# Theory of Computation 1
# Sets and Regular Expressions

**Frank Stephan**

**Department of Computer Science**

**Department of Mathematics**

**National University of Singapore**

**fstephan@comp.nus.edu.sg**

# Languages

Examples of Languages (Sets)

(a) Set of all natural numbers in binary notation: either $0$ or $1$ followed by arbitrarily many digits from $0, 1$.

(b) Set of all possible computer programs in syntax of programming language C: Tools can translate a formal description of C syntax into a syntax checker.

(c) Set of all C programs which pass a compiler without error messages: Compilers check more than just syntactical correctness.

(d) Set of all C programs which do not have bugs: No computer program can solve this task completely.

(e) Set of all texts of books written in English and published between 1066 and 1492: Exhaustive list describes this set.

Theory of Computation

How does one describe above sets? How does one modify descriptions? Do descriptions allow membership checks?

# Languages

Language = Set of Strings over an Alphabet.
Alphabet $\Sigma$, for example $\Sigma = \{0, 1, 2\}$. Always finite.

Finite languages

$L_1 = \emptyset$, no elements.

$L_2 = \{\varepsilon\}$, set consisting of empty string.

$L_3 = \{00, 01, 02, 10, 11, 12, 20, 21, 22\}$, all elements of length $2$.

$L_4 = \{\varepsilon, 0, 00, 000, 0000\}$, all strings of $0$s up to length $4$.

$L_5 = \{01, 001, 02, 002\}$, all strings consisting of one or two $0$s followed by a $1$ or $2$.

# Operations with Languages

Union:

$L \cup H = \{u : u \in L \text{ or } u \in H\}$;

$\{00, 01, 02\} \cup \{01, 11, 21\} = \{00, 01, 02, 11, 21\}$;

$\{0, 00, 000\} \cup \{00, 000, 0000\} = \{0, 00, 000, 0000\}$.

Intersection:

$L \cap H = \{u : u \in L \text{ and } u \in H\}$;

$\{00, 01, 02\} \cap \{01, 11, 21\} = \{01\}$;

$\{0, 00, 000\} \cap \{00, 000, 0000\} = \{00, 000\}$.

Set Difference:

$L - H = \{u : u \in L \text{ and } u \notin H\}$;

$\{00, 01, 02\} - \{01, 11, 21\} = \{00, 02\}$.

Concatenation:

$000 \cdot 1122 = 0001122$;

$L \cdot H = \{v \cdot w : v \in L \text{ and } w \in H\}$;

$\{0, 00\} \cdot \{1, 2\} = \{01, 001, 02, 002\}$.

# Kleene Star and Plus

Definition

$$L^* = \{\varepsilon\} \cup L \cup L \cdot L \cup L \cdot L \cdot L \cup \ldots$$
$$= \{w_1 \cdot w_2 \cdot \ldots \cdot w_n : n \geq 0 \text{ and } w_1, w_2, \ldots, w_n \in L\};$$
$$L^+ = L \cup L \cdot L \cup L \cdot L \cdot L \cup \ldots$$
$$= \{w_1 \cdot w_2 \cdot \ldots \cdot w_n : n > 0 \text{ and } w_1, w_2, \ldots, w_n \in L\}.$$

Examples

$\emptyset^* = \{\varepsilon\}$.

$\Sigma^*$ is the set of all words over $\Sigma$.

$\{0\}^* = \{\varepsilon, 0, 00, 000, 0000, \ldots\}$.

$\{00, 01, 10, 11\}^*$ are all binary words of even length.

$\varepsilon \in L^+$ iff $\varepsilon \in L$.

Notation

Often $w^*$ in place of $\{w\}^*$;

Often $w \cdot L$ in place of $\{w\} \cdot L$.

# Regular Languages

Regular expressions are either finite sets listed by their elements or obtained from other regular expressions by forming the Kleene star, Kleene plus, union, intersection, set-difference or concatenation.

A language is regular iff it can be described by a regular expression.

Regular sets have many different regular expressions.

For example, $\{0, 00\} \cdot \{1, 2\}$ and $\{01, 001, 02, 002\}$ describe the same set. Also $0^*$ and $(00)^* \cup 0 \cdot (00)^*$ describe the same set.

Intersections and set difference are traditionally not used in regular expressions, as every regular set has an expression only using union, concatenation and Kleene star.

The complement of a language $L$ is $\Sigma^* - L$.

# Quiz

Which three sets are described by two of the following regular expressions?

