# Formal Languages and Parsing

## CS 462

Jeffrey Shallit

# Preface

**Disclaimer**   Much of the information on this set of notes is transcribed directly/indirectly from the lectures of CS 462 during Winter 2022 as well as other related resources. I do not make any warranties about the completeness, reliability and accuracy of this set of notes. Use at your own risk.

For any questions, send me an email via https://notes.sibeliusp.com/contact.

You can find my notes for other courses on https://notes.sibeliusp.com/.

*Sibelius Peng*

# Contents

# CS 462 notation

- Natural numbers $\mathbb{N} = \{0, 1, 2, \ldots\}$ and we use letters $i, j, k, \ell, m, n \in \mathbb{N}$.

- Finite string/word: a map from $[0, n-1]$ (an interval) to $\Sigma$ (a finite alphabet of symbols)

  $w[i]$ is $i$th symbol of $w$

- infinite strings/words: a map from $\mathbb{N}$ to $\Sigma$. We denote infinite strings by bold-face:

$$\mathbf{w} = \mathbf{w}[0]\mathbf{w}[1]\mathbf{w}[2]\cdots$$

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

- $\Sigma^\omega$ is the set of all infinite words over $\Sigma$. Also written $\Sigma^{\mathbb{N}}$.

- $\Sigma^\infty = \Sigma^* \cup \Sigma^{\mathbb{N}}$.

Finite words typically denote by $s, t, u, v, w, x, y, z$

## 1.1 Some refreshers from CS 360/365

- $x$ is a **prefix** of $z$ if there exists $y$ such that $z = xy$

- $x$ is a **suffix** of $z$ if there exists $y$ such that $z = yx$

- $x$ is a **subword** (factor) of $z$ if there exists $w, y$ such that $z = wxy$.

- $x$ is a **subsequence** of $z$ if $x$ can be obtained from $z$ by striking out zero or more symbols.

Remark:

Does substring mean contiguous (like subword)? or noncontiguous (like subsequence)? This definition depends the author of the book.

Empty string $\epsilon$ is a first-class string like any other string and is not ruled out unless done so explicitly.

Then we have "proper" prefix, suffix, etc. If $z = xy$ and $x \neq z$, then $x$ is a **proper prefix** of $z$.

## 1.2   Some notations

A shorthand for subword:

$$w[a..b] = w[a]w[a+1] \cdots w[b]$$

Concatenation of strings:

which is not commutative in general. Because we write concatenation in a multiplicative way, we can raise strings to powers: $x^n = \underbrace{xx \cdots x}_{n \text{ times}}$, or formally

$$x^0 = \epsilon$$
$$x^n = x \cdot x^{n-1} \qquad n \geq 1$$
$$x^{m+n} = x^m x^n$$

A word is not of the form $z^n$, $n \geq 2$, $z \neq \epsilon$ is called **primitive**. The set of binary primitive words are denoted

$$P_2 = \{0, 1, 01, 10, 001, 010, 011, \ldots\}$$

One open question: is $P_2$ context-free? Probably not! But no one knows a proof.

## 1.3   Other operations on words

We define perfect shuffle on $x$ and $y$, for $|x| = |y| = n$ as

$$x \text{ Ш } y = x[1]y[1]x[2]y[2] \cdots x[n]y[n]$$

where Ш is the Russian "sha". For example,

$$\text{term Ш hoes} = \text{theorems}$$

Single symbols are denoted by $a, b, c \in \Sigma$.

The reversal: $x^R$, symbols of $x$ in reverse order. If you feel stressed, we can reverse it and get

$$(\text{stressed})^R = (\text{desserts})$$

Palindromes: $x = x^R$.

### Ordering

**lexicographic order**   We define it for the words of same length, $|x| = |y|$. Then $x < y$ means[1] there exists $i$ such that $1 \leq i \leq n = |x| = |y|$, and $x[j] = y[j]$ for $j < i$ and $x[i] < y[i]$. $x \leq y$ means $x = y$ or $x < y$.

