# Pre-lecture brain teaser

Consider the problem of a n-input *AND* function. The input (*x*) is a string n-digits long with $\Sigma = \{0, 1\}$ and has an output (*y*) which is the logical *AND* of all the elements of *x*.

Formulate a **language** that describes the above problem.

# CS/ECE-374: Lecture 2 - Regular Languages

Lecturer: Nickvash Kani
Chat moderator: Samir Khan

January 28, 2021
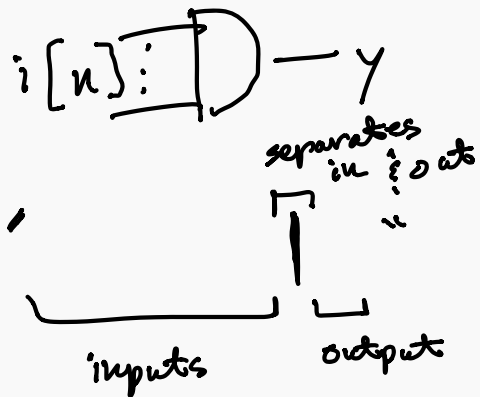
University of Illinois at Urbana Champaign

Consider the problem of a n-input *AND* function. The input ($x$) is a string n-digits long with $\Sigma = \{0, 1\}$ and has an output ($y$) which is the logical *AND* of all the elements of $x$. $n > 0$

Formulate a **language** that describes the above problem.

$i[n]: \quad \rightarrow y$

separates
in & out

inputs     output

$L_{AND} = \{ \text{"}0|0\text{"}, \text{"}1|1\text{"},$

$\text{"}00|0\text{"}, \text{"}01|0\text{"}, \text{"}10, 0\text{"}, \text{"}11|1\text{"},$

$\text{"}000|0\text{"}, \ldots \quad \ldots \quad \ldots, \text{"}111|1\text{"},$
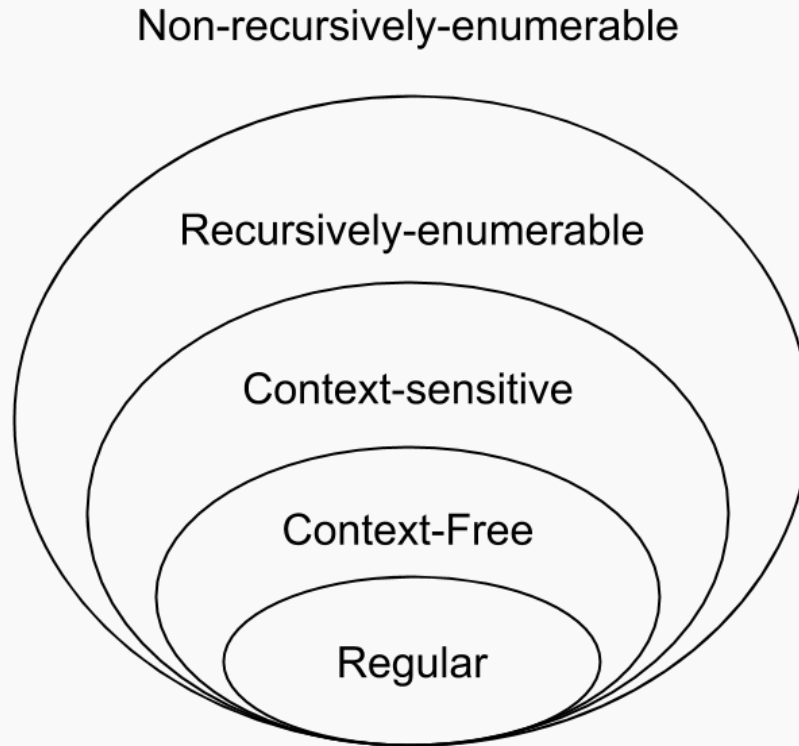
$\text{"}0000|0\text{"}, \ldots \quad -. \}$

# Pre-lecture brain teaser

Consider the problem of a n-input *AND* function. The input ($x$) is a string n-digits long with $\Sigma = \{0, 1\}$ and has an output ($y$) which is the logical *AND* of all the elements of $x$.

