# [SUPERSTRINGWITHEXPANSION] Problem

Salik Lennert Pedersen
Anders Holmgaard Opstrup
Federico Bergamin

## 1 DESCRIPTION OF THE PROBLEM

In this decisional problem we have two different alphabet $\Sigma$ and $\Gamma$, where $\Gamma = \{\gamma_1, ..., \gamma_m\}$ contains $m$ symbol, a string $s$ which is made up of symbols that belongs to the $\Sigma$ alphabet (formally $s \in \Sigma^*$), $k$ strings $t_1, ..., t_k$, which is made up of symbols from both alphabets, and in the end, we have $m$ subsets $R_1, ..., R_m \subseteq \Sigma^*$. These subsets contain symbols of $\Sigma$ alphabet, and as we can notice, they are as many as the symbols in the alphabet $\Gamma$. That's because for each $\gamma_i$ there is a subset $R_i$. We can understand better this relation explaining the core of our problem. The output of the problem will be YES if exits a sequence of words $r_1 \in R_1$, $r_2 \in R_2, ..., r_m \in R_m$, such that for all $1 \leq i \leq k$ the so-called *expansion* $e(t_i)$ is substring of $s$; the expansion of a string $t_i$ consists in the replacement of every symbols $\gamma_j \in \Gamma$ that appears in the string $t_i$ with its expansion, where the expansion of a symbol is defined as follow: $e(\gamma_j) := r_j$. That is, we have to choose for every symbol $\gamma_j \in \Gamma$ a symbol $r_j \in R_j$, and we know that $r_j \in \Sigma^*$. We have to replace the symbol $\gamma_j$ in each string $t_i$ with the symbol $r_j$ that was chosen. We have to do that for each $\gamma \in \Gamma$. In this way, using these replacements, we are transforming our string $t_i$ in its expansion $e(t_i)$.

In the end we will obtain that every expansion-string is made up of symbols of the alphabet $\Sigma$ as our string $s$. Hence, if every new constructed string $e(t_i)$ is a substring of $s$, the answer for our problem will be YES, otherwise NO.

In other words, the answer of our problem will be YES if there is a selection of the $r_i \in R_i$ such that, replacing each symbols $\gamma_i \in \Gamma$ with the respective symbol $r_i$ (we have to remember that for each $\gamma_i$ we have to chose one symbol from the subset $R_i \subseteq \Sigma^*$) in every string $t \in T$, where $|T| = k$, we will have that every new string that we obtain with this replacement is a substring of the string $s$. If a selection of $r_i \in R_i$ like this doesn't exist, the answer will be NO.

## 1.1 Simple examples

With the text of the problem we receive also some problem instances on the alphabets $\Sigma = \{a,b,\ldots,z\}$, $\Gamma = \{A,B,\ldots,Z\}$. The first line of the file contains the number k, which is the number of strigs $t$ we have. The second line contains the string $s$ and the following $k$ lines the strings $t_1,\ldots,t_k$. Finally, the last lines (at most 26, the number of letters in the alphabet) start with a letter $\gamma_j \in \Gamma$ followed by a colon and the contents of the set $R_j$ belonging to the letter, where the elements of the set are separated by commas. Hence, we can clearly see that for each symbols $\gamma_i \in \Gamma$, in this case letter, there is a subset $R_i$ and in this subset we have to chose one symbol.

For example:

```
4
abdde
ABD
DDE
AAB
ABd
A:a,b,c,d,e,f,dd
B:a,b,c,d,e,f,dd
C:a,b,c,d,e,f,dd
D:a,b,c,d,e,f,dd
E:aa,bd,c,d,e
```

In this example we can observe that there is NO possible selection of symbols $r_i$ such that every string $t$ would be a substring of s=abdde. In fact if we have a look to $t_2$=DDE and $t_3$=AAB, we can notice that they have the same structure, and in the same time in our string $s$ the only letter that is repeated twice is the d. So we conclude that A=d and D=d. Moreover from $t_4$=ABd we see that, since $t_4$ will be a substring, or B=d or B=b. If we try to choose our selection of $r_i$ following these rules, we can noticed that there is no possible solution, such that at the same time, all $e(t_i)$ are substring of $s$.

