







Article

Solving the Longest Common Subsequence Problem Concerning Non-Uniform Distributions of Letters in Input Strings

Bojan Nikolic ^{1,†} , Aleksandar Kartelj ^{2,†} , Marko Djukanovic ^{1,3,*,†} , Milana Grbic ^{1,†} , Christian Blum ^{4,†} 
and Günther Raidl ^{3,†} 

- ¹ Faculty of Natural Science and Mathematics, University of Banja Luka, 78 000 Banja Luka, Bosnia and Herzegovina; bojan.nikolic@pmf.unibl.org (B.N.); milana.grbic@pmf.unibl.org (M.G.)
² Faculty of Mathematics, University of Belgrade, 105104 Belgrade, Serbia; aleksandar.kartelj@gmail.com
³ Institute of Logic and Computation, Faculty of Informatics, TU Wien, 1040 Vienna, Austria; raidl@ac.tuwien.ac.at
⁴ Artificial Intelligence Research Institute (IIIA-CSIC), Campus UAB, 08193 Bellaterra, Spain; christian.blum@iiia.csic.es
* Correspondence: djukanovic@ac.tuwien.ac.at; Tel.: +387-65-699-683
† These authors contributed equally to this work.

Abstract: The longest common subsequence (LCS) problem is a prominent \mathcal{NP} -hard optimization problem where, given an arbitrary set of input strings, the aim is to find a longest subsequence, which is common to all input strings. This problem has a variety of applications in bioinformatics, molecular biology and file plagiarism checking, among others. All previous approaches from the literature are dedicated to solving LCS instances sampled from uniform or near-to-uniform probability distributions of letters in the input strings. In this paper, we introduce an approach that is able to effectively deal with more general cases, where the occurrence of letters in the input strings follows a non-uniform distribution such as a multinomial distribution. The proposed approach makes use of a time-restricted beam search, guided by a novel heuristic named GMPSUM. This heuristic combines two complementary scoring functions in the form of a convex combination. Furthermore, apart from the close-to-uniform benchmark sets from the related literature, we introduce three new benchmark sets that differ in terms of their statistical properties. One of these sets concerns a case study in the context of text analysis. We provide a comprehensive empirical evaluation in two distinctive settings: (1) short-time execution with fixed beam size in order to evaluate the guidance abilities of the compared search heuristics; and (2) long-time executions with fixed target duration times in order to obtain high-quality solutions. In both settings, the newly proposed approach performs comparably to state-of-the-art techniques in the context of close-to-uniform instances and outperforms state-of-the-art approaches for non-uniform instances.

Keywords: longest common subsequence problem; multi-nomial distribution; probability-based search guidance



Citation: Nikolic, B.; Kartelj, A.; Djukanovic, M.; Grbic, M.; Blum, C.; Raidl, G. Solving the Longest Common Subsequence Problem on Non-Uniform Distributions. *Mathematics* **2021**, *9*, 1515. <https://doi.org/10.3390/math9131515>

Academic Editor: Alfredo Milani

Received: 21 May 2021
Accepted: 25 June 2021
Published: 29 June 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

In the field of bioinformatics, strings are commonly used to model sequences such as DNA, RNA, and protein molecules or even time series. Strings represent fundamental data structures in many programming languages. Formally, a string s is a finite sequence of $|s|$ letters over (usually) a finite alphabet Σ . A subsequence of a string s is any sequence obtained by removing arbitrary letters from s . Similarities among several strings can be determined by considering common subsequences, which may serve for deriving relationships and possibly to distill different aspects of the of input strings, such as mutations. More specifically, one such measure of similarity can be defined as follows. Given a set of m input strings $S = \{s_1, \dots, s_m\}$, the *longest common subsequence* (LCS) problem [1] aims

at finding a subsequence of maximum length that is common for all strings from the set of input strings S . The length of the LCS for two or more input strings is a widely used measure in computational biology [2], file plagiarism check, data compression [3,4], text editing [5], detecting road intersections from GPS traces [6], file comparison (e.g., in the Unix command *diff*) [7] and revision control systems such as GIT. For a fixed m , polynomial algorithms based on dynamic programming (DP) are known [8] in the literature. These dynamic programming approaches run in $O(n^m)$ time, where n denotes the length of the longest input string. Unfortunately, these approaches quickly become impractical when m and n get large. For an arbitrary large number of input strings, the LCS problem is \mathcal{NP} -hard [1]. In practice, heuristic techniques are typically used for larger m and n . Constructive heuristics, such as the Expansion algorithm and the Best-Next heuristic [9,10], appeared first in the literature to tackle the LCS problem. Significantly better solutions are obtained by more advanced metaheuristic approaches. Most of these are based on *Beam Search* (BS), see for example, [11–15]. These approaches differ in various important aspects, which include the heuristic guidance, the branching scheme and the filtering mechanisms.

Djukanovic et al. (2019), in [16], proposed a generalized BS framework for the LCS problem with the purpose of unifying all previous BS-based approaches from the literature. By respective parametrization, each of the previously introduced BS-based approaches from the literature could be expressed, which also enabled a more direct comparison of all of them. Moreover, a heuristic guidance that approximates the expected length of an LCS on uniform random strings was proposed. This way, a new state-of-the-art BS variant that leads on most of the existing random and quasi-random benchmark instances from the literature was obtained.

Concerning exact approaches for the LCS problem, an integer linear programming model was considered in [17]. It turned out not to be competitive enough as it was not applicable to most of the commonly used benchmark instances from the literature. This was primarily due to the model size—too many binary variables and a huge number of constraints are needed even for small-sized problem instances. Dynamic programming approaches also quickly run out of memory for small-to-middle sized benchmark instances or typically return only weak solutions, if any. Chen et al. (2016) [18] proposed a parallel FAST_LCS search algorithm that mitigated some of the runtime weaknesses. Wang et al. (2011) in [14] proposed another parallel algorithm called QUICK-DP, which is based on the dominant point approach and employs a quick divide-and-conquer technique to compute the dominant points. One more parallel and space efficient algorithm based on a graph model, called the LEVELED-DAG, was introduced by Peng and Wang (2017) [19]. Li et al. (2016) in [20] suggested the TOP_MLCS algorithm, which is based on a directed acyclic layered-graph model (called the irredundant common subsequence graph) and parallel topological sorting strategies used to filter out paths representing suboptimal solutions. A more enhanced version of the latter approach that utilizes Branch-and-Bound is proposed by Wang et al. (2021) in [21]. An approximate algorithm based on the topological sorting strategy was developed and applied to COVID-19 data by Li et al. (2020) in [22]. A specific Path Recorder Directed Acyclic Graph (PRDAG) model was designed and used for developing the Path Recorder Algorithm (PRA) by Wei et al. (2020) in [23]. Recently, Djukanovic et al. (2020) in [24] proposed an A^* search that is able to outperform TOP_MLCS and previous exact approaches in terms of memory usage and the number of instances solved to optimality. Nevertheless, the applicability of this exact A^* search approach is still limited to small-sized instances. The same holds for the other exact techniques from the literature. In [24], the A^* search also served as a basis for a hybrid anytime algorithm, which can be stopped at almost any time and can then be expected to yield a reasonable heuristic solution. In this approach, classical A^* search iterations are intertwined with iterations of *Anytime column search* [25]. This hybrid was able to outperform some other state-of-the-art anytime algorithms from the literature, such as PRO_MLCS [26] and *Anytime Pack Search* [27], in terms of solution quality. To the best of our knowledge,

this hybrid approach is still the leading method for the LCS problem on a wide range of benchmark sets from the literature.

The methods so-far proposed in the literature were primarily tested on independent random and quasi-random strings, where the number or occurrences of letters in each string is similar for each letter. In fact, we are aware of just one benchmark set with different distributions (BB, see Section 4), where the input strings are constructed in such a way so that they exhibit high similarity, but still the letters' frequencies are similar. In practical applications, this assumption of uniform or close-to-uniform distribution of letters does not need to hold. Some letters may occur substantially more frequently than others. For example, if we are concerned with finding motifs in sentences of any spoken language, each letter has its characteristic frequency [28]. Text in natural languages can be modeled by a multinomial distribution over the letters. The required level of model adaptation can vary depending on the distribution assumptions such as letter dependence of a particular language. Letter frequencies in a language can also differ depending on text types (e.g., poetry, fiction, scientific documents, business documents). For example, it is interesting that the letter 'E' is the most frequent letter in both English (12.702%) [28] and German (17.40%) [29], but only the second most common letter in Russian [30]. Moreover, the letter 'N' is very frequent in German (9.78%), but not so common in English (6.749%) and Russian (6.8%). Motivated by this consideration, we develop a new BS-based algorithm exhibiting an improved performance for instances with different string distributions.

Unlike the BS proposed by Djukanovic et al. (2019) in [16], our BS variant applies a novel guidance heuristic at its core as a replacement of the so far leading approximate expected length calculation. In addition to an improved performance in the context of non-uniform instances, the novel heuristic is easier to implement than the approximate expected length calculation (which required a Taylor series expansion and an efficiently implemented divide-and-conquer approach) and there are no issues with numerical stability. Our new BS version is time restricted in the sense that its running time approximately fits a desired target time limit. Moreover, it obtains high-quality results without the need for extensive parameters tuning.

The main contributions of this article can be summarized as follows:

- We propose a novel search guidance for a BS which performs competitively on the standard LCS benchmark sets known from the literature and, in some cases, even produces new state-of-the-art results.
- We introduce two new LCS benchmark sets based on multinomial distributions, whose main property is that letters occur with different frequencies. The proposed new BS variant excels in these instances in comparison with previous solution approaches.
- A new time-restricted BS version is described. It automatically adapts the beam width over BS levels with regard to given time restrictions, such that the overall running time of BS approximately fits a desired target time limit. A tuning of the beam width to achieve comparable running times among different algorithms is hereby avoided.

In the following, we introducing some commonly used notation before giving an overview of the remainder of this article.

1.1. Preliminaries

By S , we always refer to the set of m input strings, that is, $S = \{s_1, \dots, s_m\}$, $m \geq 1$. The length of a string s is denoted by $|s|$, and its i -th letter, $i \in \{1, \dots, |s|\}$, is referred to by $s[i]$. Let n refer to the length of a longest string and n_{\min} to the length of a shortest string in S . A continuous subsequence (substring) of string s that starts with the letter at index i and ends with the letter at index j is denoted by $s[i, j]$; if $i > j$, this refers to the empty string ε . The number of occurrences of a letter $a \in \Sigma$ in string s is denoted by $|s|_a$. For a subset of the alphabet $A \subseteq \Sigma$, the number of appearances of each letter from A in s is denoted by $|s|_A$. For an m -dimensional integer vector $\vec{\theta} \in \mathbb{N}^m$ and the set of strings S , we define the set of suffix-strings $S[\vec{\theta}] = \{s_1[\theta_1, |s_1|], \dots, s_m[\theta_m, |s_m|]\}$, which induces a respective LCS subproblem. For each letter $a \in \Sigma$, the position of the first occurrence of a in $s_i[\vec{\theta}_i, |s_i|]$ is

denoted by $\vec{\theta}_{i,a}$, $i = 1, \dots, m$. Last but not least, if a string s is a subsequence of a given string r , we write $s \prec r$.