1. $\{00, 000\}^+$;

2. $\{000, 0000\}^+$;

3. $00 \cdot 0^*$;

4. $000 \cdot 0^*$;

5. $\{000, 0000\} \cup (000000 \cdot 0^*)$;

6. $\{00, 01, 02, 10, 11, 12\}$;

7. $0^*1^*2^*$;

8. $(0^*1^*2^*)^*$;

9. $(\{0, 1\} \cdot \{0, 1, 2\}^*) \cap (\{0, 1, 2\} \cdot \{0, 1, 2\})$.

# Exercises 1.6, 1.7 and 1.8

## Exercise 1.6

Assume $A$ has $3$ and $B$ has $2$ elements. How many elements do the following sets have at least and at most; it depends on the actual choice which of the bounds is realised: $A \cup B$, $A \cap B$, $A \cdot B$, $A - B$, $A^* \cap B^*$.

## Exercise 1.7

Let $A, B$ be finite sets and $|A|$ be the number of elements of $A$. Is the following formula correct:

$$|A \cup B| + |A \cap B| = |A| + |B|?$$

Prove your answer.

## Exercise 1.8

Make a regular expression without intersection and set difference for $0^*1^*0^*1^* \cap (11)^*(00)^*(11)^*(00)^*$.

# Tutorial

Register for a tutorial group

During the semester, you can make up to TEN marks in tutorials (full course = 100 marks).

(a) Write up to four exercises up in the Forum, at most one per week; these can be reserved after the lecture at the lecturer; two marks for each exercise, at most eight marks.

(b) Present in the tutorial group one exercise per week; you have in your own group priority on exercises you write up, however, you can also present other exercises not signed up by any member of your group; five presentations, one mark; eight presentations, two marks; one presentation per week is counted.

# Tests

Midterm Tests are inside the class and take up a part but not all of the lectures on Tuesday 17 September 2019 and Tuesday 29 October 2019; each midterm test counts 20 marks. The midterm tests are in the second half of the lectures on those days.

The final examination counts 50 marks. It is 28 November 2019 at 17:00 hrs and the duration is 2 hours. Please consult NUS webpages for more information and doublecheck the information there.

Summary: Tutorial 10 marks, Mitderms 20 + 20 marks, Final 50 marks.
https://nusmods.com/modules/CS4232/theory-of-computation

# Theorem of Lyndon and Schützenberger

## Theorem

If two words $v, w$ satisfy $vw = wv$ then $\exists u\, [v, w \in u^*]$;

If all words $v, w \in L$ satisfy $vw = wv$ then $\exists u\, [L \subseteq u^*]$.

## Proof

Case $v = \varepsilon$ or $w = \varepsilon$: $u = vw$.

Case $v \neq \varepsilon$ and $w \neq \varepsilon$ and $|v| = |w|$: $v = w$.

Case $v \neq \varepsilon$ and $w \neq \varepsilon$ and $|v| < |w|$: Let $k$ be greatest common divisor of $|v|, |w|$ and $u$ be the first $k$ symbols of $v$. There are $i, j$ with $v = u_1 u_2 \ldots u_i$ and $w = u_1 u_2 \ldots u_j$ for words $u_1, u_2, \ldots, u_j$ of length $k$.

Note that $v^j w^i = w^i v^j$ and $|v^j| = |w^i|$, as both have the length $ijk$. Thus $v^j = w^i$.

For each $u_h$ there is a position of $u_h$ in $w^i$ where $u_1$ is at the same position in $v^j$. Thus $u_h = u_1$.

So $v, w \in u^*$ for $u = u_1$.

# Example, Second Part

Example: Let $v = abcd$ and $w = abcdef$ and $vw = wv$. Now $k = 2$ (greatest common divisor of $4, 6$).

$$v^3 = ab\,cd\,ab\,cd\,ab\,cd;$$
$$w^2 = ab\,cd\,ef\,\,ab\,cd\,ef.$$

Now $ab = ab$ at $0, 1$, $ab = ef$ at $4, 5$, $ab = cd$ at $8, 9$. So $ab = cd = ef$ and $v, w \in (ab)^*$.