## 2 SuperStringWithExpansion is in $\mathcal{NP}$

We have to show and prove that our problem is in $\mathcal{NP}$. To do that we use the so called *guess-verify algorithms* proof, so we have to start designing a deterministic algorithm A which takes as input a problem instance X and a random sequence R.
**Important:** In this case, since we already use R to indicates the subset of the problem, we are going to call the random sequence with the capital letter **G**.

*Proof.* We have to show that there is a polynomial $p$ and a randomized $p$-bounded algorithm $A$ which satisfies the condition of the class $\mathcal{NP}$.

1. Let $\Sigma$ be an alphabet and $\Gamma$ another alphabet. We know that $|\Gamma| = m$ (cardinality of the alphabet $\Gamma$ is $m$). Let $s$ be a string made up of symbols of $\Sigma^*$ and let $t_1, ..., t_k$ be $k$ strings made up of symbols of $(\Sigma \cup \Gamma)^*$. Then let $R_1, ..., R_m$ be subsets that contains symbols of $\Sigma^*$. So the first part of the input is:

$$X = (\Sigma, \Gamma, s, t_1, ..., t_k, R_1, ..., R_m)$$

We could simplify this symbolism, considering the set $T = \{t_1, ..., t_k\}$ and $R = \{R_1, ..., R_m\}$. In this way we obtain that the first part of the input for our algorithm A now is:

$$X = (\Sigma, \Gamma, s, T, R)$$

   a) The random sequence $G$ consists of **integers** in the range $\{1, ..., n\}$, where $n = max(|R_i|)$ for $1 \le i \le m$.

   b) Now on input $((\Sigma, \Gamma, s, T, R), \mathbf{G})$, algorithm $A$ checks whether $G$ contains at least $m$ integers, and for every integer it has to check whether:

$$int_j \le |R_j|$$

   which means that $A$ has to check if the integer contains in the position $j$, where $1 \le j \le m$, is less or equal than the cardinality of the subset $R_j$. The reason of that will be clear immediately.
   If $G$ is shorter, or the first $m$ integer don't satisfy the condition above, $A$ will return NO. Otherwise $A$ uses the first $m$ integers $g_1, ..., g_m$ of $G$ to select a symbols in every subsets $R_1, ..., R_m$. In this way:

$$r_{g_i} \in R_i$$

   That is, the integer $g_i$, which is in the position $i$ inside the random string $G$, indicates the position of the symbol that we have to choose inside the subset $R_i$. That's why it has to satisfy the condition above, in fact if the integer $g_i$ is greater than the cardinality of $R_i$ we are not able to select a symbol.

   c) Now we have to choose an element for every subset $R_i$, where $1 \le i \le m$, and then using the expansion for the symbols of the alphabet $\Gamma$, we are able to "link" a letter $\gamma_i \in \Gamma$ with a symbol in $R_i$. (That's because for each letter $\gamma_i$ we have to choose one and only one symbol in the subset $R_i$, the $i$ is the same, because the subset is linked to the letter).

$$\gamma_i \in \Gamma \implies e(\gamma_i) = r_j \in R_i$$

   where $j = g_i$ of the string $G$

   Now we could substitute every letters $\gamma_i \in \Gamma$, $1 \le i \le m$, which are in the $k$ strings $t$ with their expansion that we obtained thanks to the formula above. After that, we will have $k$ strings $e(t_i)$ (which are the expansion of our strings) made up of only symbols in $\Sigma^*$.

Then we have to check if every string that we obtain $e(t_l)$, $1 \le l \le k$, is a substring of our string $s$. If every string is a substring the algorithm A returns YES, otherwise it returns NO.