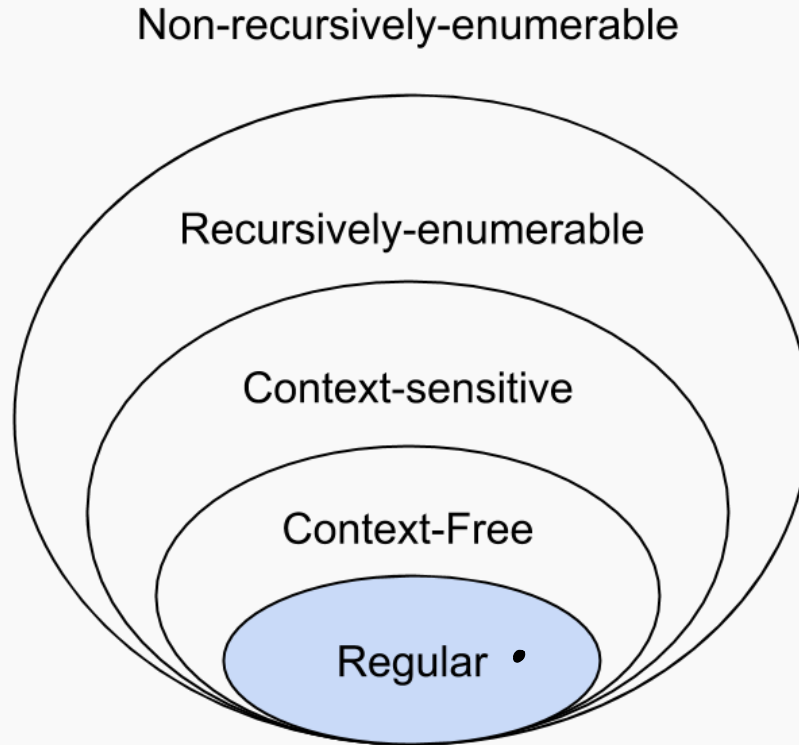Formulate a **language** that describes the above problem.

**This is an example of a regular language which we'll be discussing today.**

Non-recursively-enumerable

Recursively-enumerable

Context-sensitive
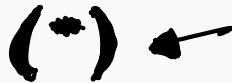
Context-Free

Regular

# Regular Languages

**Theorem (Kleene's Theorem )**

*A language is regular if and only if it can be obtained from finite languages by applying the three operations:*

- *Union*
- *Concatenation*
- *Repetition* $(\ ^*)$ ←

*a finite number of times.*

A class of simple but useful languages.

The set of regular languages over some alphabet $\Sigma$ is defined inductively.

## Base Case

- $\emptyset$ is a regular language.
- $\{\epsilon\}$ is a regular language.
- $\{a\}$ is a regular language for each $a \in \Sigma$. Interpreting $a$ as string of length 1.

$$\{``ab", ``ba"\}$$

**Inductive step:**

We can build up languages using a few basic operations:

- If $L_1, L_2$ are regular then $L_1 \cup L_2$ is regular.
- If $L_1, L_2$ are regular then $L_1 L_2$ is regular.
- If $L$ is regular, then $L^* = \cup_{n \geq 0} L^n$ is regular.
  The $\cdot^*$ operator name is *Kleene star*.
- If $L$ is regular, then so is $\overline{L} = \Sigma^* \setminus L$.

Regular languages are closed under operations of union, concatenation and Kleene star.

Have basic operations to build regular languages.

Important: Any language generated by a finite sequence of such operations is regular.

