

# Finding Diverse Strings and Longest Common Subsequences in a Graph

Luca Lombardo

Seminar for the course of Bioinformatics

# Structure of the Presentation

- 1 Introduction to the problems
  - Diverse Strings Problems
  - $\Sigma$ -DAG for LCSs
- 2 Exact Algorithms for Bounded Number of Diverse Strings
  - MAX-MIN DIVERSE STRING SET problem
  - MAX-SUM DIVERSE STRING SET problem
- 3 Approximation Algorithms for Unbounded Number of Diverse Strings
  - Approximation Algorithm for MAX-MIN DIVERSE STRING SET problem
  - MAX-SUM FARTHEST  $r$ -STRING problem
- 4 FPT Algorithms for Bounded Number and Length of Diverse Strings
- 5 Hardness Results
  - Hardness of Diverse String Set for Unbounded  $K$
  - Hardness of Diverse String LCSs for Unbounded  $K$

# Longest Common Subsequence (LCS)

## Definition

Given a set of  $m$  strings  $S = \{S_1, S_2, \dots, S_m\}$ , a **common subsequence** (CS) is a sequence that appears in all  $m$  strings. A **longest common subsequence** LCS is a common subsequence of maximum length. We denote the set of all LCSs of  $S$  as  $LCS(S)$ .

## Goal

The goal is to find a diverse set of solutions to the LCS problem under the Hamming distance.

## Definition

Given two strings  $X, Y \in \Sigma^r$ , the **Hamming distance** between  $X$  and  $Y$ , denoted with  $d_H(X, Y)$ , is the number of positions at which the corresponding symbols differ.

# A simple example

Add here the example of the paper

# Efficient methods for finding a diverse set of solutions

More formally, let's consider the following two diversity measures for a multiset  $\mathcal{X} = \{X_1, X_2, \dots, X_K\} \subseteq \Sigma^r$  of solutions, allowing repetitions:

$$D_{d_H}^{\text{sum}}(\mathcal{X}) = \sum_{1 \leq i < j \leq K} d_H(X_i, X_j) \quad \text{MAX-SUM DIVERSITY} \quad (1)$$

$$D_{d_H}^{\text{min}}(\mathcal{X}) = \min_{i < j} d_H(X_i, X_j) \quad \text{MAX-MIN DIVERSITY} \quad (2)$$

## Notation

A subset  $\mathcal{X} \subseteq \Sigma^r$  is  $\Delta$ -diverse w.r.t.  $D_{d_H}^\tau$  if  $D_{d_H}^\tau(\mathcal{X}) \geq \Delta$  for some  $\Delta \geq 0$ .

Where for  $\tau \in \{\text{sum}, \text{min}\}$ ,  $D_{d_H}^\tau$  denotes one of the two diversity measures.

## Two problems

### Problem 1: DIVERSE LCSs WITH DIVERSITY MEASURE $D_{d_H}^\tau$

*Input:* A set  $S = \{S_1, S_2, \dots, S_m\}$  of  $m \geq 2$  strings over  $\Sigma$ , an integer  $K \geq 1$  and  $\Delta \geq 0$ .

*Question:* Is there some set  $\mathcal{X} \subseteq LCS(S)$  such that  $|\mathcal{X}| = K$  and  $D_{d_H}^\tau(\mathcal{X}) \geq \Delta$ ?

### Problem 2: DIVERSE STRING SET

*Input:*  $K, r, \Delta \in \mathbb{Z}$  and a  $\Sigma$ -DAG  $G$  for a set  $L(G) \subseteq \Sigma^r$  of strings.

*Question:* Decide if there exists some subset  $\mathcal{X} \subseteq L(G)$  such that  $|\mathcal{X}| = K$  and  $D_{d_H}^\tau(\mathcal{X}) \geq \Delta$ .