1.2. Overview

This article is organized as follows: Section 2 provides theoretical aspects concerning the calculation of the probability that a given string is a subsequence of a random string chosen from a multinomial distribution. Section 3 describes the BS framework for solving the LCS problem, as well as the novel heuristic guidance. Moreover, the time-restricted BS variant is also proposed. In Section 4, a comprehensive experimental study and comparison is conducted. Section 5 extends the experiments by considering instances derived from a textual corpus. Finally, Section 6 draws conclusions and outlines interesting future work.

2. Theoretical Aspects of Different String Distributions

Most papers in the literature are dedicated to the development and improvement of methods for finding an LCS of instances on strings that come from a uniform distribution. In our work, we propose new methods for the more general case, where strings are assumed to come from a multinomial distribution $MN(p_1, \dots, p_\eta)$ of strings. More precisely, for an alphabet $\Sigma = \{a_1, \dots, a_\eta\}$, $\eta > 1$, as a sample space for the letter of the strings, a multinomial distribution $MN(p_1, \dots, p_\eta)$ is determined by specifying a (real) number p_i for each letter a_i such that p_i represents the probability of seeing letter a_i and $\sum_{i=1}^{\eta} p_i = 1$. Note that the uniform distribution is a special case of the multinomial distribution $MN(p_1, \dots, p_\eta)$, with $p_1 = \dots = p_n = \frac{1}{\eta}$.

Assuming that the selection of each letter in a string is independent, each string can be considered a random vector composed of independent random variables, resulting in its probability distribution being completely determined by a given multinomial distribution. By a random string, in this paper, we refer to a string whose letters are chosen randomly in accordance with the given multinomial distribution.

Let r be a given string. We now aim to determine the probability that a random string s , chosen from the same multinomial distribution $MN(p_1, \dots, p_\eta)$ as string r , is a subsequence of the string r . We denote this probability by $P(s \prec r)$. In the next theorem, we propose a new recurrence relation to calculate this probability.

Theorem 1. *Let r be a given string and s be a random string chosen from the same multinomial distribution. Then,*

$$P(s \prec r) = \begin{cases} 1, & \text{if } |s| = 0; \\ 0, & \text{if } |s| > |r|; \\ P(s[1] = r[1]) \cdot P(s[2, |s|] \prec r[2, |r|]) + P(s[1] \neq r[1]) \cdot P(s \prec r[2, |r|]), & \text{otherwise.} \end{cases} \tag{1}$$

Proof. It is clear from the definition of a subsequence that the empty string is a subsequence of every string and that a string cannot be a subsequence of a shorter one. Therefore, the cases $|s| = 0$ and $|s| > |r|$ are trivial. In the remaining case ($1 \leq |s| \leq |r|$),

$$P(s \prec r) = P(s[1] = r[1]) \cdot P(s[2, |s|] \prec r[2, |r|]) + P(s[1] \neq r[1]) \cdot P(s \prec r[2, |r|])$$

follows from the law of total probability. \square

The probability $P(s \prec r)$ in recurrence relation (1) is dependent not only on the length of string r , but also on the letter distribution of this string. Therefore, it is hard to come up with a closed-form expression for the general case of a multinomial distribution $MN(p_1, \dots, p_\eta)$. One way to deal with this problem is to consider some special cases of the multinomial distribution, for which closed-form expressions may be obtained.

2.1. Multinomial Distribution—Special Case 1: Uniform Distribution

The most frequently used form of the multinomial distribution considered in the literature is the uniform distribution. Since in this case, every letter has the same occurrence probability, probability $P(s \prec r)$ in the recurrence relation (1) depends only on the lengths $k = |s|$ and $l = |r|$ and can be written more simply as $P(k, l)$. This case is covered by Mousavi and Tabataba in [12], where the recurrence relation (1) is reduced as follows:

$$P(k, l) = \begin{cases} 1, & \text{if } k = 0; \\ 0, & \text{if } k > l; \\ \frac{1}{\eta} \cdot P(k - 1, l - 1) + \frac{\eta - 1}{\eta} \cdot P(k, l - 1), & \text{otherwise.} \end{cases} \tag{2}$$

Probabilities $P(k, l)$ can be calculated using dynamic programming as described by Mousavi and Tabataba in [12].

2.2. Multinomial Distribution—Special Case 2: Single Letter Exception

Let one letter $a_j \in \Sigma$ have occurrence probability $p \in (0, 1)$, $p \neq 1/\eta$ and each other letter $a_i, i \in \{1, \dots, \eta\} \setminus \{j\}$ have occurrence probability $(1 - p)/(\eta - 1)$. For this multinomial distribution, recurrence relation (1) reduces to:

$$P(s \prec r) = \begin{cases} 1, & \text{if } |s| = 0; \\ 0, & \text{if } |s| > |r|; \\ q \cdot P(s[2, |s|] \prec r[2, |r|]) + (1 - q) \cdot P(s \prec r[2, |r|]), & \text{otherwise.} \end{cases} \tag{3}$$

where

$$q := \begin{cases} p, & \text{if } r[1] = a_j; \\ \frac{1-p}{\eta-1}, & \text{otherwise.} \end{cases}$$

Note that, besides lengths $|s|$ and $|r|$, (3) depends only on whether or not a letter in the string r is equal to a_j .

2.3. Multinomial Distribution—Special Case 3: Two Sets of Letters

We now further generalize the previous case. Let $\{\Sigma_1, \Sigma_2\}$ be a partitioning of the alphabet Σ , that is, let $\Sigma_1, \Sigma_2 \subseteq \Sigma$ be nonempty sets such that $\Sigma = \Sigma_1 \cup \Sigma_2$ and $\Sigma_1 \cap \Sigma_2 = \emptyset$. Let us assume that every letter in Σ_1 has the same occurrence probability and also that every letter in Σ_2 has the same occurrence probability. We define

$$p_i := \begin{cases} \frac{p}{|\Sigma_1|}, & \text{if } a_i \in \Sigma_1; \\ \frac{1-p}{\eta-|\Sigma_1|}, & \text{if } a_i \in \Sigma_2, \end{cases}$$

where $p \in (0, 1)$ is the probability mass assigned to the set Σ_1 . For this multinomial distribution, recurrence relation (1) reduces to

$$P(s \prec r) = \begin{cases} 1, & \text{if } |s| = 0; \\ 0, & \text{if } |s| > |r|; \\ q \cdot P(s[2, |s|] \prec r[2, |r|]) + (1 - q) \cdot P(s \prec r[2, |r|]), & \text{otherwise,} \end{cases} \tag{4}$$

where

$$q := \begin{cases} \frac{p}{|\Sigma_1|}, & \text{if } r[1] \in \Sigma_1; \\ \frac{1-p}{\eta-|\Sigma_1|}, & \text{if } r[1] \in \Sigma_2. \end{cases}$$

This probability therefore depends on whether or not a letter in r belongs to the set Σ_1 or not.

2.4. The Case of Independent Random Strings

Another approach to calculating the probability that a string s is a subsequence of a string r is based on the assumption that both s and r are random strings chosen from the same multinomial distribution and are independent as a random vector. Using this setup, we established a recurrence relation for calculating probability $P(s \prec r)$.

Theorem 2. Let r and s be random independent strings chosen from the same multinomial distribution $MN(p_1, \dots, p_\eta)$. Then,

$$P(s \prec r) = \begin{cases} 1, & \text{if } |s| = 0; \\ 0, & \text{if } |s| > |r|; \\ \left(\sum_{i=1}^{\eta} p_i^2 \right) \cdot P(s[2, |s|] \prec r[2, |r|]) + \\ \quad \left(1 - \sum_{i=1}^{\eta} p_i^2 \right) \cdot P(s \prec r[2, |r|]), & \text{otherwise.} \end{cases} \quad (5)$$

Proof. The first two cases are trivial, so it remains to show the last case. Using the law of total probability, we obtain

$$P(s \prec r) = P(s[1] = r[1]) \cdot P(s[2, |s|] \prec r[2, |r|]) + P(s[1] \neq r[1]) \cdot P(s \prec r[2, |r|]).$$

Probability $P(s[1] = r[1])$ can be calculated with another application of the law of total probability, using the assumption that random strings s and r are mutually independent:

$$\begin{aligned} P(s[1] = r[1]) &= \sum_{i=1}^{\eta} P(r[1] = a_i) \cdot P(s[1] = r[1] \mid r[1] = a_i) \\ &= \sum_{i=1}^{\eta} P(r[1] = a_i) \cdot P(s[1] = a_i) = \sum_{i=1}^{\eta} p_i^2. \end{aligned}$$

□

Except for the obvious dependency on the multinomial distribution $MN(p_1, \dots, p_\eta)$, probability $P(s \prec r)$ is determined by the lengths of strings s and r , only. Therefore, as in the case of the uniform distribution, we can abbreviate this probability with $P(k, l)$, where $k = |s|$ and $l = |r|$. This allows us to pre-compute a probability matrix for all relevant values of k and l by means of dynamic programming.

3. Beam Search for Multinomially Distributed LCS Instances

In this section, we propose a new BS variant, which is characterized as follows:

- It makes use of a novel heuristic for search guidance that combines two complementary scoring functions by means of a convex combination.
- Our algorithm is a time-restricted BS variant, which dynamically adapts the beam width depending on the progress over the levels. In that way, the algorithm does not require a time-consuming tuning of the beam size with regard to the instance size.

3.1. Beam Search Framework

Beam search (BS) is a well-known search heuristic that is widely applied to many problems from various research fields, such as scheduling [31], speech recognition [32], machine learning tasks [33], packing problems [34], and so forth. It is a reduced version of breadth-first-search (BFS), where, instead of expanding all not-yet-expanded nodes from the same level, only a specific number of up to $\beta > 0$ nodes that appear most promising are selected and considered for expansions. In this way, BS keeps the search tree polynomial in size. The selection of the up to β nodes for further expansion is made according to a problem-specific heuristic guidance function h . The effectivity of the search thus substantially depends on this function. More specifically, BS works as follows. First, an initial beam B is set up with a root node r representing an initial state, in the case of

the LCS problem, this is the empty partial solution. At each major iteration, all nodes from beam B are expanded in all possible ways by considering all feasible actions. The obtained child nodes are kept in the set of extensions V_{ext} . Note that, for some problems, efficient filtering techniques can be applied to discard nodes from V_{ext} that are dominated by other nodes, that is, nodes that cannot yield better solutions. It is controlled by an internal parameter k_{filter} . This (possibly filtered) set of extensions is then sorted according to the nodes' values obtained from the guidance heuristic h , and the top β nodes (or less if V_{ext} is smaller) then form the beam B of the next level. The whole process is repeated level-by-level until B becomes empty. In general, to solve a combinatorial optimization problem, information about the longest (or shortest) path from the root node to a feasible goal node is kept to finally return a solution that maximizes or minimizes the problem's objective function. The pseudocode of such a general BS is given in Algorithm 1.