### Lemma
*Let $L_1, L_2, \ldots,$ be regular languages over alphabet $\Sigma$. Then the language $\cup_{i=1}^{\infty} L_i$ is not necessarily regular.*

$$L = \bigcup_{i=1}^{\infty} L_i$$

Example:

Define $L_i = \{0^i 1^i\}$ * One string in each $L_i$

$$L = \bigcup_{i=1}^{\infty} L_i = \{0^n 1^n \mid n \geqslant 0\} = \{01\}^* $$

$0$, $1$

$00011$

$0000000011$

$= 0^* 1^*$

$= \{01\}^*$   $010101$

$01$
$0011$
$000111$

not regular

# Some simple regular languages

**Lemma**
*If w is a string then L = {w} is regular.*

**Example:** {*aba*} or {*abbabbab*}. Why?

$$L_a = \{``a"\}$$

$$L_b = \{``b"\}$$

$$\{aba\} = L_a \, L_b \, L_a$$

**Lemma**
*If w is a string then L = {w} is regular.*

**Example:** $\{aba\}$ or $\{abbabbab\}$. Why?

**Lemma**
*Every finite language L is regular.*

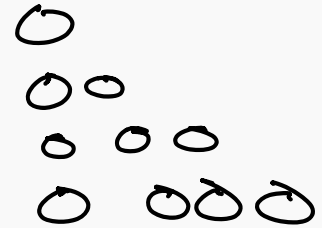Examples: $L = \{a, abaab, aba\}$. $L = \{w \mid |w| \leq 100\}$. Why?

$$\mathbb{R} = \{0, \quad , '\}$$

$$0.000000$$
$$0.12345\ldots$$

1. $L_1 = \left\{ 0^i \mid i = 0, 1, \ldots, \infty \right\}$. The language $L_1$ is regular. T/F?

$L_0 = \{"0"\}$ $\qquad$ $L_1 = L_0$

1. $L_1 = \left\{ 0^i \mid i = 0, 1, \ldots, \infty \right\}$. The language $L_1$ is regular. T/F?

2. $L_2 = \left\{ 0^{17i} \mid i = 0, 1, \ldots, \infty \right\}$. The language $L_2$ is regular. **TRUE**
   T/F?

$$L_0 = \{0\}$$
$\uparrow$
reg

$$L_{170} = \{ 0 \, L_0 \overset{17 \ times}{\cdots} L_0$$
$\uparrow$
reg

$$L_2 \{ L_{17 \cdot 0s} \}$$
1 operation
$\downarrow$

$\uparrow$
reg

1. $L_1 = \left\{ 0^i \mid i = 0, 1, \ldots, \infty \right\}$. The language $L_1$ is regular. T/F?

2. $L_2 = \left\{ 0^{17i} \mid i = 0, 1, \ldots, \infty \right\}$. The language $L_2$ is regular. T/F?

3. $L_3 = \left\{ 0^i \mid i \text{ is divisible by } 2, 3, \text{ or } 5 \right\}$. $L_3$ is regular. T/F?

$$L_0 = \{ "0" \} \qquad L_0 L_0 = \{ "00" \}$$

$$L_{i/2} = (L_0 L_0)^* \longleftarrow reg$$

$$L_{i/3} = (L_0 L_0 L_0)^* \longleftarrow reg$$

$$L_{i/5} = (L_0 L_0 L_0 L_0 L_0)^* \longleftarrow reg$$

$$L_3 = L_{i/2} \cup L_{i/3} \cup L_{i/5}$$

all binary strings

4. $L_4 = \{w \in \{0,1\}^* \mid w$ has at most ~~374~~ 3 1s$\}$. $L_4$ is regular. **TRUE**
   T/F?

$w$ is a binary string that has at most 3 1's

$L_0 = \{\text{"}0^*\text{"}\}$    $L_\epsilon = \{\epsilon\}$    $L_{\epsilon 1} = L_\epsilon \cup L_1$