**radio order**   $x < y$ in radix order, if $|x| < |y|$ or $|x| = |y|$ and $x < y$ in lexicographic order. For example,

$$\{0, 1, 2\}^* = \{\epsilon, 0, 1, 2, 00, 01, 02, 10, \ldots\}$$

**cyclic shift of a string**   One example is `eat, ate, tea`

If $x, y$ are cyclic shifts of each other, we say they are conjugates. Formally, $x, y$ are conjugates if there exists $u, v$ such that $x = uv$ and $y = vu$.

---

[1] need underlying order on $\Sigma$. For example, $a < b < c < \cdots, 0 < 1 < 2 < \cdots$

**Borders**   A word $w$ is **bordered** if it has a proper nonempty prefix that is also a suffix. Otherwise, it's **unbordered**. One example is `entanglement`, whose border is `ent`. Also, we can have overlapping border: `alfalfa`.

## 1.4   Properties of infinite words

**periodicity of infinite words**   Let $x \in \Sigma^+$, finite nonempty words over $\Sigma$. Then we can define

$$x^\omega = xxx \cdots$$

If $z = x^\omega$ for some $x$, we say $z$ is **purely periodic**. If $z = yx^\omega$ for some finite $y$, then $z$ is **ultimately periodic**.

<div style="text-align: right; font-size: 3em;">2</div>

# Combinatorics on words

## 2.1 Equations in words

Suppose we have an equation from number theory,

$$x^2 + xy = y^2 - 1$$

and let's find solution in natural numbers:

$$x = 0 \quad y = 1$$
$$x = 1 \quad y = 2$$
$$x = 3 \quad y = 5$$

Then we can guess the solutions are $x = F_{2n}, y = F_{2n+1}$ for $n \geq 0$.

Now we can consider equations in words: $x, y, z \in \Sigma^+$ (nonempty)

1. $xy = yx$ characterizes commuting words

2. $xy = yz$ characterizes bordered words

For the second equation, one solution would be $x = \mathtt{alf}, y = \mathtt{alfa}, z = \mathtt{lfa}$.

---

**Theorem 2.1**

Suppose $x, y, z \in \Sigma^+$, $xy = yz$ if and only if $\exists u \in \Sigma^+ \mathrm{m}\ v \in \Sigma^*, e \geq 0$ such that

$$x = uv$$
$$z = vu$$
$$y = (uv)^e u = u(vu)^e$$

---

This theorem gives complete characterization to the equation.

> **Proof:**
> $\Leftarrow$ is easy to see:
> $$xy = uv(uv)^e v = (uv)^e uvu = yz$$
> For $\Rightarrow$, we prove by induction on $|y|$.

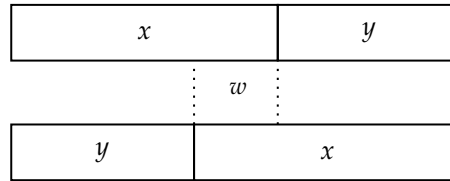Base case $|y| = 1$. Let $y = a$, a single symbol. Then we have

$$xa = az$$

and then we find that $\exists x', z'$ such that $x = ax'$ and $z = z'a$. Then

$$ax'a = az'a$$

So $x' = z'$. Then we can take $u = a, v = x' = z', e = 0$. Then we are done with the base case.
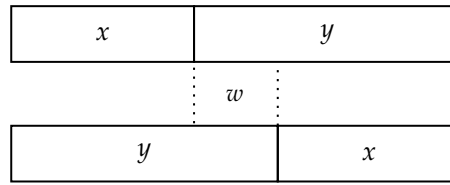
Now induction step. We discuss by cases (imposing length conditions) to break the symmetry.

**Case I**   $|x| \geq |y|$.



We define $w$ (could be empty) as in the picture. Then let $u = y, v = w, e = 0$.

**Case II**   $|x| < |y|$.



We define $w$ as in the picture. We observe that $w \neq \epsilon$, otherwise $|x| = |y|$. Also $x \neq \epsilon, z \neq \epsilon$. Then we observe that

$$y = wz = xw$$

which is our original equation with $w$ playing the role of $y$. In order to apply induction, we need $|w| < |y|$, which is achieved by $x \neq \epsilon$. So induction says $\exists u, v, e, x = uv, z = vu, w = (uv)^e u$. Sub it back in, we get

$$wz = y = (uv)^e uvu = (uv)^{e+1} u$$

$\square$