Second Part: Let $v$ be shortest nonempty word of $L$ and $u$ be shortest word with $v \in u^*$.
Let $w \in L$ be arbitrary.
There is $\tilde{u}$ with $v, w \in \tilde{u}^*$.
Now $\tilde{u}^i = u^j = v$ for some $i, j$.
Thus $\tilde{u}, u \in \hat{u}^*$ as in Part 1 for some $\hat{u}$ where $|\hat{u}|$ is greatest common divisor of $|u|$ and $|\tilde{u}|$.
By choice of $u$, $\hat{u} = u$ and $w \in u^*$.
So $L \subseteq u^*$.

# Structural Induction

Theorem

Let $\mathbf{P}$ be a property of sets such that the following holds:

- Every finite set (including $\emptyset$) satisfies $\mathbf{P}$;

- If $\mathbf{L}, \mathbf{H}$ satisfy $\mathbf{P}$ so does $\mathbf{L} \cup \mathbf{H}$;

- If $\mathbf{L}, \mathbf{H}$ satisfy $\mathbf{P}$ so does $\mathbf{L} \cdot \mathbf{H}$;

- If $\mathbf{L}$ satisfies $\mathbf{P}$ so does $\mathbf{L}^*$.

Then all regular sets satisfy $\mathbf{P}$.

Proof.

Let $\mathbf{L}(\sigma)$ be the set generated by the regular expression $\sigma$. Here finite sets and the operations union, concatenation and Kleene star are permitted for regular expressions.

Now it is shown that there is no shortest regular expression $\sigma$ such that $\mathbf{L}(\sigma)$ does not satisfy $\mathbf{P}$.

# Shortest Expression

Assume that $\sigma$ is a shortest expression not satisfying $\mathbf{P}$; if there are several shortest ones, $\sigma$ is just any of these.

- If $\sigma$ is a list of a finite set then $\mathbf{P}$ satisfies $\mathbf{L}(\sigma)$;

- If $\sigma = (\rho \cup \tau)$ then $\rho, \tau$ are shorter than $\sigma$ and $\mathbf{L}(\rho), \mathbf{L}(\tau)$ satisfy $\mathbf{P}$ and so does $\mathbf{L}(\sigma) = \mathbf{L}(\rho) \cup \mathbf{L}(\tau)$;

- If $\sigma = (\rho \cdot \tau)$ then $\rho, \tau$ are shorter than $\sigma$ and $\mathbf{L}(\rho), \mathbf{L}(\tau)$ satisfy $\mathbf{P}$ and so does $\mathbf{L}(\sigma) = \mathbf{L}(\rho) \cdot \mathbf{L}(\tau)$;

- If $\sigma = \tau^*$ then $\tau$ is shorter than $\sigma$ and $\mathbf{L}(\tau)$ satisfies $\mathbf{P}$ and so does $\mathbf{L}(\sigma) = (\mathbf{L}(\tau))^*$.

So there is no case in which $\mathbf{L}(\sigma)$ would not satisfy $\mathbf{P}$, thus this $\sigma$ does not exist and there is no regular expression $\sigma$ for which $\mathbf{L}(\sigma)$ does not satisfy $\mathbf{P}$. All regular languages satisfy $\mathbf{P}$.

# Strengthening the Theorem

Theorem

Let $P$ be a property of sets such that the following holds:

- The empty set and the set $\{\varepsilon\}$ satisfy $P$;

- For every letter $a$, the set $\{a\}$ satisfies $P$;

- If $L, H$ satisfy $P$ so does $L \cup H$;

- If $L, H$ satisfy $P$ so does $L \cdot H$;

- If $L$ satisfies $P$ so does $L^*$.

Then all regular sets satisfy $P$.

This strengthening is just based on the fact that every finite set of strings can be formed using concatenation and union from the sets containing a single letter word, the set containing the empty word and the empty set.

# Polynomial and Exponential Growth

## Definition

A language $L$ has polynomial growth iff there is a polynomial $p$ such that for all $n$ there are in $L$ at most $p(n)$ many words shorter than $n$.

A language $L$ has exponential growth iff there are constants $h, k$ such that $L$ contains for all $n$ at least $2^n$ words which are shorter than $h \cdot n + k$.

## Theorem

Every regular set has either polynomial or exponential growth.

This will be proven by structural induction.

# Examples

Every finite set has polynomial growth, as one plus the number of elements is a polynomial which is an upper bound as required.

The set $0^* 1^*$ has polynomial growth as there are $n(n+1)/2$ many words shorter than $n$ in this set.

The set $0^* \cup 1^*$ has polynomial growth as there are at most $2n$ many words shorter than $n$ in this set.

The set $\{00, 11\}^* \cdot \{222222\}$ has exponential growth as it has for all $n$ at least $2^n$ words shorter than $7 + 2n$.