$L_1 = \{\text{"}1^*\text{"}\}$    $L_q = L_0^* L_{\epsilon 1} L_0^* L_{\epsilon 1} L_0^* L_{\epsilon 1} L_0^*$
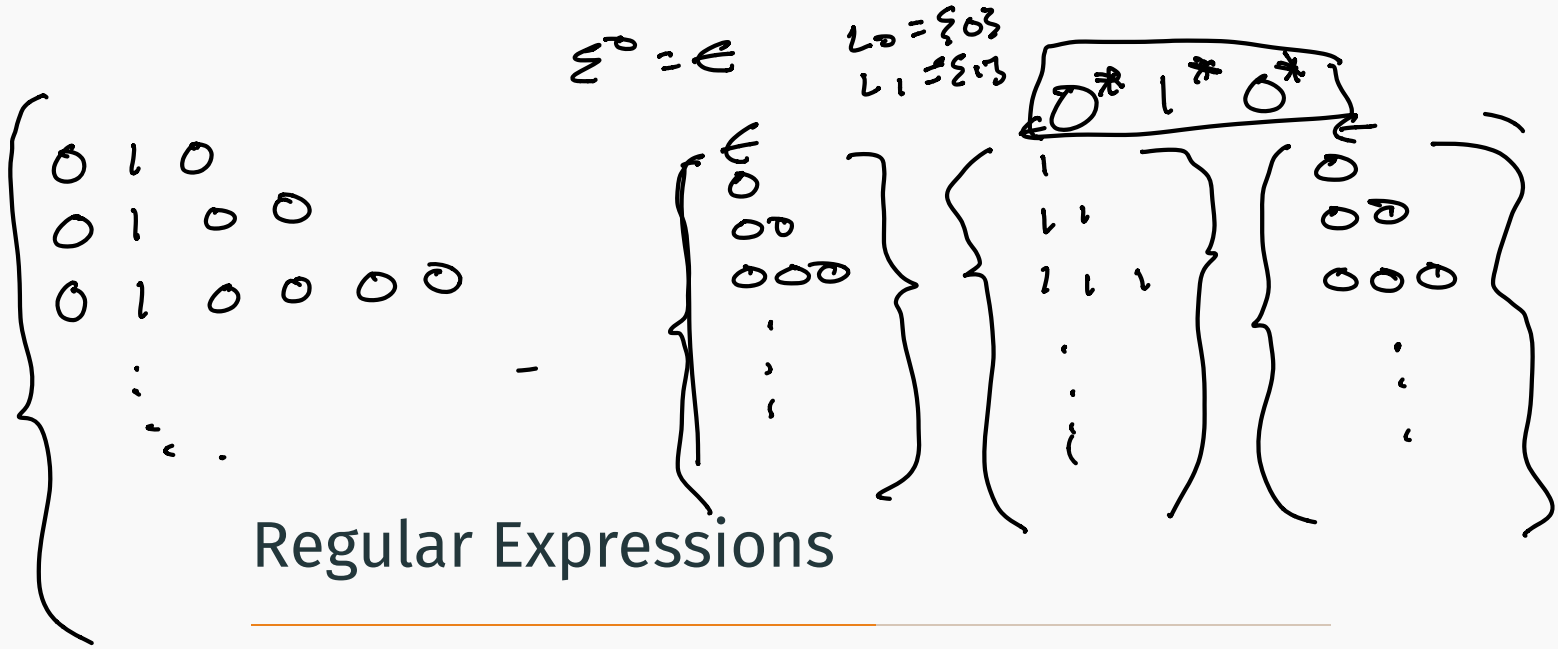
reg $\to$ $L_{4-3} = L_0^* L_1 L_0^* L_1 L_0^* L_1 L_0^* \leftarrow$ all string with 3 1's