# $\Sigma$ -Labeled Directed Acyclic Graphs ( $\Sigma$ -DAGs)

## Definitions

- **Alphabet ( $\Sigma$ ):** Set of symbols.
- **String Set (Language):**  $L = \{X_1, X_2, \dots, X_n\} \subseteq \Sigma^*$ , with:
  - ▶ **Total Length:**  $||L|| = \sum_{X \in L} |X|$
  - ▶ **Max Length:**  $\text{maxLen}(L) = \max_{X \in L} |X|$
- **r-String:** Any string  $X$  where  $|X| = r$

## $\Sigma$ -DAG Structure

A graph  $G = (V, E, s, t)$  with:

- $V$ : vertices,  $E$ : labeled edges  $(v, c, w)$  with  $c \in \Sigma$
- **Source**  $s$  and **Sink**  $t$  such that paths exist from  $s$  to all vertices
- **Size:**  $\text{size}(G)$ , the number of its labeled edges

# $\Sigma$ -Labeled Directed Acyclic Graphs ( $\Sigma$ -DAGs)

Any path  $P = (e_1, e_2, \dots, e_n)$  of outgoing edges *spells out* a string  $\text{str}(P) = c_1 c_2 \dots c_n \in \Sigma^n$  where  $c_i$  is the label of edge  $e_i$ .

## Language Representation

A  $\Sigma$ -DAG represents  $L(G) \subset \Sigma^*$ : all strings spelled from paths  $s \rightarrow t$ .  
Equivalent to an NFA over  $\Sigma$  with initial  $s$ , final  $t$ , and no  $\epsilon$ -edges.

## Remarks

- For any set  $L$  of strings, a  $\Sigma$ -DAG  $G$  exists s.t.  $L(G) = L$  and  $\text{size}(G) \leq ||L||$ . Construction time:  $O(||L|| \log |\Sigma|)$ .
- If  $G$  represents a set  $L$  of  $r$ -strings ( $L \subseteq \Sigma^r, r \geq 0$ ), all paths  $s \rightarrow v$  have the same length  $d \leq r$ .



# $\Sigma$ -DAGs for LCSs and Diverse String Sets

## Lemma ( $\Sigma$ -DAG for LCSs)

For any constant  $m \geq 1$  and set  $S = \{S_1, \dots, S_m\} \subseteq \Sigma^*$  of  $m$  strings, there exists a  $\Sigma$ -DAG  $G$  of polynomial size in  $\ell := \maxlen(S)$  such that  $L(G) = \text{LCS}(S)$  and can be computed in polynomial time /

## Consequence of the Lemma

- **If** MAX-MIN (or MAX-SUM) DIVERSE STRING SET solvable in  $f(M, K, r, \Delta)$ ,
- **Then** MAX-MIN (or MAX-SUM) DIVERSE LCSs on  $S \subseteq \Sigma^r$  solvable in  $O(|\Sigma| \cdot \ell^m + f(\ell^m, K, r, \Delta))$  time.

where  $\ell = \maxlen(S)$

# Algorithm for MAX-MIN DIVERSE STRING SET problem

## Dynamic Programming Approach

**Pattern of Path Tuple:** For each  $d \leq r$  and  $K$ -tuple of length- $d$  paths  $P = (P_1, \dots, P_K)$ , define pattern  $\text{Pattern}(P) = (\mathbf{w}, \mathbf{Z})$ , where:

- $\mathbf{w} = (w_1, \dots, w_K)$ :  $K$ -tuple of vertices representing endpoints of  $P_i$  paths.
- $\mathbf{Z} = (Z_{i,j})$ : Upper triangular matrix of Hamming distances, where  $Z_{i,j} = \min\{\Delta, d_H(\text{str}(P_i), \text{str}(P_j))\}$ .