Algorithm 1 Beam Search.

```

1: Input: A problem instance, heuristic  $h$ ,  $\beta > 0$ ,  $k_{\text{filter}}$ 
2: Output: A heuristic solution
3:  $B \leftarrow \{r\}$ 
4: while  $B \neq \emptyset$  do
5:    $V_{\text{ext}} \leftarrow \emptyset$ 
6:   for  $v \in B$  do
7:     if  $v$  is a goal node then
8:       if node represents new best solution, store it
9:     else
10:      add not-yet-visited child nodes of  $v$  to  $V_{\text{ext}}$ 
11:    end if
12:  end for
13:  if  $k_{\text{filter}} \geq 0$  then
14:     $V_{\text{ext}} \leftarrow \text{Filter}(V_{\text{ext}}, k_{\text{filter}})$  // optionally filter dominated nodes
15:  end if
16:   $B \leftarrow \text{SelectBetaBest}(V_{\text{ext}}, \beta, h)$ 
17: end while
18: return best found solution

```

3.2. State Graph for the LCS Problem

The state graph for the LCS problem that is used by all BS variants is already well known in the literature, see, for example [16,24]. It is defined as a directed acyclic graph $G = (V, A)$, where a node $v = (\vec{\theta}^v, l^v) \in V$ represents the set of partial solutions, which:

1. have the same length l^v ;
2. induce the same subproblem denoted by $S[\vec{\theta}^v]$ w.r.t. the position vector $\vec{\theta}^v$.

We say that a partial solution s induces a subproblem $S[\vec{\theta}^v]$ iff $s_i[1, \vec{\theta}_i^v - 1]$ is the smallest prefix of s_i among all prefixes that has s as a subsequence.

An arc $a = (v_1, v_2) \in A$ exists between two nodes $v_1 \neq v_2 \in V$ and carries label $\ell(a) \in \Sigma$, iff

1. $l^{v_2} = l^{v_1} + 1$;
2. the partial solution that induces v_2 is obtained by appending $\ell(a)$ to the partial solution inducing v_1 .

The root node $r = ((1, \dots, 1), 0)$ of G refers to the original LCS problem on input string set S and can be said to be induced by the empty partial solution ε .

For deriving the successor nodes of a node $v \in V$, we first determine the subset $\Sigma_v \subseteq \Sigma$ of the letter that feasibly extends the partial solutions represented by v . The candidates for letter $a \in \Sigma_v$ are therefore all letters $a \in \Sigma$ that appear at least once in each string in

the subproblem given by strings $S[\vec{\theta}^v]$. This set Σ_v may be reduced by determining and discarding dominated letters. We say that letter $a \in \Sigma_v$ dominates letter $b \in \Sigma_v$ iff

$$\vec{\theta}_{i,a}^v \leq \vec{\theta}_{i,b}^v \quad \forall i \in \{1, \dots, m\}. \tag{6}$$

Dominated letters can be safely omitted since they lead to suboptimal solutions. Let $\Sigma_v^{\text{nd}} \subseteq \Sigma_v$ be the set of feasible and non-dominated letters. For each letter $a \in \Sigma_v^{\text{nd}}$, graph G contains a successor node $v' = (\vec{\theta}^{v'}, l^v + 1)$ of v , where $\vec{\theta}_i^{v'} = \vec{\theta}_{i,a}^{v'} + 1, i \in \{1, \dots, m\}$ (remember that $\vec{\theta}_{i,a}^{v'}$ denotes the position of the first appearance of letter a in string s_i from position $\vec{\theta}_i^{v'}$ onward). A node v that has no successor node, that is, when $\Sigma_v^{\text{nd}} = \emptyset$, is called a *non-extensible* node, or *goal* node. Among all goal nodes v we are looking for one representing a longest solution string, that is, a goal node with the largest l^v . Note that any path from the root node r to any node in $v \in V$ represents the feasible partial solution obtained by collecting and concatenating the labels of the traversed arcs. Thus, it is not necessary to store actual partial solutions s in the nodes. In the graph G , any path from root r to a non-extensible node represents a common, non-extensible subsequence of S . Any longest path from r to a goal node represents an optimal solution to problem instance S . As an example of a full state graph of an instance, see Figure 1.

We still have to explain the filtering of dominated nodes from the set V_{ext} , that is, procedure `Filter` in Algorithm 1. We adopt the efficient *restricted filtering* proposed in [13], which is parameterized by a filter size $k_{\text{filter}} > 0$. The idea is to select only the (up to) k_{filter} best nodes from V_{ext} and to check the dominance relation (6) for this subset of nodes in combination with all other nodes in V_{ext} . If the relation is positively evaluated, the dominated node is removed from V_{ext} . Note that parameter settings $k_{\text{filter}} = 0$ and $k_{\text{filter}} = |V_{\text{ext}}|$ represent the two extreme cases of no filtering and full filtering, respectively. A filter size of $0 < k_{\text{filter}} < |V_{\text{ext}}|$ may be meaningful, as full filtering may be too costly in terms of running time for larger beam widths.

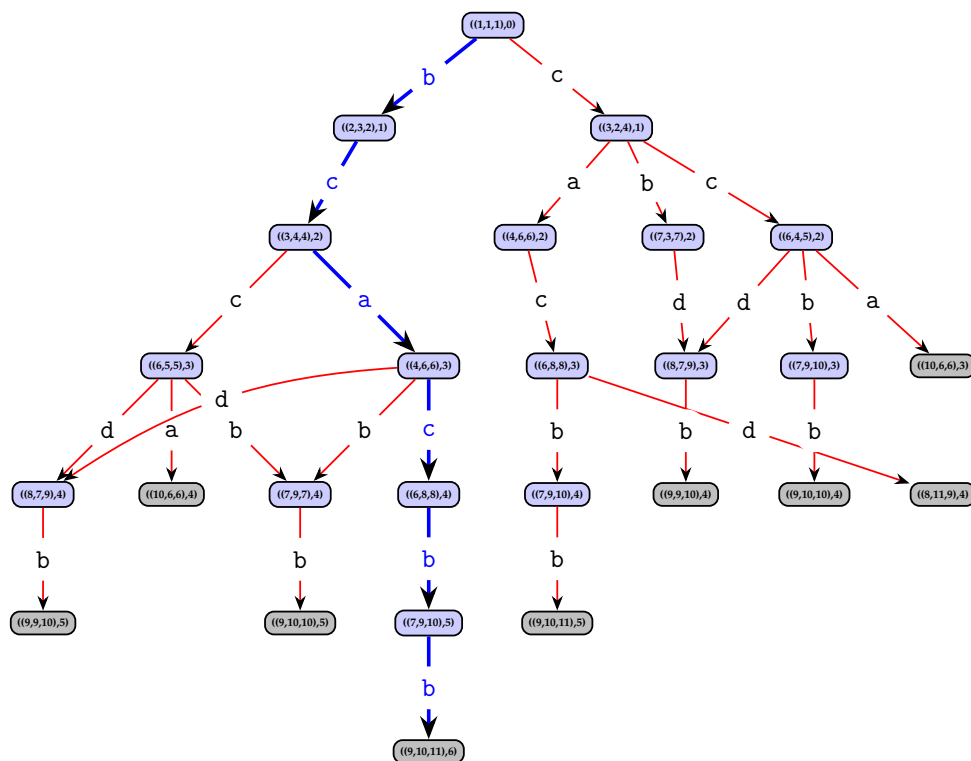


Figure 1. State graph for the LCS problem instance on strings $\{s_1 = \text{bcaacbdba}, s_2 = \text{cbccadcbdd}, s_3 = \text{bbccabcbba}\}$ and alphabet $\Sigma = \{a, b, c, d\}$. Light gray nodes are non-extensible goal nodes. The longest path in this state graph is shown in blue, leads from the root to node $((1, 10, 11), 6)$ and corresponds to the solution $s = \text{bcacbb}$, having length six.

3.3. Novel Heuristic Guidance

We now present a new heuristic for evaluating nodes in the BS in order to rank them and to select the beam of the next level. This heuristic, called the Geometric mean probability sum (GMPSUM), in particular aims at unbalanced instances and is a convex combination of the following two scores.

- The Geometric mean (GM) score is based on the geometric mean and the geometric standard deviation of the letters' occurrences across all input strings of the respective subproblem. It is calculated on a per letter basis and aggregated into a single numeric value;
- The Probability sum (PSUM) score is based on the previously introduced probability matrix $P(k, l)$ for the arbitrary unbalanced multinomial distribution case, see recurrence relation (5) or, in the special cases, any of the recurrence relations (3) and (4) might be used instead.

More specifically, for a given node v and a letter $a \in \Sigma$, let us define

$$C_a(S[\vec{\theta}^v]) = (s_1[\vec{\theta}_1^v|s_1|]_a, \dots, s_m[\vec{\theta}_m^v|s_m|]_a)$$

as the vector indicating for each remaining string of the respective subproblem the number of occurrences of letter $a \in \Sigma$. By $\mu_g(\cdot)$ and $\sigma_g(\cdot)$, we denote the geometric mean and geometric standard deviation, respectively, which are calculated for $\vec{x} = (x_1, \dots, x_m) \in \mathbb{R}^m$ by

$$\begin{aligned} \mu_g(\vec{x}) &= \sqrt[m]{x_1 \cdot \dots \cdot x_m}, \\ \sigma_g(\vec{x}) &= e \sqrt{\frac{\sum_{i=1}^m \left(\ln \frac{x_i}{\mu_g(\vec{x})} \right)^2}{m}}. \end{aligned}$$

By $UB_1(v)$, we denote the known upper bound on the length of an LCS for the subproblem represented by node v from [11], which is calculated as

$$UB_1(v) = \sum_{a \in \Sigma} \min_{i=1, \dots, m} C_a(S[\vec{\theta}^v])_i.$$

Then, the GM score is calculated as

$$GM(v) = GM(S[\vec{\theta}^v]) = \sum_{a \in \Sigma} \frac{\mu_g(C_a(S[\vec{\theta}^v]))}{\sigma_g(C_a(S[\vec{\theta}^v]))} \cdot \frac{\min_{i=1, \dots, m} C_a(S[\vec{\theta}^v])_i}{UB_1(v)}. \tag{7}$$

Overall, the GM score is thus a weighted average of the adjusted geometric means ($\mu_g(\cdot)/\sigma_g(\cdot)$) of the number of letter occurrences, and the weight of each letter is determined by normalizing the minimal number of the letter occurrences across all strings with the sum of the minimal number occurrences across all letters. The motivation behind this calculation is three-fold:

1. Letters with higher average numbers of occurrences across the strings will increase the chance of finding a longer common subsequence (composed of these letters).
2. Higher deviations around the mean naturally reduce this chance.
3. The minimal numbers of occurrences of a letter across all input strings is an upper bound on the length of common subsequences that can be formed by this single letter. Therefore, by normalizing it with the sum of all minimal letter occurrences, an impact of each letter in the overall summation is quantified.

The GM score is relevant if its underlying sampling geometric mean and standard deviation are based on a sample of sufficient size. In all our experiments, the minimal number of input strings is therefore ten. Working on samples of smaller sizes would likely make the GM score not that useful.