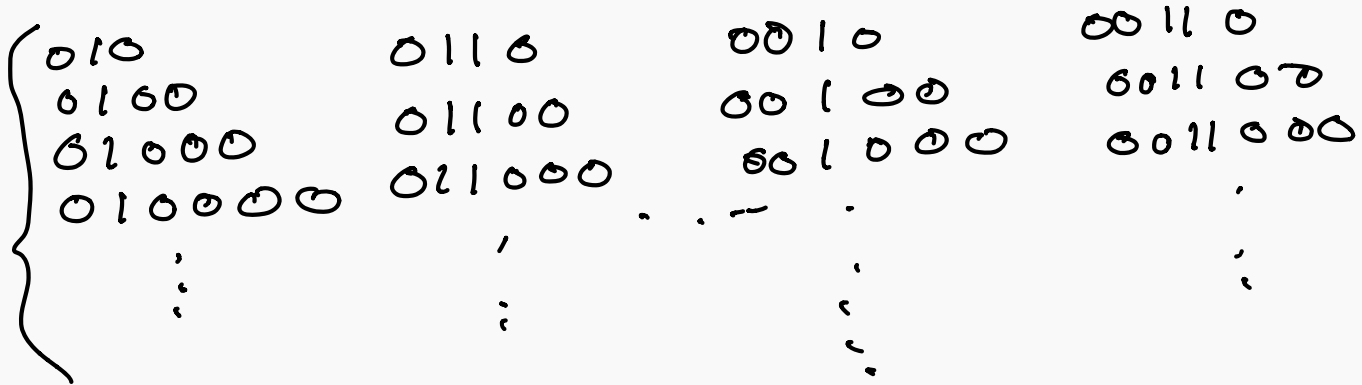reg $\to$ $L_{4-2} = L_0^* L_1 L_0^* L_1 L_0^*$    $L_{4-1} = \ldots$    $\leftarrow$ all strings with 2 1's

reg $\nearrow$ $L_{4-0} = \ldots$

$0^* 1$

reg $\to$ $L_q = L_{4-0} \cup L_{4-1} \cup L_{4-2} \cup L_{4-3}$

$\varepsilon^0 = \varepsilon$

$L_0 = \{\varepsilon\}$
$L_1 = \{\varepsilon\}$

$0^* \, 1^* \, 0^*$

$$\left\{ \begin{array}{l} 0\ 1\ 0 \\ 0\ 1\ 0\ 0 \\ 0\ 1\ 0\ 0\ 0\ 0 \\ \quad \vdots \end{array} \right. \qquad - \qquad \left\{ \begin{array}{l} \varepsilon \\ 0 \\ 00 \\ 000 \\ \vdots \end{array} \right\} \quad \left\{ \begin{array}{l} 1 \\ 1\ 1 \\ 1\ 1\ 1 \\ \vdots \end{array} \right\} \quad \left\{ \begin{array}{l} 0 \\ 0\ 0 \\ 0\ 0\ 0 \\ \vdots \end{array} \right\}$$

## Regular Expressions

$$L_1 L_2 = \{ xy \mid x \in L_1, \ y \in L_2 \}$$

$$\left\{ \begin{array}{l} 010 \\ 0100 \\ 01000 \\ 010000 \\ \vdots \end{array} \right. \qquad \begin{array}{l} 0110 \\ 01100 \\ 011000 \\ \vdots \end{array} \qquad \begin{array}{l} 0010 \\ 00100 \\ 001000 \\ \vdots \end{array} \qquad \begin{array}{l} 00110 \\ 001100 \\ 0011000 \\ \vdots \end{array}$$

A way to denote regular languages

- simple patterns to describe related strings
- useful in
    - text search (editors, Unix/grep, emacs)
    - compilers: lexical analysis
    - compact way to represent interesting/useful languages
    - dates back to 50's: Stephen Kleene
      who has a star names after him [1].

---

[1]Kleene, Stephen C.: "Representation of Events in Nerve Nets and Finite Automata". In Shannon, Claude E.; McCarthy, John. Automata Studies, Princeton University Press. pp. 3–42., 1956.

# Inductive Definition

A regular expression **r** over an alphabet $\Sigma$ is one of the following:

**Base cases:**

- $\emptyset$ denotes the language $\emptyset$
- $\epsilon$ denotes the language $\{\epsilon\}$.
- *a* denote the language $\{a\}$.