The set $\{0000, 1111\}^*$ has exponential growth as it has for all $n$ at least $2^n$ words shorter than $1 + 4n$.

The set $\{00, 11\}^* \cdot \emptyset$ is empty and has polynomial growth.

# Rules for Growth

Finite sets have polynomial growth.

If $L$ and $H$ have polynomial growth then so do $L \cup H$ and $L \cdot H$.

If $L$ or $H$ have exponential growth then so does $L \cup H$.

The sets $L \cdot \emptyset$ and $\emptyset \cdot L$ have polynomial growth.

If $L$ and $H$ are not empty and at least one of them has exponential growth so does $L \cdot H$.

If $L$ contains $v, w$ with $vw \neq wv$ then $L^*$ has exponential growth else $L^*$ has polynomial growth.

Let $P(L)$ say that the language $L$ has either polynomial or exponential growth. Then the rules imply that all finite sets satisfy $P$ and that, whenever $L, H$ satisfy $P$ so do $L \cup H$, $L \cdot H$ and $L^*$. Thus all regular sets satisfy $P$ by structural induction.

# Quiz

Does $L \cap H$ have exponential growth whenever $L$ and $H$ have exponential growth?

Does $\{0101, 010101\}^*$ have exponential growth?

Does $\{000, 001, 011, 111\}^* \cdot \{0000, 1111\}$ have exponential growth?

Does the (non-regular) set $\{w \in \{0, 1\}^* : w$ has at most $\log(|w|)$ many $1s\}$ have polynomial growth?

Does the set $\{w \in \{0, 1\}^* : w$ has at most $\log(|w|)$ many $1s\}$ have exponential growth?

Is there a polynomial $p$ such that every set of polynomial growth has at most $p(n)$ elements shorter than $n$ for every $n$?

# Rules for Regular Expressions

**(a)** $L \cup L = L$, $L \cap L = L$, $(L^*)^* = L^*$, $(L^+)^+ = L^+$;

**(b)** $(L \cup H)^* = (L^* \cdot H^*)^*$ and if $\varepsilon \in L \cap H$ then $(L \cup H)^* = (L \cdot H)^*$;

**(c)** $(L \cup \{\varepsilon\})^* = L^*$, $\emptyset^* = \{\varepsilon\}$ and $\{\varepsilon\}^* = \{\varepsilon\}$;

**(d)** $L^+ = L \cdot L^* = L^* \cdot L$ and $L^* = L^+ \cup \{\varepsilon\}$;

**(e)** $(L \cup H) \cdot K = (L \cdot K) \cup (H \cdot K)$ and $K \cdot (L \cup H) = (K \cdot L) \cup (K \cdot H)$;

**(f)** $(L \cup H) \cap K = (L \cap K) \cup (H \cap K)$ and $(L \cap H) \cup K = (L \cup K) \cap (H \cup K)$;

**(g)** $(L \cup H) - K = (L - K) \cup (H - K)$ and $(L \cap H) - K = (L - K) \cap (H - K)$.

# Inequality Rules

**(a)** $L \cdot L$ can be different from $L$: $\{0\} \cdot \{0\} = \{00\}$;

**(b)** $(L \cap H)^* \subseteq L^* \cap H^*$;
Properness: $L = \{00\}$, $H = \{000\}$, $(L \cap H)^* = \{\varepsilon\}$,
$L^* \cap H^* = \{000000\}^*$;

**(c)** If $\{\varepsilon\} \cup (L \cdot H) = H$ then $L^* \subseteq H$;
Properness: $L = \{\varepsilon\}$, $H = \{0\}^*$;

**(d)** If $L \cup (L \cdot H) = H$ then $L^+ \subseteq H$;
Properness: $L = \{\varepsilon\}$, $H = \{0\}^*$;

**(e)** $(L \cap H) \cdot K \subseteq (L \cdot K) \cap (H \cdot K)$;
Properness: $(\{0\} \cap \{00\}) \cdot \{0, 00\} = \emptyset \subset \{000\} = (\{0\} \cdot \{0, 00\}) \cap (\{00\} \cdot \{0, 00\})$;

**(f)** $K \cdot (L \cap H) \subseteq (K \cdot L) \cap (K \cdot H)$;
Properness: $\{0, 00\} \cdot (\{0\} \cap \{00\}) = \emptyset \subset \{000\} = (\{0, 00\} \cdot \{0\}) \cap (\{0, 00\} \cdot \{00\})$.