In addition to the GM score, we consider the PSUM score. For a given node v , we define

$$l_{\max}(v) = \min_{i=1, \dots, m} (|s_i| - \vec{\theta}_i^v + 1).$$

Then, the PSUM score is calculated by

$$\text{PSUM}(v) = \text{PSUM}(S[\vec{\theta}^v]) = \sum_{k=1}^{l_{\max}(v)} \prod_{i=1}^m P(k, |s_i| - \vec{\theta}_i^v + 1). \quad (8)$$

Unlike the GM score, which considers mostly general aspects of an underlying probability distribution, PSUM better captures more specific relations among input strings. It represents the sum of probabilities that a string of length k will be a common subsequence for all remaining input strings relevant for further extensions. Index k goes from one to $l_{\max}(v)$, that is, the length of the shortest possible non-empty subsequence up to the length of the longest possible one, which corresponds to the size of the shortest input string residual. The motivation behind using a simple (non-weighted) summation across all potential subsequence lengths is three-fold:

1. The exact length of the resulting subsequence is not known in advance. Note that, in the case of the HP heuristic proposed in [12], the authors heuristically determine an appropriate value of k for each level in the BS.
2. The summation across all k provides insight into the overall potential of node v —approximating the integral on the respective continuous function. Note that it is not required for this measure to have an interpretation in absolute terms since throughout the BS it is used strictly to compare different alternative extensions on the same level of the BS tree.
3. A more sophisticated approach that assigns different weights to the different k values would impose the challenge of deciding these specific weights. This would bring us back to the difficult task of an expected length prediction, which would be particularly hard when now considering the arbitrary multinomial distribution.

Finally, the total GMPSUM score is calculated by the linear combination

$$\text{GMPSUM}(v, \lambda) = \lambda \cdot \text{GM}(v) + (1 - \lambda) \cdot \text{PSUM}(v), \quad (9)$$

where $\lambda \in [0, 1]$ is a strategy parameter. Based on an empirical study with different benchmark instances and values for parameter λ , we came up with the following rules of thumb to select λ .

1. Since GM and PSUM have complementary focus, that is, they capture and award (or implicitly penalize) different aspects of the extension potential, their combined usage is indeed meaningful in most cases, that is, $0 < \lambda < 1$.
2. GM tends to be a better indicator when instances are more regular, that is, when each input string better fits the overall string distribution.
3. PSUM tends to perform better when instances are less regular, that is, when input strings are more dispersed around the overall string distribution.

Regarding the computational costs of the GMPSUM calculation, the GM score calculation requires $O(|\Sigma| \cdot m)$ time. This can be concluded from (7), where the most expensive part is the iteration through all letters from Σ and finding the minimal number of the letter occurrences across all m input strings ($\mu_g(\cdot)$ and $\sigma_g(\cdot)$ have the same time complexity). Note that the number of occurrences of each letter across all possible suffixes of all m input string positions is calculated in advance, before starting the beam search, and is stored in an appropriate three-dimensional array, see [24]. The worst case computational complexity of this step is $O(|\Sigma| \cdot m \cdot n_{\max})$. This is because the number of occurrences of a given letter across all positions inside the given input string can be determined in a single linear pass. Since this is done only at the start, and the expected number of GM calls is much higher than n_{\max} , this up front calculation can be neglected in the overall computational complexity

The PSUM score given by (8) takes $O(n_{min} \cdot m)$ time to be calculated due to a definition of $l_{max}(\cdot)$. Similar to GM, the calculation of matrix P is performed in pre-processing—its computational complexity corresponds to the number of entries, that is, $O(n_{max} \cdot n_{max})$, see (5).

Finally, the total computational complexity of GMPSUM can be concluded to be $O((|\Sigma| + n_{min}) \cdot m)$. The total computational complexity of the beam search is therefore a product of the number of calls of GMPSUM $O(n_{min} \cdot \beta \cdot |\Sigma|)$ and the time complexity of GMPSUM. Note that the number of GMPSUM calls equals the number of nodes created within a BS run. Since the LCS length, that is, the number of BS levels, is unknown, we use here n_{min} as the upper bound. Overall, the BS guided by GMPSUM runs in $O(n_{min} \cdot \beta \cdot |\Sigma| \cdot m \cdot (|\Sigma| + n_{min}))$ time if no filtering is performed. In the case of filtering, at each level of the BS, $O(\beta \cdot k_{filter} \cdot m)$ time is required, which gives $O(n_{min} \cdot \beta \cdot k_{filter} \cdot m)$ total time for executing the filtering within the BS. According to this, the BS guided by GMPSUM and utilizing (restricted) filtering requires $O(n_{min} \cdot \beta \cdot m \cdot (k_{filter} + |\Sigma|^2 + |\Sigma| \cdot n_{min}))$ time.

3.4. A Time-Restricted BS

In this section, we extend the basic BS from Algorithm 1 to a time-restricted beam search (TRBS). This BS variant is motivated by the desire to compare different algorithms with the same time-limit. The core idea we apply is to dynamically adapt the beam width in dependence of the progress over the levels.

Similarly to the standard BS from Algorithm 1, TRBS is parameterized with the problem instance to solve, the guidance heuristic h , and the filtering parameter k_{filter} . Moreover, what was previously the constant beam width β now becomes only the initial value. The goal is to achieve a runtime that comes close to a target time t_{max} , now additionally specified as an input. At the end of each major iteration, that is, level, if $t_{max} < +\infty$, that is, the time limit is actually enabled, the beam width for the next level is determined as follows.

1. Let t_{iter} be the time required for the current iteration.
2. We estimate the remaining number of major iterations (levels) by taking the maximum of lower bounds for the subinstances induced by the nodes in V_{ext} . More specifically,

$$LB_{max}(V_{ext}) = \max_{(v,a) \in V_{ext} \times \Sigma} \min_{i=1, \dots, m} C_a(S[\vec{\theta}^v])_i. \tag{10}$$

Thus, for each node $v \in V_{ext}$ and each letter a we consider the minimal number of occurrences of the letter across all string suffixes $S[\vec{\theta}^v]$ and select the one that is maximal. In other words, this LCS lower bound is based on considering all common subsequences in which a single letter is repeated as often as possible. In the literature, this procedure is known by the name *Long-run* [35] and provides a $|\Sigma|$ -approximation.

3. Let t_{rem} be the actual time still remaining in order to finish at time t_{max} .
4. Let $\overline{t_{rem}} = t_{iter} \cdot LB_{max}(V_{ext})$ be the expected remaining time when we would continue with the current beam width and the time spent at each level would stay the same as it was when measured for the current level.
5. Depending on the discrepancy between the actual and expected remaining time, we possibly increase or decrease the beam width for the next level:

$$\beta \leftarrow \begin{cases} \lfloor \beta \cdot 1.2 \rfloor & \text{if } t_{rem} / \overline{t_{rem}} > 1.1; \\ \min(100, \lfloor \beta / 1.2 \rfloor) & \text{if } t_{rem} / \overline{t_{rem}} < 0.9; \\ \beta & \text{otherwise.} \end{cases} \tag{11}$$

In this adaptive scheme, the thresholds for the discrepancy to increase or decrease the beam width, as well as the factor by which the beam width is modified, were determined empirically. Note that there might be better estimates of the LCS length than LB_{max} ; however, this estimate is inexpensive to obtain, and even if it underestimates or overestimates the LCS length in the early phases, gradually it converges toward the actual LCS length

as the algorithm progresses. This allows TRBS to smoothly adapt its expected remaining runtime to the desired one. Note that we only adapt the beam width and do not set it completely anew based on the runtime measured for the current level in order to avoid erratic changes of the beam width in the case of a larger variance of the level's runtimes. Based on preliminary experiments, we conclude that the proposed approach in general works well in achieving the desired time limit, while not changing β dramatically up and down in the course of a whole run. However, of course, how close we get to the time limit being met depends on the actual length of the LCS. For small solution strings, the approach has fewer opportunities to adjust β and then tends to overestimate the remaining time, thus utilizing less time than desired.

4. Experimental Results

In this section we evaluate our algorithms and compare them with the state-of-the-art algorithms from the literature. The proposed algorithms are implemented in C# and executed on machines with Intel i9-9900KF CPUs with @ 3.6 GHz and 64 Gb of RAM under Microsoft Windows 10 Pro OS. Each experiment was performed in single-threaded mode. We conducted two types of experiments:

- short-runs: these are limited-time scenarios—that is, BS configurations with $\beta = 600$ are used—executed in order to evaluate the quality of the guidance of each of the heuristics towards promising regions of the search space.
- long-runs: these are fixed-duration scenarios (900 s) in which we compare the time-restricted BS guided by the GMPSUM heuristic with the state-of-the-art results from the literature. The purpose of these experiments is the identification of new state-of-the-art solutions, if any.

4.1. Benchmark Sets

All relevant benchmark sets from the literature were considered in our experiments:

- Benchmark sets RAT, VIRUS and RANDOM, each one consisting of 20 single instances, are well known from the related literature [36]. The first two sets are biologically motivated, originating from the NCBI database. In the case of the third set, instances were randomly generated. The input strings in these sets are 600 characters long. Moreover, they contain instances based on alphabets of size four and 20.
- Benchmark set ES, introduced in [37], consists of randomly distributed input strings whose length varies from 1000 to 5000, while alphabet sizes range from two to 100. This set consists of 12 groups of instances.
- Benchmark set BB, introduced in [38], is different to the others because the input strings of each instance are generated such that there is a high similarity between them. For this purpose, first, a randomly generated base string was generated. Second, all input strings were generated based on the base string by probabilistically introducing small mutations such as delete/update operations of each letter. This set consists of eight groups (each one containing 10 single instances).
- Benchmark set BACTERIA, introduced in [39], is a real-world benchmark set used in the context of the constrained longest common subsequence problem. We made use of these instances by simply ignoring all pattern strings (constraints). This set consists of 35 single instances.
- Finally, we introduce two new sets of instances:

- The input strings of the instances of benchmark set POLY were generated such that the number of occurrences of each letter in each input string were determined by a multinomial distribution with known probabilities $p_1, \dots, p_\eta > 0$, such that $\sum_i p_i = 1$; see [40] for how to sample such distributions. More specifically, we used the multinomial distribution with $p_i = \frac{1}{2^i}, i = 1, \dots, |\Sigma| - 1$ and

$$p_\eta = 1 - \sum_{i=1}^{|\Sigma|-1} \frac{1}{2^i} \text{ for generating the input strings. The number of the occurrences}$$

of different letters was very much unbalanced in the obtained input strings. This set consists of 10 instances for each combination of the input string length $n \in \{100, 500, 1000\}$ and the number of input strings $m \in \{10, 50\}$, which makes a total of 60 problem instances.

- Benchmark set ABSTRACT, which will be introduced in Section 5, is a real-world benchmark set whose input strings are characterized by close-to-polynomial distributions of the different letters. The input strings originate from abstracts of scientific papers written in English.

4.2. Considered Algorithms

All considered algorithms make use of the state-of-the-art BS component. In order to test the quality of the newly proposed GMPSUM heuristic for the evaluation of the partial solutions at each step of BS, we compared it to the other heuristic functions that were proposed for this purpose in the literature: EX [16], POW [13], and HP [12]. The four resulting BS variants are labeled BS-GMPSUM, BS-EX, BS-POW, and BS-HP, respectively. These four BS variants were applied with the same parameter settings ($\beta = 600$ and $k_{\text{filter}} = 100$) in the short-run scenario in order to ensure that all of them used the same amount of resources.