**Inductive cases:** If $\mathbf{r_1}$ and $\mathbf{r_2}$ are regular expressions denoting languages $R_1$ and $R_2$ respectively then,

- $(\mathbf{r_1} + \mathbf{r_2})$ denotes the language $R_1 \cup R_2$
- $(\mathbf{r_1} \cdot \mathbf{r_2}) = r_1 \cdot r_2 = (\mathbf{r_1 r_2})$ denotes the language $R_1 R_2$
- $(\mathbf{r_1})^*$ denotes the language $R_1^*$

# Regular Languages vs Regular Expressions

### Regular Languages

$\emptyset$ regular

$\{\epsilon\}$ regular

$\{a\}$ regular for $a \in \Sigma$

$R_1 \cup R_2$ regular if both are

$R_1 R_2$ regular if both are

$R^*$ is regular if $R$ is

### Regular Expressions

$\emptyset$ denotes $\emptyset$

$\epsilon$ denotes $\{\epsilon\}$

$\mathbf{a}$ denote $\{a\}$

$\mathbf{r_1 + r_2}$ denotes $R_1 \cup R_2$

$\mathbf{r_1 \cdot r_2}$ denotes $R_1 R_2$

$\mathbf{r^*}$ denote $R^*$

Regular expressions denote regular languages — they explicitly show the operations that were used to form the language

- For a regular expression **r**, $L(\mathbf{r})$ is the language denoted by **r**. Multiple regular expressions can denote the same language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$

- For a regular expression r, $L(r)$ is the language denoted by
  r. Multiple regular expressions can denote the same
  language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$
- Two regular expressions $r_1$ and $r_2$ are equivalent if
  $L(r_1) = L(r_2)$.

- For a regular expression $r$, $L(r)$ is the language denoted by $r$. Multiple regular expressions can denote the same language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$
- Two regular expressions $r_1$ and $r_2$ are equivalent if $L(r_1) = L(r_2)$.
- Omit parenthesis by adopting precedence order: $*$, concatenate, $+$.
  **Example:** $r^*s + t = ((r^*)s) + t$

- For a regular expression $r$, $L(r)$ is the language denoted by $r$. Multiple regular expressions can denote the same language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$
- Two regular expressions $r_1$ and $r_2$ are equivalent if $L(r_1) = L(r_2)$.
- Omit parenthesis by adopting precedence order: $*$, concatenate, $+$.
  **Example:** $r^*s + t = ((r^*)s) + t$
- Omit parenthesis by associativity of each of these operations.
  **Example:** $rst = (rs)t = r(st)$,
  $r + s + t = r + (s + t) = (r + s) + t$.

- For a regular expression r, $L(r)$ is the language denoted by r. Multiple regular expressions can denote the same language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$
- Two regular expressions $r_1$ and $r_2$ are equivalent if $L(r_1) = L(r_2)$.
- Omit parenthesis by adopting precedence order: $*$, concatenate, $+$.
  **Example:** $r^*s + t = ((r^*)s) + t$
- Omit parenthesis by associativity of each of these operations.
  **Example:** $rst = (rs)t = r(st)$,
  $r + s + t = r + (s + t) = (r + s) + t$.
- Superscript $+$. For convenience, define $r^+ = rr^*$. Hence if $L(r) = R$ then $L(r^+) = R^+$.

$\varepsilon^0 = \epsilon$

no epsilon

- For a regular expression **r**, $L(\mathbf{r})$ is the language denoted by **r**. Multiple regular expressions can denote the same language!
  **Example:** $(0 + 1)$ and $(1 + 0)$ denote same language $\{0, 1\}$
- Two regular expressions $\mathbf{r_1}$ and $\mathbf{r_2}$ are equivalent if $L(\mathbf{r_1}) = L(\mathbf{r_2})$.
- Omit parenthesis by adopting precedence order: $*$, concatenate, $+$.
  **Example:** $r^*s + t = ((r^*)s) + t$
- Omit parenthesis by associativity of each of these operations.
  **Example:** $rst = (rs)t = r(st)$,
  $r + s + t = r + (s + t) = (r + s) + t$.
- Superscript $+$. For convenience, define $\mathbf{r}^+ = \mathbf{rr}^*$. Hence if $L(\mathbf{r}) = R$ then $L(\mathbf{r}^+) = R^+$.