2. Now we have to show that the two conditions of the class $\mathcal{NP}$ are met:

   a) Let us first assume that the true answer is YES, so there is a selection of symbols $r_1, ..., r_m$ from the subsets $R_1, ..., R_m$ such that, using the expansion for the symbols $\gamma \in \Gamma$ first, and then for the $k$ strings $t$, we will obtain that every string $e(t)$ (the expanded one) is a substring of $s$. So the symbols $r_1, ..., r_m$ is our solution. That means, that our subsets $R_1, ..., R_m$ contains one symbol each to expand $\gamma_1, ..., \gamma_m \in \Gamma$ in the right way, and we have to choose exactly that symbol in every subset $R_i$. Hence, there will be a string $G = g_1...g_m$ (which is the random string and contains the position of the element we will select) such that we will select the right symbol inside every subset $R$, so we are able to obtain that every string $e(t_i)$ is a substring of $s$. In this case if $G$ is given to $A$, $A$ will correctly return YES. Hence, there is a string of length at most $m$ which lead to a correct answer YES. To calculate the probability to obtain this string we try to see how to built this string that satisfy our problem: we show this probability in the two cases: the first, where all the subsets have exactly $n$ elements (which is useful to understand how change the probability in regard to the number of elements of the subsets and the number of the subsets) and the second, in which we consider every cardinality of the subset. We have to choose exactly one element in each subset, so our $g_i$ has to be the exactly position of the right symbols. In the first case for each $g_i$ we have to choose 1 position between $n$, ans so we have:

   $$G = g_1 g_2 ... g_{m-1} g_m$$

   $$\mathbf{P}[\text{G is the right string}] = \underbrace{\left(\frac{1}{n}\right) \times ... \times \left(\frac{1}{n}\right)}_{m \text{ times}} = \left(\tfrac{1}{n}\right)^m$$

   Otherwise, if we want to use subset with no fixed cardinality, we have this probability:

   $$G = g_1 g_2 ... g_{m-1} g_m$$

   $$\mathbf{P}[\text{G is the right string}] = \left(\tfrac{1}{|R_1|}\right) \times ... \times \left(\tfrac{1}{|R_m|}\right)$$

   We could see that this probability is small, especially if $n$ and $m$ are bigger, but it's positive. Hence, the first condition is satisfied.

   b) Now, we assume that the true answer is NO, so there is no possible selection of symbols $r_1, ..., r_m$ from $R_1, ..., R_m$ such that every expanded-strings $e(t_i)$ are substring of $s$. In other words, for each random string $G$ we pass to $A$, there is no

possible sequence of position $g_1...g_m$ such that let us select the right symbols to obtain that every expanded-strings $e(t_i)$ are substring of $s$. Hence, regardless the random string $G$, the algorithm $A$ will always answer NO in this case, so also the second condition is satisfied.

3. Now we have to show that the running time of the algorithm $A$ that we've built is $p$-bounded for some polynomial $p$. We could see that our algorithm have to:

   - check if $G$ contains at least $m$ integer: *O(m)*;

   - for each integer check if the condition is satisfied: so for each integer has to calculate the cardinality of the subset $R_i$ related to the position $g_i$ and compare to the integer value. If we are going to use arrays we have to compare an integer to the `length` of the array. So it will be *O(1)* for each element in $G$. For $m$ element it will be *O(m)*;

   - select the exact element in $R_i$ using the position $g_i$ for every subset. If we will work with arrays we have to enter one cell, so *O(1)*, but for all the letters of the alphabet $\Gamma$ is *O(m × 1) = O(m)*;

   - read all the $k$ strings and substitute all the $\gamma$ with their expansions. The worst case is when the length of the strings are almost $n$, so we have to read all the $k$ strings and substitute. If we think that a substitution is *O(1)*, then it will be *O(k × n)*;

   - in the end check if the all $k$ expanded-strings are substring of $s$. Using the simplest algorithm, we have that it will be *O(n+q)*, where $n$ is the length of the string $s$ and $q$ the length of strings $t$. Hence, for all $k$ strings $t$ it will be *O(k(n+q))*

Therefore, we can see that our algorithm give the result in time:

$$m + m + m + kn + k(n + q) = 3m + kn + k(n + q)$$

In the worst case, we could have $n$ symbols in $\Gamma$, and $n$ strings $t$, so the algorithm will be *O(n$^2$)*, which is polynomial.

$\square$

# 3  SuperStringWithExpansion is $\mathcal{NP}$−complete

In this section we are going to prove that our SuperStringWithExpansion problem is $\mathcal{NP}$−complete.

*Proof.* We could reduce from 1-in-3 Satisfiability, which is a $\mathcal{NP}$−complete problem, as follows. Firstly we have to define two alphabets for our problem:

$$\Sigma = \{0, 1, 2\} \text{ and } \Gamma = \{x_1, ..., x_n\}$$

Then for each clause of the form $c = x_1 \vee x_2 \vee x_3$ in the instance **X** of 1-in-3 Satisfiability, we generate a new string $t$ that contains the literals of $c$ in this way:

$$t = x_1 x_2 x_3$$

and as we could see, the string $t$ contains only literals $\in \Gamma^*$.