In the long-run scenario, we tested the proposed time-restricted BS (TRBS) guided by the novel GMPSUM heuristic, which is henceforth labeled as TRBS-GMPSUM. Our algorithm was compared to the current state-of-the-art approach from the literature: A* + ACS [24]. These two algorithms were compared in the following way:

- Concerning A* + ACS, the results for benchmark sets RANDOM, VIRUS, RAT, ES and BB were taken from the original paper [24]. They were obtained with a computation time limit of 900 s per run. For the new benchmark sets—that is, POLY and BACTERIA—we applied the original implementation of A* + ACS with a time limit of 900 s on the above-mentioned machine.
- TRBS-GMPSUM was applied with a computation time limit of 600 s per run to all instances of benchmark sets RANDOM, VIRUS, RAT, ES and BB. Note that we reduced the computation time limit used in [24] by 50% because the CPU of our computer was faster than the one used in [24]. In contrast, the time limit for the new instances was set to 900 s. Regarding restricted-filtering, the same setting ($k_{\text{filter}} = 100$) as for the short-run experiments was used.

4.3. Tuning of Parameter λ

As already mentioned in the previous section, the values of parameters β (beam width) and k_{filter} (for filtering dominated nodes) were adopted from [16] (for the short-run executions) and from [24] (for the long-run executions). Regarding the strategy parameter λ that controls the impact of both scores in the GMPSUM score, we performed short-run evaluations across a discrete set of possible values: $\lambda \in \{0, 0.25, 0.5, 0.75, 1\}$. The conclusion was that the best performing values were $\lambda = 0$ for BB, $\lambda = 0.5$ for VIRUS and BACTERIA, $\lambda = 0.75$ for RANDOM, RAT and POLY, and $\lambda = 1$ for ES. The same settings for λ were used in the context of the long-run experiments. The settings of the parameters of our algorithms BS-GMPSUM and TRBS-GMPSUM are displayed in Tables 1 and 2.

Table 1. Values of parameters β and k_{filter} for all benchmark sets.

Param.	BS-GMPSUM	TRBS-GMPSUM
β	600	20,000
k_{filter}	100	100

Table 2. The settings of parameter λ per each benchmark set.

Benchmark Sets	BS-GMPSUM	TRBS-GMPSUM
BB	0	0
VIRUS & BACTERIA	0.5	0.5
RANDOM & RAT & POLY	0.75	0.75
ES	1.0	1.0

4.4. Summary of the Results

Before studying the results for each benchmark set in detail, we present a summary of the results in order to provide the reader with a broad picture of the comparison. More specifically, the results of the short-run scenarios are summarized in Table 3, while those for the long-run scenarios are given in Table 4. Table 3 displays the results in such a way that each line corresponds to a single benchmark set. The meanings of the columns are as follows: the first column contains the name of the benchmark set, while the second column provides the number of instances—respectively, instance groups—in the set. Then there are four blocks of columns, one for each considered BS variant. The first column of each block shows the obtained average solution quality ($\overline{|s|}$) over all instances of the benchmark set. The second column indicates the number of instances—respectively, instance groups—for which the respective BS variant achieves the best result (#b.). Finally, the third column provides the average running time (\bar{t}) in seconds over all instances of the considered benchmark set.

Table 3. Short-run results summary.

Benchmark Set	BS-Ex			BS-POW			BS-HP			BS-GMPSUM			
	Name	#	$\overline{ s }$	# b.	\bar{t}	$\overline{ s }$	# b.	\bar{t}	$\overline{ s }$	# b.	\bar{t}	$\overline{ s }$	# b.
Random	20	108.9	16	2.7	108.1	6	1.4	108.15	6	1.1	108.95	16	6.7
RAT	20	102.8	13	2.6	101.6	4	1.2	100.95	2	0.9	102.9	14	5.5
VIRUS	20	115.85	11	2.6	114.1	6	1.5	115.35	6	1.1	116.3	17	7.4
BB	8	407.13	2	8.5	430.13	6	6.3	422.94	4	3.6	424.86	5	26.9
ES	12	242.18	8	23	241.51	0	15.6	241.14	0	13.8	242.12	4	118.8
Poly	6	232.67	0	5.6	232.27	0	3.3	231.53	0	2.7	233.02	6	6.7
Bacteria	35	809.97	12	14.7	814.86	15	8.2	830.69	22	7.9	832.09	18	29.3
All	121		62			37			40			80	

The following conclusions can be drawn:

- Concerning the fully random benchmark sets RANDOM and ES, in which input strings were generated uniformly at random and are independent, it was already well-known previously that the heuristic guidance EX performs strongly. Nevertheless, it can be seen that BS-GMPSUM performs nearly as well as BS-EX, and clearly better than the remaining two BS variants.
- In the case of the quasi-random instances of benchmark sets VIRUS and RAT, BS-GMPSUM starts to show its strength by delivering the best solution qualities in 31 out of 40 cases. The second best variant is BS-EX, which still performs very well, and is able to achieve the best solution qualities in 24 out of 40 cases.
- For the special BB benchmark set, in which input strings were generated in order to be similar to each other, GMPSUM was revealed to perform comparably to the best variant BS-POW.
- Concerning the real-world benchmark set BACTERIA, BS-GMPSUM is able to deliver the best results for 18 out of 35 groups, which is slightly inferior to the BS-HP variant with 22 best-performances, and superior to variants BS-EX (12 cases) and BS-POW (15 cases). Concerning the average solution quality obtained for this benchmark set, BS-GMPSUM is able to deliver the best among all considered approaches.
- Concerning the multinominal non-uniformly distributed benchmark set POLY, BS-GMPSUM clearly outperforms all other considered BS variants. In fact, BS-GMPSUM is

able to find the best solutions for all six instance groups. Moreover, it beats the other approaches in terms of the average solution quality.

- Overall, BS-GMPSUM finds the best solutions in 80 (out of 121) instances or instance groups, respectively. The second best variant is BS-EX, which is able to achieve the best performance in 62 cases. In contrast, BS-Hp and BS-POW are clearly inferior to the other two approaches. We conclude that BS-GMPSUM performs well in the context of different letter distributions in the input strings, and it is worth trying this variant first when nothing is known about the distribution in the considered instance set.
- Overall, the running times of all four BS variants are comparable. The fastest one is BS-HP, while BS-GMPSUM requires somewhat more time compared to the others since it makes use of a heuristic function that combines two functions.

Table 4 provides a summary of the long-run scenarios, that is, it compares the current state-of-the-art algorithm $A^* + ACS$ with TRBS-GMPSUM. As the benchmark instances are the same as in the short-run scenarios, the first two table columns are the same as in Table 4. Then there are two blocks of columns, presenting the results of $A^* + ACS$ and TRBS-GMPSUM in terms of the average solution quality over all instances of the respective benchmark set ($\overline{|s|}$), and the number of instances (or instance groups) for which the respective algorithm archived the best result (#b.).

Table 4. Long-run results summary.

Benchmark Set		$A^* + ACS$		TRBS-GMPSUM 600 s/900 s	
Name	#	$\overline{ s }$	# b.	$\overline{ s }$	# b.
Random	20	109.9	20	109.7	16
RAT	20	104.3	17	104.4	18
VIRUS	20	117.0	14	117.3	19
BB	8	412.81	3	430.28	6
ES	12	243.82	9	243.73	4
Poly	6	234.13	4	234.23	5
Bacteria	35	829.26	10	862.63	33
All	121		77		101

The following can be concluded based on the results obtained for the long-run scenarios:

- Concerning RANDOM and ES, $A^* + ACS$ is—as expected—slightly better than TRBS-GMPSUM in terms of the number of best results achieved. However, when comparing the average performance, there is hardly any difference between the two approaches: 109.9 vs. 109.7 for the RANDOM benchmark set, and 243.82 vs. 243.73 for the ES benchmark set.
- In the context of benchmark sets RAT and VIRUS, TRBS-GMPSUM improves over the state-of-the-art results by a narrow margin. This holds both for the number of best results achieved and for the average algorithm performance.
- Concerning benchmark set BB, TRBS-GMPSUM significantly outperforms $A^* + ACS$. In six out of eight groups it delivers the best average solution quality, while $A^* + ACS$ does so only for three cases.
- The same holds for the real-world benchmark set BACTERIA, that is, TRBS-GMPSUM achieves the best results for 33 out of 35 instances, in contrast to only 10 instances in the case of $A^* + ACS$. Moreover, the average solution quality obtained is much better for TRBS-GMPSUM, namely 862.63 vs. 829.26.
- Finally, the performances of both approaches for benchmark set POLY are very much comparable.
- Overall, we can conclude that TRBS-GMPSUM is able to deliver the best results in 101 out of 121 cases, while $A^* + ACS$ does so only in 77 cases. This is because TRBS-GMPSUM provides a consistent solution quality across instances characterized by various kinds of letter distributions. It can therefore be stated that TRBS-GMPSUM is a new state-of-the-art algorithm for the LCS problem.

In summary, for the 32 random instances—respectively, instance groups—from the literature (sets RANDOM and ES) $A^* + ACS$ performs quite strongly due to the presumed randomness of the instances. However, the new TRBS-GMPSUM approach is not far behind. A weak point of $A^* + ACS$ becomes obvious when instances are not generated uniformly at random. In the 40 cases with quasi-random input strings (sets RAT and VIRUS), TRBS-GMPSUM performs the best in 37 cases, while $A^* + ACS$ does so in 31 cases. When input strings are similar to each other—see the eight instance groups of set BB— $A^* + ACS$ performs weakly compared to TRBS-GMPSUM. This tendency is reinforced in the context of the instances of set POLY (six instance groups) for which TRBS-GMPSUM clearly outperforms $A^* + ACS$ in all cases. The same holds for the real-world benchmark set BACTERIA. The overall conclusion yields that TRBS-GMPSUM works very well on a wide range of different instances. Moreover, concerning the instances from the previous literature (80 instances/groups), our TRBS-GMPSUM approach is able to obtain new state-of-the-art results in 13 cases. This will be shown in the next section.

4.5. New State-of-the-Art Results for Instances from the Literature

Due to space restrictions, we provide the complete set of results for each problem instance in the Supplementary Materials of the repository whose the link is given in the Data Availability Statement at the end of the paper.

The tables reporting the new state-of-the-art results are organized as follows. The first column contains the name of the corresponding benchmark set, while the following two columns identify, respectively, the instance (in the case of RAT and VIRUS) and the instance group (in the case of BB and ES). Afterwards, there are two columns that provide the best result known from the literature. The first of these columns provides the result, and the second column indicates the algorithm (together with the reference) that was the first one to achieve this result. Next, the tables provide the results of BS-EX, BS-POW, BS-HP and BS-GMPSUM in the case of the short-run scenario, the results of $A^* + ACS$ and TRBS-GMPSUM, respectively, in the case of the long-run scenario. Note that computation times are only given for the short-run scenario, because time served as a limit in the long-run scenario.

Concerning the short-run scenario (Table 5), BS-GMPSUM was able to produce new best results in 17 cases. This even includes four cases of the benchmark set ES, which was generated uniformly at random. The four cases of sets VIRUS and RAT are remarkable, in which the currently best-known solution was improved by two letters (see, for example, the case of set RAT and the instance $|\Sigma| = 4, m = 40$ and $n = 600$). Concerning the more important long-run scenario (Table 6), the best-known results so far were improved in 14 cases. Especially remarkable is the case of set BB, for which an impressive improvement of around 24 letters was achieved.

Table 5. New best results for the instances from the literature in the short-run scenario.