# Some examples of regular expressions

1. $(0+1)^*$: All binary strings

$L_0 =$

$L_1 =$

$(L_0 \cup L_1)^*$

1. $(0 + 1)^*$:    could be
                      $\epsilon$
2. $(0 + 1)^*001(0 + 1)^*$: ⟵

   All strings with 001 as substring

   0101 001

1. $(0 + 1)^*$:
2. $(0 + 1)^*001(0 + 1)^*$:
3. $0^* + \boxed{(0^*10^*10^*10^*)}^*$: $\longleftarrow$ # of 1's divisible by 3

$$L = \begin{cases} 0 \quad 000100001\ 0010 \\ \\ 0+1\ 1\ 111\ 1 \end{cases}$$

1. $(0 + 1)^*$:

2. $(0 + 1)^*001(0 + 1)^*$:

3. $0^* + (0^*10^*10^*10^*)^*$:

4. $(\epsilon + 1)(01)^*(\epsilon + 0)$:

*(handwritten annotations)*

1 0 1 0 1   $\epsilon$

$\epsilon$   0 1 0 0 0 1 0

alternating $\left\{ \begin{array}{l} \epsilon \\ 1\ 0\ 1\ 0\ 1\ 0\ 1\ 0 \\ 1\ 0 \end{array} \right.$

1. All strings that end in 1011? ←

$$(0+1)^* \cdot 1011$$

# Creating regular expressions

1. All strings that end in 1011?

2. All strings except 11?

All string 3 characters or greater

$$\cup \{00, 01, 10\}$$

$$\cup \{0, 1\} \cup \{\epsilon\}$$

$$\epsilon + 0 + 1 + 00 + 01 + 10 + (0+1)^3 (0+1)^*$$

1. All strings that end in 1011?

2. All strings except 11?

3. All strings that do not contain 000 as a subsequence?

$abc$ is a subsequence of string $w$ if the letters 'a' 'b' 'c' appear in the string in that order but not neccessarily consecutively

$Ex = gabtvc$
$= gcatub$

$$1^*(\epsilon + 0) \ 1^* \ (\epsilon + 0) \ 1^*$$

$00: \quad \epsilon \quad 0 \quad \epsilon \quad 0 \quad \epsilon \ = \ 00$

$0: \ \epsilon \quad \epsilon \quad \epsilon \quad 0 \quad \epsilon \ = \ 0$

1. All strings that end in 1011?

2. All strings except 11?

3. All strings that do not contain 000 as a subsequence?

4. All strings that do not contain the substring 10?

All zeros must come before all 1's

$0^* 1^*$

Consider the problem of a n-input *AND* function. The input (*x*) is a string n-digits long with $\Sigma = \{0, 1\}$ and has an output (*y*) which is the logical *AND* of all the elements of *x*.

Formulate the regular expression which describes the above language:

$$L = \{ \text{``}010\text{''} \quad \text{``}111\text{''} \quad \text{``}0010\text{''} \quad \text{``}0010\text{''} \quad \text{``}1010\text{''} \quad \text{``}1111\text{''} \quad \ldots$$

Case where $y = 1$ $\qquad 1^* \; 1^n 1$

Case where $y = 0$ $\qquad (0+1)^* 0 \; (1+0)^* 0$

$$\boxed{1^+ \; 1^n 1 \; + \; (0+1)^* 0 \; (1+0)^* 0}$$

# Regular expressions in programming

# One last expression….

# Bit strings with odd number of 0s and 1s

The regular expression is

$$(00 + 11)^*(01 + 10)$$
$$\left(00 + 11 + (01 + 10)(00 + 11)^*(01 + 10)\right)^*$$

The regular expression is

$$\left(00 + 11\right)^*(01 + 10)$$
$$\left(00 + 11 + (01 + 10)(00 + 11)^*(01 + 10)\right)^*$$

(Solved using techniques to be presented in the following lectures...)