Then for each letters of the alphabet $\Gamma$ we have to define the subsets $R_1, ..., R_n$, and we will do it in this way:

$$R_1 = R_2 = ... = R_n = \{0, 1\}$$
$$\text{so we obtain:}$$
$$x_1 : R_1 = \{0, 1\}$$
$$x_2 : R_2 = \{0, 1\}$$
$$.$$
$$.$$
$$.$$
$$x_n : R_n = \{0, 1\}$$

In the end we have to define the string $s$, and we decide to construct $s$ as:

$$s = 00120102100$$

Essentially,for each clause in the instance **X** we create a string made up of the literals of the clause and for all the variables $\{x_1, ..., x_n\}$ we create a subset that contains two elements. In the end we create a string $s$ made up of letters of $\Sigma$. Hence, our transformation $T$ constructs $k$ strings $t$ of 3 letters, a string $s$ of 11 letters, and $n$ subset of 2 elements. Therefore we could see that the running time of our transformation is linear in the input size, in fact for for all $k$ clauses of length 3 of the instance **X** there will be $k$ strings $t$ of length 3, and for each boolean variables it creates a subset with two element, so $n$ subset in total.

Now suppose that the original set of clauses of the instance **X** is satisfiable, that means that there is an assignment of truth values to the boolean variables $\{x_1, ..., x_n\}$ such that every clause of length 3 in **X** is satisfied and contains only one literal true. We claim that exist also a selection of the element in the subsets $R_1, ..., R_n$ such that, if we expand the letters in $\Gamma$ and substitute them in the strings $t$, all of these are substring of the string $s$. We could expand each letter $x_1, ..., x_n$ using the value of assignment used in the clauses of the instance **X**. In this way, substituting the expanded letters in the strings $t$, we could see all of strings are in this form:

$$t = 001 \text{ or } t = 010 \text{ or } t = 100$$

That's because each clauses with this truth assignment contains only one literal true. We could see, also that all the string will be substring of the string $s$. Hence, also the SUPERSTRINGWITH-EXPANSION problem will return YES.

On the other hand, let us assume that there is a selection of the element in $R_1, ..., R_n$ such that, if we substitute the letters in $\Gamma$ with their expansions in the strings $t$, we obtain that all strings $t$ are substring of $s$. That means that we've assigned a value (or 0 or 1) to all the variables $\{x_1, ..., x_n\}$ such that each string $t$ should have a structure like 001 or 010 or 100. Otherwise it cannot be a substring of $s$ and so the answer will be NO. Since this fact, we could see that we have a truth assignment to the value such that all the clauses contains only one literal equal to 1. Using the same truth values also for the instance **X**, we will obtain that each clause is satisfied and contains only one literal true. Hence, also for the 1-IN-3 SATISFIABILITY will

return YES.

That shows that we can reduce from 1-IN-3 SATISFIABILITY, then our problem SUPERSTRING-WITHEXPANSION is $\mathcal{NP}$−complete. □

# 4 IMPLEMENTATION OF THE DECODER

We have to implement a decoder that read a SWE files, which contains an instance of our SUPERSTRINGWITHEXPANSION problems (where $\Sigma$ is the lower-case alphabet and $\Gamma$ the upper-case one), that check if these are written in a correctly way, and in particular rejected the files in which the subset $R_i$ are not specified in a proper way.
Our program has to read line to line the SWE file and check if everything is correct. Hence, we have to check these things:

- the first line has to be a number, which represents the number of the strings $t$;

- the second line has to be a string made up of only lower-case characters, we cannot accept numbers, symbols and upper-case letters;

- each string $t$ has to be a string which could make up of both upper-case and lower-case letters, but we cannot accept numbers and symbols. Moreover, the length of the string $t$ has to be less or equal to the length of the string $s$;

- each subset has to begin with a capital letter followed by the column and then there has to be lower-case letter separated by the comma. We have to check if this "grammar" is respected, hence, we cannot accept symbols, blank spaces, upper-case letter except at the beginning. Moreover we have to check that there exists a subset for all the upper-case letter that appears on the strings $t$.

The program code (in Java) is attached in the project folder.