Instance (Group)			Literature Best $ s $		BS-EX		BS-POW		BS-HP		BS-GMPSUM		
Benchmark Set	$ \Sigma $	m	n	$ s $	Alg.	$ s $	t	$ s $	t	$ s $	t	$ s $	t
RAT	4	20	600	172	BS-EX	172	2.3	170	0.9	168	0.5	173	2.5
RAT	4	40	600	152	BS-EX	152	1.8	150	1	145	0.5	154	3.4
RAT	4	200	600	123	BS-EX	123	2.7	123	0.7	122	0.8	124	9.9
RAT	20	20	600	54	BS-EX	54	2.5	54	1.7	54	1.2	55	3.5
RAT	20	40	600	49	BS-EX	49	3	49	1.1	49	1.2	50	4.6
VIRUS	4	25	600	194	BS-EX	194	2.2	192	1.2	194	0.7	195	3.1
VIRUS	4	40	600	170	BS-EX	170	2.2	170	1.2	169	0.9	172	3.8
VIRUS	4	60	600	166	BS-EX	166	2.4	165	0.8	166	0.7	168	5.1
VIRUS	4	100	600	158	BS-EX	158	2.3	155	1.2	158	0.9	160	7.8
VIRUS	4	150	600	156	BS-EX	156	2.4	147	1.2	156	0.7	157	11
VIRUS	4	200	600	155	BS-HP	154	2.6	148	1.4	155	1.2	156	14.8
VIRUS	20	40	600	50	BS-EX	50	2.9	49	1.9	50	0.9	51	5.5
BB	2	100	1000	560.7	BS-POW	536.6	6.1	560.7	5.7	558.9	1.9	560.8	23.7
ES	2	10	1000	615.06	BS-EX	615.06	4.4	614.2	1.4	612.5	0.9	615.1	5.1
ES	10	50	1000	136.32	BS-EX	136.32	3.9	135.52	2.1	135.22	1.4	136.34	9.9
ES	25	10	2500	235.22	BS-POW	231.12	19.1	235.22	10.5	233.34	8	235.58	29
ES	100	10	5000	144.9	BS-POW	144.18	91.9	144.9	75.9	143.62	71.6	145.1	185.4

Table 6. New best results for the instances from the literature in the long-run scenario.

Instance (Group)			Literature Best $ s $		A* + ACS	TRBS-GMPSUM	
Benchmark Set	$ \Sigma $	m	n	$ s $	Alg.	$ s $	$ s $
RAT	4	20	600	174	A* + ACS	174	175
RAT	4	40	600	154	A* + ACS	154	156
RAT	20	25	600	52	A* + ACS	52	53
VIRUS	4	10	600	228	A* + ACS	228	229
VIRUS	4	15	600	206	A* + ACS	206	207
VIRUS	4	60	600	168	A* + ACS	168	169
VIRUS	4	80	600	163	A* + ACS	163	164
VIRUS	4	100	600	160	A* + ACS	160	162
VIRUS	4	150	600	157	A* + ACS	157	158
BB	2	100	1000	563.6	APS	547.1	571.1
BB	4	100	1000	390.2	APS	344.3	391.8
ES	2	10	1000	618.9	A* + ACS	618.9	619.1
ES	10	50	1000	137.5	A* + ACS	137.5	137.6
ES	25	10	2500	236.6	A* + ACS-DIST	235	238

4.6. Results for Benchmark Sets Poly and Bacteria

The tables reporting the results for benchmark set POLY are structured in the same way as those described above in the context of the other benchmark sets. The difference is that instance groups are identified by means of $|\Sigma|$ (first column), m (second column), and n (third column). The best results per instance group—that is, per table row—are displayed in bold font.

The results of the short-run scenario for benchmark set POLY are given in Table 7. According to the obtained results, a clear winner is BS-GMPSUM, which obtains the best average solution quality for all six instance groups. This indicates that GMPSUM is clearly better as a search guidance than the other three heuristic functions for this benchmark set. As previously mentioned, this is due to the strongly non-uniform nature of the instances, that is, the intentionally generated imbalance of the number of occurrences of different letters in the input strings. Nevertheless, the absolute differences between the results of BS-GMPSUM and BS-EX are not so high. The results of the long-run executions for benchmark set POLY are provided in Table 8. It can be observed that TRBS-GMPSUM and the state-of-the-art technique A* + ACS perform comparably.

Table 7. Short-run results for benchmark set POLY.

Instance Group			BS-EX		BS-POW		BS-HP		BS-GMPSUM	
$ \Sigma $	m	n	$ s $	t	$ s $	t	$ s $	t	$ s $	t
4	10	100	43.2	0.5	43.2	0.3	43.1	0.3	43.3	0.1
4	10	500	232.5	4.1	232.7	2.6	231.3	2.1	233	2.5
4	10	1000	470.7	8.8	470.1	5.4	467.3	4.2	470.9	10.3
4	50	100	35.7	0.6	35.6	0.4	35.5	0.3	35.8	0.3
4	50	500	201.4	6.1	200.8	3.5	200.4	3	202.3	6.2
4	50	1000	412.5	13.2	411.2	7.4	411.6	6.3	412.8	20.9

Table 8. Long-run results for benchmark set POLY.

Instance Group			A* + ACS		TRBS-GMPSUM	
$ \Sigma $	m	n	$ s $	$ s $	$ s $	t
4	10	100	43.4	43.4		580.7
4	10	500	234.3	234.3		890.5
4	10	1000	473.9	473.4		896.2
4	50	100	35.9	35.9		83.6
4	50	500	203	203.5		883.8
4	50	1000	414.3	414.9		892.3

Remember that, as in the case of POLY, the instances of benchmark set BACTERIA have been used for the first time in a study concerning the LCS problem. They were initially proposed in a study concerning the constrained LCS problem [39]. The results are again presented in the same way as described before. This set consists of 35 instances. Therefore, each line in Table 9 (short-run scenario) and Table 10 (long-run scenario) deals with one single instance, which is identified by $|\Sigma|$ (always equal to 4), m (varying between 2 and 383), n_{\min} (the length of the shortest input string) and n_{\max} (the length of the longest input string). The best results are indicated in bold font. The results obtained for the short-run scenario allow us to observe that BS-HP performs very well for this benchmark set. In fact, it obtains the best solution in 22 out of 35 cases. However, BS-GMPSUM is not far behind with 18 best solutions. Moreover, BS-GMPSUM obtains a slightly better average solution quality than BS-HP. Concerning the long-run scenario, as already observed before, TRBS-GMPSUM clearly outperforms A* + ACS. In fact, the differences are remarkable in some cases, such as, for example, instance number 32 (fourth but last line in Table 10) for which TRBS-GMPSUM obtains a solution of value 1241, while A* + ACS finds—in the same computation time—a solution of value 1204.

Table 9. Short-run results for benchmark set BACTERIA.

$ \Sigma $	Instance			BS-Ex		BS-POW		BS-HP		BS-GMPSUM	
	m	n_{\min}	n_{\max}	$ s $	t	$ s $	t	$ s $	t	$ s $	t
4	383	610	1553	256	31.1	252	13.9	279	16.8	271	94.1
4	3	1458	1458	1365	2.1	1365	1	1365	1.8	1365	7.7
4	33	1349	1577	610	17.6	605	10.6	755	10.8	689	36.5
4	106	1252	1520	503	25.1	483	12	515	12.2	514	61.8
4	2	1502	1502	1499	0	1499	0	1499	0	1499	0.1
4	12	1274	1413	659	13.3	636	8.5	627	6.9	659	18.9
4	15	1302	1515	598	13.3	602	8.5	655	7.7	678	20.7
4	13	1479	1557	811	15.8	752	10.1	1061	10	883	21.7
4	13	1308	1507	1037	17.6	1039	11.1	862	8.6	882	25.9
4	44	873	1543	493	16.3	473	9.3	470	7.8	494	29.6
4	4	1408	1530	1204	9	1271	6.3	1271	5.8	1271	15.9
4	173	1234	1847	502	34.7	463	15	541	18.3	525	97.5
4	13	1446	1551	681	14.5	713	9.5	794	8.6	785	22.2
4	88	1360	1545	583	27.3	570	13.8	667	15.1	601	67
4	2	1540	1548	1522	0.2	1522	0.1	1522	0.1	1522	0.3
4	3	1395	1424	1141	11.2	1141	6.8	1141	6.1	1141	15
4	4	1410	1488	886	9.8	1123	8	1123	6.7	1123	17.4
4	51	1266	1522	681	25.2	552	12.3	667	12	641	48.7
4	2	1461	1539	1354	0.9	1354	0.5	1354	1.7	1354	8.6
4	13	1246	1411	687	13	662	7.4	609	6.7	699	19.6
4	4	1434	1478	876	9.6	1112	8	1112	6.9	1112	16.3
4	18	1023	1438	464	11.9	468	7.6	458	5.8	475	14.2
4	2	1454	1460	1431	0.2	1431	0.1	1431	0.1	1431	0.3
4	8	1401	1533	1024	15.9	1061	9.8	858	7.5	864	18.8
4	33	990	1483	410	12.1	492	8.8	467	6.9	456	16.4
4	29	1422	1549	587	16.3	581	9.8	634	8.9	590	26.3
4	20	571	1394	438	9.6	405	5.4	401	4.5	431	11.7
4	96	1270	1565	516	24	467	11	531	12.3	522	55.5
4	10	1322	1455	1026	16	1026	9.7	796	7.1	1026	19.5
4	26	1334	1596	617	16.7	584	9.4	640	8.6	631	26.2
4	195	1345	1547	503	38.2	448	15.3	537	19.2	524	100.9
4	8	1454	1532	1221	16.4	1241	10	1241	8.6	1241	25.5
4	8	1359	1612	555	18.9	555	11.2	600	10.3	627	38.4
4	89	455	1587	251	11.3	214	4.7	233	4.8	239	18.2
4	2	1465	1469	1358	0.7	1358	0.4	1358	0.5	1358	6.8

Table 10. Long-run results for benchmark set BACTERIA.

Σ	Instance			A* + ACS	TRBS-GMPSUM	t
	m	n_{\min}	n_{\max}	s	s	
4	383	610	1553	265	273	887.2
4	3	1458	1458	1365	1365	810.9
4	33	1349	1577	670	723	899
4	106	1252	1520	518	532	897.1
4	2	1502	1502	1499	1499	0.1
4	12	1274	1413	665	694	899.7
4	15	1302	1515	680	708	899.6
4	13	1479	1557	842	883	899.5
4	13	1308	1507	870	1043	899.7
4	44	873	1543	514	501	897.3
4	4	1408	1530	1204	1271	898.4
4	173	1234	1847	520	528	895.6
4	13	1446	1551	732	816	899.6
4	88	1360	1545	557	634	897.9
4	2	1540	1548	1522	1522	0.3
4	3	1395	1424	1141	1141	899.7
4	4	1410	1488	1059	1123	899.4
4	51	1266	1522	659	871	898.9
4	2	1461	1539	1354	1354	851.9
4	13	1246	1411	716	727	899.6
4	4	1434	1478	1030	1112	899.2
4	18	1023	1438	481	488	898.4
4	2	1454	1460	1431	1431	0.3
4	8	1401	1533	1040	1063	899.2
4	33	990	1483	449	510	899.1
4	29	1422	1549	643	661	899
4	20	571	1394	439	432	899.7
4	96	1270	1565	529	546	897.2
4	10	1322	1455	1026	1026	899.7
4	26	1334	1596	654	676	899.3
4	195	1345	1547	514	544	894.2
4	8	1454	1532	1204	1241	898.3
4	8	1359	1612	624	644	897.9
4	89	455	1587	250	252	898.2
4	2	1465	1469	1358	1358	829.6

4.7. Statistical Significance of the So-Far Reported Results

In this section, we study the results of the short-run and long-run executions from a statistical point of view. In order to do so, Friedman’s test was performed simultaneously considering all four algorithms in the case of the short-run scenario, and the two considered algorithms in the case of the long-run scenario. All these tests and the resulting plots were generated using R’s *scamp* package [41].