# Characterising Kleene Star

Corollary 1.17. For any set $L$, the following statements characterise $L^*$ and $L^+$:

**(a)** $L^*$ is the smallest set $H$ such that $\{\varepsilon\} \cup (L \cdot H) = H$;

**(b)** $L^*$ is the smallest set $H$ such that $\{\varepsilon\} \cup (L \cdot H) \subseteq H$;

**(c)** $L^+$ is the smallest set $H$ such that $L \cup (L \cdot H) = H$;

**(d)** $L^+$ is the smallest set $H$ such that $L \cup (L \cdot H) \subseteq H$.

In the above, one could also use $H \cdot L$ in place of $L \cdot H$.

# Exercise 1.18

Which three of the following sets are not equal to any of the other sets:

**(a)** $\{01, 10, 11\}^*$;

**(b)** $((\{0, 1\} \cdot \{0, 1\}) - \{00\})^*$;

**(c)** $(\{01, 10\} \cdot \{01, 10, 11\} \cup \{01, 10, 11\} \cdot \{01, 10\})^*$;

**(d)** $(\{01, 10, 11\} \cdot \{01, 10, 11\})^* \cup \{01, 10, 11\} \cdot (\{01, 10, 11\} \cdot \{01, 10, 11\})^*$;

**(e)** $\{0, 1\}^* - \{0, 1\} \cdot \{00, 11\}^*$;

**(f)** $((\{01\}^* \cup \{10\})^* \cup \{11\})^*$;

**(g)** $(\{\varepsilon\} \cup (\{0\} \cdot \{0, 1\}^* \cap \{1\} \cdot \{0, 1\}^*))^*$.

Explain your answer.

# Exercise 1.19

Make a regular expression which contains all those decimal natural numbers which start with $3$ or $8$ and have an even number of digits and end with $5$ or $7$.

Make a further regular expression which contains all odd ternary numbers without leading $0$s; here a ternary number is a number using the digits $0, 1, 2$ with $10$ being three, $11$ being four and $1212$ being fifty. The set described should contain the ternary numbers $1, 10, 12, 21, 100, 102, 111, 120, 122, 201, \ldots$ which are the numbers $1, 3, 5, 7, 9, 11, 13, 15, 17, 19, \ldots$ in decimal.

# Exercise 1.20

Let $S$ be the smallest class of languages such that

- every language of the form $u^*$ for a nonempty word $u$ is in $S$;

- the union of two languages in $S$ is again in $S$;

- the concatenation of two languages in $S$ is again in $S$.

Prove by structural induction the following properties of $S$:

**(a)** Every language in $S$ is infinite;

**(b)** Every language in $S$ has polynomial growth.

Lay out all inductive steps explicitly without only citing results in this lecture.

# Exercise 1.21

Let $L$ satisfy the following statement: For all $u, v, w \in L$, either $uv = vu$ or $uw = wu$ or $vw = wv$. Which of the following statements are true for all such $L$:

- All $x, y \in L$ satisfy $xy = yx$;
- All sufficiently long $x, y \in L$ satisfy $xy = yx$;
- The language $L$ has polynomial growth.

Give an answer to these questions and prove them.

# Exercises 1.22-1.23

In the following, each digit, the symbol $\varepsilon$, the symbol $\cdot$, the symbol $\cup$, the comma and each set bracket and each normal bracket have length $1$ and the length of the expression is the number of all the symbols in it (counting repetitions).

## Exercise 1.22

Let $L$ consist of all words which contain each of the letters $0, 1, 2, 3$ exactly once. Make a regular expression generating $L$ which has at most length $100$.

## Exercise 1.23

Make a regular expression for the set $\{w \in \{0\}^* : |w| \le 9\}$ which has at most length 26.

# Exercises 1.24-1.26

Let $V$ be the set of vowels, $W$ be the set of consonants and $S$ be the set of punctuation marks and $T$ be the set of spacings (blancs and new lines and so on).

Exercise 1.24. Make a regular expression (using above sets) of all words which contain at least two vowels and before, after and between vowels is exactly one consonant.

Exercise 1.25. Make a regular expression of all sentences where each sentence consists of words containing one vowel and arbitrarily many consonants and between two words are spacings and after the last word is a punctuation mark.

Exercise 1.26. Make a regular expressions generating texts of sentences separted by spacings where sentences are as above with the only difference that words can have one or two vowels and up to four consonants.