

The Optimal Copypasta Problem

avaxzat

July 16, 2016

1 Introduction

1.1 The Copypasta Problem

Let Σ be a finite alphabet. A *context* is a tuple (b, c, s) where b , c and s are strings over Σ , respectively called the *output buffer*, *clipboard* and *selection buffer*. We let ϵ denote the empty string. A *copypasta* is a string of the form $P_1; \dots; P_n$ where the P_k are any of $\text{Write}(\sigma, i)$, $\text{Select}(i, j)$, Copy or $\text{Paste}(i)$, for all $\sigma \in \Sigma$ and integers i, j . Informally, these statements have the following meaning:

- $\text{Write}(\sigma, i)$. Write symbol σ to position i in the output buffer.
- $\text{Select}(i, j)$. Select the range of characters from index i up to and including index j in the output buffer and copy them to the selection buffer, overwriting its previous contents.
- Copy . Copy the selection buffer to the clipboard.
- $\text{Paste}(i)$. Insert the string in the clipboard at index i in the output buffer.

The optimal copypasta problem can now be stated as follows:

Given a string t over Σ . Find the smallest copypasta P such that $(\epsilon, \epsilon, \epsilon) \vdash P$ rewrites to (t, c, s) according to the rules of Figure 1, where c and s are arbitrary strings over Σ .

1.2 Basic properties

Proposition 1.1. *Let t be a string over Σ containing a total of n unique characters. The length of the optimal copypasta for t will be at least n .*

Proof. Every unique character in t needs to be written to the output buffer at least once, so any copypasta for t must contain at least n Write statements. \square

Proposition 1.2. *The optimal copypasta for t has length at most $|t|$.*

Proof. This is exactly the length of the copypasta that explicitly writes out t using Write statements. \square

2 A simple algorithm

Algorithm 1 shows the pseudocode.

$$\begin{array}{c}
\frac{(b, c, \epsilon) \vdash \text{Write}(\sigma, i)}{(b_1 \dots b_{i-1} \sigma b_i \dots b_n, c, \epsilon)} \text{E-WRITE} \\
\\
\frac{(b, c, s) \vdash \text{Select}(i, j)}{(b, c, b_i \dots b_j)} \text{E-SELECT} \\
\\
\frac{(b, c, s) \vdash \text{Copy}}{(b, s, s)} \text{E-COPY} \\
\\
\frac{(b, c, s) \vdash \text{Paste}(i)}{(b_1 \dots b_{i-1} c b_i \dots b_n, c, s)} \text{E-PASTE} \\
\\
\frac{(b, c, s) \vdash A}{\frac{(b', c', s') \quad (b, c, s) \vdash A; B}{(b', c', s') \vdash B}} \text{E-SEQ}
\end{array}$$

Figure 1: Rewrite rules for the copy-paste problem

Data: a string t , $|t| = n$
Result: a copy-paste for t

```

1  $i \leftarrow 1$ ;
2 while  $i \leq n$  do
3   if  $t[1 : i - 1]$  and  $t[i : n]$  have a common substring of length at least 3 which repeats at index  $i$ 
4     then
5       Let  $s$  be the longest common substring of  $t[1 : i - 1]$  and  $t[i : n]$  which repeats at index  $i$ .
6       Let  $j$  be the index of the first occurrence of  $s$  in  $t[1 : i - 1]$ .
7       Output  $\text{Select}(j, j + |s| - 1)$ ;  $\text{Copy}$ ;  $\text{Paste}(i)$ .
8        $i \leftarrow i + |s|$ ;
9     else
10      Output  $\text{Write}(t_i, i)$ .
11       $i \leftarrow i + 1$ ;
12   end
13 end

```

Algorithm 1: A simple solver for the copy-paste problem