Given that, in all cases, the test rejected the hypothesis that the algorithms perform equally, pairwise comparisons were performed using the Nemenyi post-hoc test [42]. The corresponding critical difference (CD) plots considering all benchmark sets together are shown in Figures 2 and 3a. Each algorithm is positioned in the segment according to its average ranking with regard to the average solution quality over all (121) the considered instance groups. The critical difference was computed with a significance level of 0.05. The performances of those algorithms whose difference is below the CD are regarded as performing statistically in an equivalent way—that is, no difference of statistical significance

can be detected. This is indicated in the figures by bold horizontal bars joining the respective algorithm markers.

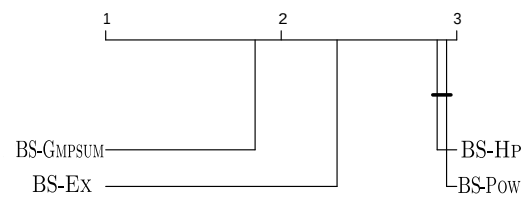


Figure 2. Critical difference (CD) plot over all considered benchmark sets (short-run executions).

Concerning short-run executions, BS-GMPSUM is clearly the overall best-performing algorithm, with statistical significance. BS-EX is in the second position. Moreover, the difference between BS-HP and BS-POW is not statistically significant. Concerning the long-run scenario, the best average rank is obtained by TRBS-GMPSUM. In the case of benchmark set BACTERIA, the difference between TRBS-GMPSUM and A* + ACS is significant, see Figure 3b. For the other benchmark sets, the two approaches perform statistically equivalently.

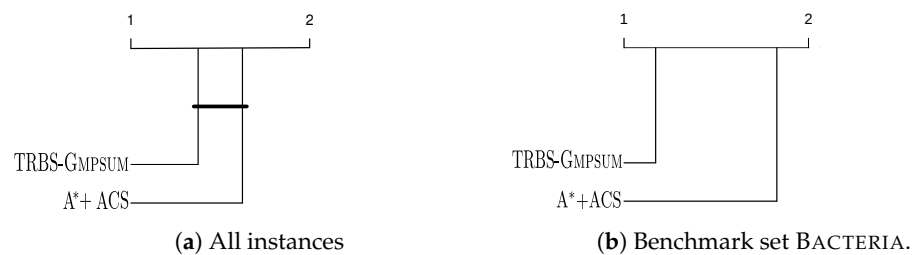


Figure 3. Critical difference (CD) plots concerning the long-run scenario.

5. Textual Corpus Case Study

In the previous section, we showed that the proposed method is highly competitive with state-of-the-art methods and generally outperforms them on instances sampled from non-uniform distributions. In order to further investigate the behavior of the proposed method in real-world instances with non-uniform distribution, we performed a case study on a corpus of textual instances originating from abstracts of scientific papers written in English. This set will henceforth be called ABSTRACT. It is known that letters in the English language are polynomially distributed [28]. The most frequent letter is *e*, with a relative frequency of 12.702%. The next most common letter is *t* (9.056%), followed by *a* (8.167%), and *o* (7.507%), and so forth.

In order to make a meaningful choice of texts we followed [43], where the authors measured the similarity between scientific papers, mainly from the field of artificial intelligence, by making use of various algorithms and metrics. By using tf-idf statistics with cosine similarity, their algorithm identified similar papers from a large paper collection. After that, the similarity between the papers proposed by their algorithm was manually checked and tagged by an expert as either similar (positive) or dissimilar (negative). The results of this research can be found at <https://cwi.ugent.be/respapersim> (accessed on 14 March 2021).

Keeping in mind that the LCS problem is also a measure of text similarity, we decided to check whether the abstracts of similar papers have longer common subsequences than abstracts of dissimilar papers. Therefore, the purpose of this case study is twofold: (1) to execute the LCS state-of-the-art methods along with the method proposed in this paper and to compare their performances on this specific instance set; and (2) to check whether the abstracts of similar papers have a higher LCS than those of dissimilar papers.

Based on these considerations, we formed two groups of twelve papers each, named POS and NEG. Group POS contains twelve papers which were identified as similar, while group NEG contains papers which are not similar to each other. We extracted abstracts

from each paper and pre-processed them in order to remove all letters except for those letters from the English alphabet. In addition, each uppercase letter was replaced with its lowercase pair.

For each of the two groups, we created a set of test instances as follows. For each $k \in \{10, 11, 12\}$ we generated $\binom{12}{k}$ different instances containing k input strings (considering all possible combinations). This resulted in the following set of instances for both POS and NEG:

- One instance containing all 12 abstracts as input strings.
- Twelve instances containing 11 out of 12 abstracts as input strings.
- Sixty-six instances containing 10 out of 12 abstracts as input strings.

Repeating our experimental setup presented in the previous section, we performed both short- and long-runs for the described instances. The obtained results for the short-run scenarios are shown in Table 11. The table is organized into five blocks of columns. The first block provides general information about the instances: NEG vs. POS, number of input strings (column with heading m), and the total number of instances (column #). The remaining four blocks contain the results of BS-EX, BS-POW, BS-HP and BS-GMPSUM, respectively. For each considered group of instances and each method, the following information about the obtained results is shown:

- $\overline{|s|}$: solution quality of the obtained LCS for the considered group of instances;
- #b.: number of cases in which the method reached the best result for the considered group of instances;
- \bar{t} : average execution time in seconds for the considered group of instances.

Table 11. Short-run results for the textual corpus instances (ABSTRACT).

Instance Set			BS-Ex			BS-Pow			BS-HP			BS-GMPSUM		
Name	m	#	$\overline{ s }$	#b.	t	$\overline{ s }$	#b.	\bar{t}	$\overline{ s }$	#b.	\bar{t}	$\overline{ s }$	#b.	\bar{t}
NEG	12	1	128	0	14.6	123	0	10.8	126	0	11.8	130	1	11.4
NEG	11	12	132.08	7	15	127	0	11.2	129.42	0	12.2	132.58	8	11.7
NEG	10	66	136.47	29	14.9	132.5	0	11.3	134.82	4	11.7	137.27	50	11.6
POS	12	1	134	1	15.2	128	0	11.9	131	0	11.8	133	0	7.4
POS	11	12	137.67	5	15	131.58	0	11.2	135.92	1	11.5	138.42	11	7.2
POS	10	66	143.33	42	14.5	135.85	0	10.7	141.53	10	11.5	143.14	39	7.2
All Negative		79		36			0			4			59	
All Positive		79		48			0			11			50	
All		158		84			0			15			109	

The results from Table 11 clearly indicate that the best results for instances based on group NEG are obtained by BS-GMPSUM. More precisely, BS-GMPSUM works best for the instance with 12 input strings, for eight out of 12 instances with 11 input strings and for 50 out of 66 instances with 10 input strings. In contrast, the second-best approach (BS-EX) reached the best result for 29 out of 66 instances with 10 input strings and seven out of 12 instances with 11 input strings. The remaining two methods were less successful for this group of instances. For the instances derived from group POS, BS-GMPSUM also achieved very good results. More specifically, BS-GMPSUM obtained the best results in almost all instances with 11 input strings. For instances with 10 input strings, BS-EX obtained the best results in 42 out of 66 cases, with BS-GMPSUM performing comparably (best result in 39 out of 66 cases). For the instance with 12 strings, the best solution was found by the BS-EX. Similarly to the instances from the NEG group, BS-HP and BS-POW are clearly less successful.

A summary of these results is provided in the last three rows of Table 11. Note that, in total, this table deals with 158 problem instances: 79 regarding group NEG, and another 79 regarding group POS. The summarized results show that the new GMPSUM guidance is, overall, more successful than its competitors. More precisely, BS-GMPSUM achieved the best results in 59 out of 79 cases concerning NEG, and in 50 out of 79 cases concerning POS.

Moreover, it can be observed that the average LCS length regarding the POS instances is greater than that regarding the NEG instances, across all m values.

Table 12 contains information about the long-run executions. The results obtained by $A^* + ACS$ and TRBS-GMPSUM are shown. The table is organized in a similar way to Table 11, with the exception that it does not contain information about execution times, since computation time served as the stopping criterion. As can be seen from the overall results at the bottom of Table, TRBS-GMPSUM obtains more best results than $A^* + ACS$ for both groups of instances (NEG and POS). More precisely, it obtained the best result for the instances with 12 input strings, both in the case of POS and NEG, while $A^* + ACS$ achieved the best result only in the case of the POS instance with 12 input strings. Concerning the results for the instances with 11 input strings, it can be noticed that—in the case of the NEG instances—TRBS-GMPSUM delivers 11 out of 12 best results, while the $A^* + ACS$ method does so only in two out of twelve cases. Regarding the POS instances with 11 input strings, the difference becomes smaller. More specifically, TRBS-GMPSUM achieves nine out of 12 best results, while $A^* + ACS$ achieved six out of 12 best results. A corresponding comparison can be made for the instances with 10 input strings. For the instances concerning group NEG, TRBS-GMPSUM delivers the best results for 60 out of 66 instances, while $A^* + ACS$ can find the best results only in 31 cases. Finally, in the case of the POS instances, the best results were achieved in 45 out of 66 cases by TRBS-GMPSUM, and in 41 out of 66 cases by $A^* + ACS$. The long-run results also indicate that abstracts of similar papers are characterized by generally longer LCS measures.

Table 12. Long-run results for the textual corpus instances (ABSTRACT).

Instance Set			$A^* + ACS$		TRBS-GMPSUM		
Name	m	#	$\overline{ s }$	#best	$\overline{ s }$	#best	\bar{t}
NEG	12	1	129	0	130	1	895.7
NEG	11	12	133.25	2	134.33	11	897.2
NEG	10	66	138.32	31	139.12	60	897.8
POS	12	1	136	1	136	1	896.5
POS	11	12	140.17	6	140.42	9	896.8
POS	10	66	145.33	41	145.52	45	897.4
All Negative		79		33		72	
All Positive		79		48		55	
All		158		81		127	

6. Conclusions and Future Work

In this paper, we considered the prominent longest common subsequence problem with an arbitrary set of input strings. We proposed a novel search guidance, named GMPSUM, for tree search algorithms. This new guidance function was defined as a convex combination of two complementary heuristics: (1) the first one is suited for instances in which the distribution of letters is close to uniform-at-random; and (2) the second is convenient for all cases in which letters are non-uniformly distributed. The combined score produced by these two heuristics provides a guidance function which navigates the search towards promising regions of the search space, on a wide range of instances with different distributions. We ran short-run experiments in which beam search makes use of a comparable number of iterations under different guidance heuristics. The conclusion was that the novel guidance heuristic performs statistically equivalently to the best-so-far heuristic from the literature on close-to-random instances. Moreover, it was shown that it significantly outperforms the known search guidance functions on instances with a non-uniform letter frequency per input string. This capability of the proposed heuristic to deal with a non-uniform scenario was validated on two newly introduced benchmark sets: (1) POLY, whose input strings are generated from a multinomial distribution; and (2) ABSTRACT, which are real-world instances whose input strings follow a multinomial distribution and originate from abstracts of scientific papers written in English. In a second part of the experimentation, we performed long-run executions. For this purpose we combined the GMPSUM guidance function with a time-restricted BS that dynamically

adapts its beam width during execution such that the overall running time is very close to the desired time limit. This algorithm was able to significantly outperform the best approach from the literature ($A^* + ACS$). More specifically, the best-known results from the literature were at least matched for 63 out of the 80 considered instance groups. Moreover, regarding the two new benchmark sets (POLY and BACTERIA), the time-restricted BS guided by GMPSUM was able to deliver solutions that were equally as good as, and in most cases better than, those of $A^* + ACS$, in 38 out of 41 instance groups.

