

Tutorial 1

Computer Language Processing (COMP 4901U)

Monday September 13

- ▶ Recap on lexical analysis
- ▶ Exercise solving in Zoom breakout rooms

Overview

Today we will have a deeper look at tokenization.

What does a tokenizer do?

How can we automatically generate tokenizers?

Overview

Today we will have a deeper look at tokenization.

What does a tokenizer do?

⇒ It transforms a stream of symbols into a stream of tokens.

How can we automatically generate tokenizers?

Overview

Today we will have a deeper look at tokenization.

What does a tokenizer do?

⇒ It transforms a stream of symbols into a stream of tokens.

How can we automatically generate tokenizers?

⇒ We will use regular languages and automata.

Tokenizer

First, let us define a tokenizer as an ordered set of token names and regular expressions

$$\langle Token_1 := e_1, Token_2 := e_2, \dots \rangle$$

where earlier token classes have higher priority than later ones.

E.g.

$$\langle ID := \text{letter} (\text{letter} \mid \text{digit})^*, LE := \leq, LT := <, EQ := = \rangle$$

Recap: Ambiguity in tokenization

Recall that tokenization differs from matching using a single regular expression (say $(e_1 \mid e_2 \mid \dots)^*$).

Rather, the result of tokenizing an input stream of symbols is a stream of tokens. Each token maps to a subsequence of the input stream, and none of the tokens' subsequences overlap.

E.g.

$$\text{i0} \leq \text{size} \xRightarrow{\text{tokenize}} \begin{matrix} ID & LE & ID \\ \text{i0} & \leq & \text{size} \end{matrix}$$

Recap: Ambiguity in tokenization

Recall that tokenization differs from matching using a single regular expression (say $(e_1 \mid e_2 \mid \dots)^*$).

Rather, the result of tokenizing an input stream of symbols is a stream of tokens. Each token maps to a subsequence of the input stream, and none of the tokens' subsequences overlap.

E.g.

$$i0 \leq size \xRightarrow{\text{tokenize}} \begin{matrix} ID & LE & ID \\ i0 & \leq & size \end{matrix}$$

▷ How do we avoid ambiguities?

$$i0 \leq size \xRightarrow{???} \begin{matrix} ID & LT & EQ & ID & ID & ID \\ i0 & < & = & s & iz & e \end{matrix}$$

Tokenization rules

Given an input string w the tokenizer will match tokens on a prefix u of $w = uv$, output the matching token and repeat the process on the remaining string v .

To disambiguate between different possible tokenizations we employ two additional rules:

- ▶ *Longest match*: If we find matching tokens for prefixes of varying lengths, we pick the longer prefix.
- ▶ *Token priority*: If multiple tokens match a prefix of the same length, we pick the token that has higher priority.

A simple tokenization example

Exercise 1

▷ Given the tokenizer

$$\langle T_1 := a(ab)^*, T_2 := b^*(ac)^*, T_3 := cba, T_4 := c^+ \rangle$$

tokenize the following input strings:

c a c c a b a c a c c b a b c

A simple tokenization example

Exercise 1

▷ Given the tokenizer

$$\langle T_1 := a(ab)^*, T_2 := b^*(ac)^*, T_3 := cba, T_4 := c^+ \rangle$$

tokenize the following input strings:

c a c c a b a c a c c b a b c

c c c a a b a b a c c b a b c c b a b a c

A simple tokenization example

Exercise 1

▷ Given the tokenizer

$$\langle T_1 := a(ab)^*, T_2 := b^*(ac)^*, T_3 := cba, T_4 := c^+ \rangle$$

tokenize the following input strings:

$$\underbrace{c}_{T_4} \underbrace{a c}_{T_2} \underbrace{c}_{T_4} \underbrace{a}_{T_1} \underbrace{b a c a c}_{T_2} \underbrace{c b a}_{T_3} \underbrace{b}_{T_2} \underbrace{c}_{T_4}$$

$$\underbrace{c c c}_{T_4} \underbrace{a a b a b}_{T_2} \underbrace{a c}_{T_2} \underbrace{c b a}_{T_3} \underbrace{b}_{T_2} \underbrace{c c}_{T_4} \underbrace{b}_{T_2} \underbrace{a}_{T_4} \underbrace{b a c}_{T_2}$$

A simple tokenization example

Exercise 1

- ▷ Given the tokenizer

$$\langle T_1 := a(ab)^*, T_2 := b^*(ac)^*, T_3 := cba, T_4 := c^+ \rangle$$

tokenize the following input strings:

$$\underbrace{c}_{T_4} \underbrace{a c}_{T_2} \underbrace{c}_{T_4} \underbrace{a}_{T_1} \underbrace{b a c a c}_{T_2} \underbrace{c b a}_{T_3} \underbrace{b}_{T_2} \underbrace{c}_{T_4}$$

$$\underbrace{c c c}_{T_4} \underbrace{a a b a b}_{T_2} \underbrace{a c}_{T_1} \underbrace{c b a}_{T_3} \underbrace{b}_{T_2} \underbrace{c c}_{T_4} \underbrace{b}_{T_2} \underbrace{a}_{T_1} \underbrace{b a c}_{T_3}$$

- ▷ Are there alternative tokenizations if we disregard the longest match rule?