## DP Table of Weights

$$\text{Weights} : V^K \times (\Delta \cup \{0\})^{K \times K} \rightarrow \{0, 1\} \quad (3)$$

Boolean matrix where  $\text{Weights}(\mathbf{w}, \mathbf{Z}) = 1$  iff  $(\mathbf{w}, \mathbf{Z})$  matches pattern  $\text{Pattern}(P)$  for some  $K$ -tuple of paths of length  $d$  from  $s$  to  $\mathbf{w}$ .

# Algorithm for MAX-MIN DIVERSE STRING SET problem

---

## Algorithm 1

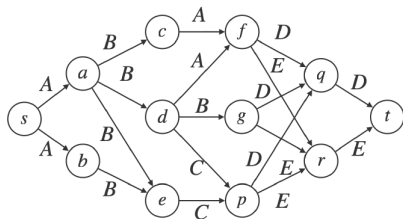
---

```
1: Set  $\text{Weights}(\mathbf{s}, Z) = 0$  for all  $Z \in (\Delta \cup \{0\})^{K \times K}$  and  $\text{Weights}(\mathbf{s}, \mathbf{0}) \leftarrow 1$ 
2: for  $d \leftarrow 1, \dots, r$  do
3:   for  $\mathbf{v} \leftarrow (v_1, \dots, v_K) \in (V_d)^K$  do
4:     for  $(v_1, c_1, w_1) \in E^+(v_1), \dots, (v_K, c_K, w_K) \in E^+(v_K)$  do
5:       Set  $\mathbf{w} = (w_1, \dots, w_K)$ 
6:       for  $U \in (\Delta \cup \{0\})^{K \times K}$  such that  $\text{Weights}(\mathbf{v}, U) = 1$  do
7:         Set  $Z = (Z_{i,j})_{i < j}$  with  $Z_{i,j} \leftarrow \min\{\Delta, U_{i,j} + \mathbb{1}\{c_i \neq c_j\}\}$   $\forall i, j \in K$ 
8:         Set  $\text{Weights}(\mathbf{w}, Z) \leftarrow 1$  ▷ Update
9:       end for
10:    end for
11:  end for
12: end for
13: Answer Yes if  $\text{Weights}(\mathbf{t}, Z) = 1$  and  $D_{d_H}^{\min}(Z) \geq \Delta$  for some  $Z$ , else No
```

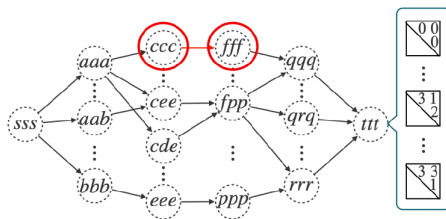
---

For any  $K \geq 1$  and  $\Delta \geq 0$ , it solves the MAX-MIN DIVERSE STRING SET problem in  $O(\Delta^{K^2} K^2 M^K (\log |V| + \log \Delta))$  time and space when an input string set  $L$  is represented by a  $\Sigma$ -DAG  $G$  with  $\text{size}(G) = M$

# Example Run of the Algorithm



(a) An input  $\Sigma$ -DAG  $G_1$  for  $LCS(X_1, Y_1)$



(b) Example run of Algorithm 1 on  $G_1$

**Figure:** (a) An input  $\Sigma$ -DAG  $G_1$  over  $\Sigma = \{A, B, C, D, E\}$  for the set of all longest common subsequences of two strings  $X_1 = ABABCDDEE$  and  $Y_1 = ABCBAEEDD$  and (b) an example run of Algorithm 1 based on dynamic programming with  $K = 3$  on an input  $G_1$

# Algorithm for MAX-SUM DIVERSE STRING SET problem

# Algorithm for MAX-SUM DIVERSE STRING SET problem

# MAX-MIN DIVERSE STRING SET problem

# MAX-SUM FARTHEST r-STRING problem



# FPT Algorithms for Bounded Number and Length of Diverse Strings

# FPT Algorithms for Bounded Number and Length of Diverse Strings

# Hardness of Diverse String Set for Unbounded $K$

# Hardness of Diverse String LCSs for Unbounded $K$