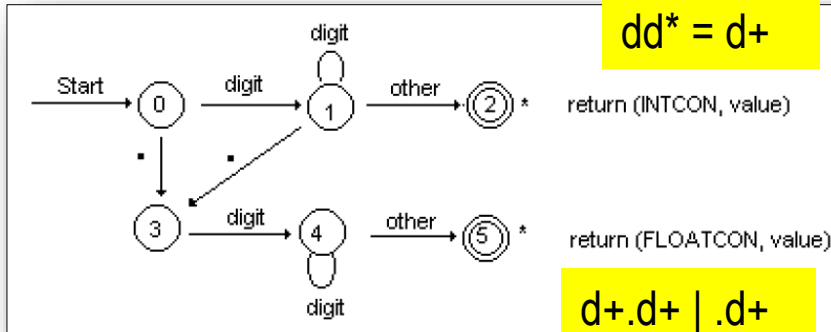
$Q \setminus \Sigma$	digit	period	other(#)
-> 0	1	3	-
1	1	3	2
2*	-	-	-
3	4	-	-
4	4	-	5
5*	-	-	-

Other = $\wedge \text{digit} \ \&\& \ \wedge$.
 $\Sigma = (\text{digit}, ., \text{other})$

Remember Kleene Theorem

- Model 3: Automaton:**

- Functional way to define the evolution of an acceptable string in a language.
- Can be visual such as:



$dd^* = d^+$

$d+.d+ \mid .d^+$

$q_1 = \delta(q_0, \text{digit})$

$q_3 = \delta(q_0, .), \delta(q_1, .)$

$Q = \text{set of states} = \{q_0 \dots q_5\}$

Finite Automaton:

$Q_f = \{q_2, q_5\}$

- Mathematical representation** of transitions that transform inputs in outputs that describes a language.
$$L(G) = A(\Sigma, Q, q_0, Q_f, \Delta)$$
- This notation includes the **alphabet** (Σ) that, **starting** in a **state** (q_0) can perform words in the **end** (Q_f), by **productions** (Δ) between **states** (Q).

Note: Empty string (ϵ) is also acceptable.

$L(G) = \text{Language from a Grammar (RE)}$

Formalization (1)

DFA

NFA

- **FA:**

$$\text{FA} = (\Sigma, Q, q_0, Q_f, \Delta)$$

- Where:
- Σ = Alphabet = $\{a_0, a_1, \dots, a_n\}$
- Q = Set of states = $\{q_0, q_1, \dots, q_m\}$
- q_0 = Initial state;
- Q_f = Final states;
- $\Delta = P$ = Production rules: $\{q_y = \delta(q_x, a_m), \dots, q_z = \delta(q_y, a_n)\}$

NFA vs DFA

1. **DFA**: Deterministic, once you are in a state and read a symbol, you know exactly where to go.
2. **NFA**: Indeterministic because:
 - It is possible to have **more than one** transition when read a symbol;
 - Null (**Epsilon**) transitions: you can go to another state reading **nothing**.



Compilers – Art. 7

Concluding

Review

- *Importance of Kleene Theorem.*

Some Questions

1. *Why to use different models?*
2. *How to chose a model for a specific representation?*
3. *Can you see advantages and disadvantages of each model?*



Source:
https://static.wixstatic.com/media/7594af_51a81a8ccc5f418281f52c8bdd2dd618~mv2.jpg



Open questions...

- Any doubts / questions?
- How we are until now?



az.allevants.in/events3/banners/26b150363d5757da8578fa6a1481585a368f12d4247adfe95837e6ea6c5ab2af-rimg-w1200-h549-gmir.jpg?v=1569691155

Image URL: <https://cdn-6c5ab2af-rimg-w1200-h549-gmir.jpg?v=1569691155>





Compilers – Art. 7

Thank you for your attention