In future work we plan to adapt GMPSUM to other LCS-related problems such as the constrained longest common subsequence problem [44], the repetition-free longest common subsequence problem [45], the LCS problem with a substring exclusion constraint [46], and the longest common palindromic subsequence problem [47]. It would also be interesting to incorporate this new guidance function into the leading hybrid approach $A^* + ACS$ to possibly further boost the obtained solution quality.

Author Contributions: B.N. was responsible for conceptualization, methodology, and writing. A.K. was responsible for software implementation, writing and methodology. M.D. was responsible for conceptualization, writing and visualization. M.G. was responsible for resources and writing. C.B. was responsible for writing, editing and supervision. G.R. was responsible for concept polishing, supervision, and funding acquisition. All authors have read and agreed to the published version of the manuscript.

Funding: The research was fully funded by the Austrian Science Fund (FWF) under project No. W1260-N35.

Data Availability Statement: The reported results can be found at <https://github.com/milanagrbc/LCSonNuD> (accessed on 22 April 2021).

Acknowledgments: This research was partially supported by the Ministry for Scientific and Technological Development, Higher Education and Information Society, Government of Republic of Srpska, B&H under the Project “Development of artificial intelligence methods for solving computer biology problems”, project No. 19.032/-961-24/19. Marko Djukanovic was funded by the Doctoral Program Vienna Graduate School on Computational Optimization (VGSCO), Austrian Science Foundation, project No. W1260-N35. Christian Blum was funded by project CI-SUSTAIN of the Spanish Ministry of Science and Innovation (PID2019-104156GB-I00).

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

LCS	Longest Common Subsequence
BS	Beam Search
ACS	Anytime column search
APS	Anytime pack search
$A^* + ACS$	The hybrid of A^* and ACS
TRBS	Time restricted BS
GM	Geometric mean score
PSUM	Probability sum score
GMPSUM	Geometric mean probability sum score
TRBS-GMPSUM	TRBS guided by GMPSUM
BS-GMPSUM	BS guided by GMPSUM

References

1. Maier, D. The Complexity of Some Problems on Subsequences and Supersequences. *J. ACM* **1978**, *25*, 322–336. [[CrossRef](#)]
2. Minkiewicz, P.; Darewicz, M.; Iwaniak, A.; Sokołowska, J.; Starowicz, P.; Bucholska, J.; Hryniewicz, M. Common Amino Acid Subsequences in a Universal Proteome—Relevance for Food Science. *Int. J. Mol. Sci.* **2015**, *16*, 20748–20773. [[CrossRef](#)]
3. Storer, J. *Data Compression: Methods and Theory*; Computer Science Press: Rockville, MD, USA, 1988.
4. Beal, R.; Afrin, T.; Farheen, A.; Adjeroh, D. A new algorithm for “the LCS problem” with application in compressing genome resequencing data. *BMC Genom.* **2016**, *17*, 544. [[CrossRef](#)]

5. Kruskal, J.B. An overview of sequence comparison: Time warps, string edits, and macromolecules. *SIAM Rev.* **1983**, *25*, 201–237. [[CrossRef](#)]
6. Xie, X.; Liao, W.; Aghajan, H.; Veelaert, P.; Philips, W. Detecting Road Intersections from GPS Traces Using Longest Common Subsequence Algorithm. *ISPRS Int. J. Geo-Inf.* **2017**, *6*, 1. [[CrossRef](#)]
7. Bergroth, L.; Hakonen, H.; Raita, T. A survey of longest common subsequence algorithms. In Proceedings of the SPIRE 2000—The 7th International Symposium on String Processing and Information Retrieval, Coruna, Spain, 27–29 September 2000; pp. 39–48.
8. Gusfield, D. *Algorithms on Strings, Trees, and Sequences*; Computer Science and Computational Biology; Cambridge University Press: Cambridge, UK, 1997.
9. Fraser, C.B. Subsequences and Supersequences of Strings. Ph.D. Thesis, University of Glasgow, Glasgow, UK, 1995.
10. Huang, K.; Yang, C.; Tseng, K. Fast Algorithms for Finding the Common Subsequences of Multiple Sequences. In Proceedings of the ICS 2004—The 9th International Computer Symposium, Funchal, Portugal, 13–16 January 2004.
11. Blum, C.; Blesa, M.J.; López-Ibáñez, M. Beam search for the longest common subsequence problem. *Comput. Oper. Res.* **2009**, *36*, 3178–3186. [[CrossRef](#)]
12. Mousavi, S.R.; Tabataba, F. An improved algorithm for the longest common subsequence problem. *Comput. Oper. Res.* **2012**, *39*, 512–520. [[CrossRef](#)]
13. Tabataba, F.S.; Mousavi, S.R. A hyper-heuristic for the longest common subsequence problem. *Comput. Biol. Chem.* **2012**, *36*, 42–54. [[CrossRef](#)]
14. Wang, Q.; Korkin, D.; Shang, Y. A fast multiple longest common subsequence (MLCS) algorithm. *IEEE Trans. Knowl. Data Eng.* **2011**, *23*, 321–334. [[CrossRef](#)]
15. Djukanovic, M.; Raidl, G.R.; Blum, C. Anytime algorithms for the longest common palindromic subsequence problem. *Comput. Oper. Res.* **2020**, *114*, 104827. [[CrossRef](#)]
16. Djukanovic, M.; Raidl, G.; Blum, C. A Beam Search for the Longest Common Subsequence Problem Guided by a Novel Approximate Expected Length Calculation. In Proceedings of the LOD 2019—The 5th International Conference on Machine Learning, Optimization, and Data Science, Siena, Italy, 10–13 September 2019.
17. Blum, C.; Festa, P. Longest Common Subsequence Problems. In *Metaheuristics for String Problems in Bioinformatics*; Wiley: Hoboken, NJ, USA, 2016; Chapter 3, pp. 45–60.
18. Chan, H.T.; Yang, C.B.; Peng, Y.H. The Generalized Definitions of the Two-Dimensional Largest Common Substructure Problems. In Proceedings of the 33rd Workshop on Combinatorial Mathematics and Computation Theory, Taipei, Taiwan, 13–14 May 2016; pp. 1–12.
19. Peng, Z.; Wang, Y. A Novel Efficient Graph Model for the Multiple Longest Common Subsequences (MLCS) Problem. *Front. Genet.* **2017**, *8*, 104. [[CrossRef](#)] [[PubMed](#)]
20. Li, Y.; Wang, Y.; Zhang, Z.; Wang, Y.; Ma, D.; Huang, J. A novel fast and memory efficient parallel MLCS algorithm for long and large-scale sequences alignments. In Proceedings of the IEEE 32nd International Conference on Data Engineering, Helsinki, Finland, 16–20 May 2016; pp. 1170–1181.
21. Wang, C.; Wang, Y.; Cheung, Y. A branch and bound irredundant graph algorithm for large-scale MLCS problems. *Pattern Recognit.* **2021**, *119*, 108059. [[CrossRef](#)]
22. Li, Y.; Liu, B.; Wang, Z.; Cui, J.; Yao, K.; Lv, P.; Shen, Y.; Xu, Y.; Guan, Y.; Ma, X. COVID-19 Evolves in Human Hosts. *arXiv* **2020**, arXiv:2003.05580.
23. Wei, S.; Wang, Y.; Yang, Y.; Liu, S. A path recorder algorithm for Multiple Longest Common Subsequences (MLCS) problems. *Bioinformatics* **2020**, *36*, 3035–3042. [[CrossRef](#)]
24. Djukanovic, M.; Raidl, G.R.; Blum, C. Finding Longest Common Subsequences: New anytime A* search results. *Appl. Soft Comput.* **2020**, *95*, 106499. [[CrossRef](#)]
25. Vadlamudi, S.G.; Gaurav, P.; Aine, S.; Chakrabarti, P.P. Anytime column search. In Proceedings of the AI'12—The 25th Australasian Joint Conference on Artificial Intelligence, Sydney, Australia, 4–7 December 2012; pp. 254–265.
26. Yang, J.; Xu, Y.; Sun, G.; Shang, Y. A New Progressive Algorithm for a Multiple Longest Common Subsequences Problem and Its Efficient Parallelization. *IEEE Trans. Parallel Distrib. Syst.* **2013**, *24*, 862–870. [[CrossRef](#)]
27. Vadlamudi, S.G.; Aine, S.; Chakrabarti, P.P. Anytime pack search. *Nat. Comput.* **2016**, *15*, 395–414. [[CrossRef](#)]
28. Lewand, R.E. *Cryptological Mathematics*; American Mathematical Soc.: Providence, RI, USA, 2000; Volume 16.
29. Beutelspacher, A. *Kryptologie*; Springer: Berlin, Germany, 1996; Volume 7.
30. Bakaev, M. Impact of familiarity on information complexity in human-computer interfaces. In Proceedings of the MATEC Web of Conferences, EDP Sciences, Munich, Germany, 6–8 June 2016; Volume 75, p. 08003.
31. Ow, P.S.; Morton, T.E. Filtered beam search in scheduling. *Int. J. Prod. Res.* **1988**, *26*, 35–62. [[CrossRef](#)]
32. Ney, H.; Haeb-Umbach, R.; Tran, B.H.; Oerder, M. Improvements in beam search for 10000-word continuous speech recognition. In Proceedings of the ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing, San Francisco, CA, USA, 23–26 March 1992; Volume 1, pp. 9–12.
33. Kumar, A.; Vembu, S.; Menon, A.K.; Elkan, C. Beam search algorithms for multilabel learning. *Mach. Learn.* **2013**, *92*, 65–89. [[CrossRef](#)]
34. Araya, I.; Riff, M.C. A beam search approach to the container loading problem. *Comput. Oper. Res.* **2014**, *43*, 100–107. [[CrossRef](#)]

35. Jiang, T.; Li, M. On the approximation of shortest common supersequences and longest common subsequences. *SIAM J. Comput.* **1995**, *24*, 1122–1139. [[CrossRef](#)]
36. Shyu, S.J.; Tsai, C.Y. Finding the longest common subsequence for multiple biological sequences by ant colony optimization. *Comput. Oper. Res.* **2009**, *36*, 73–91. [[CrossRef](#)]
37. Easton, T.; Singireddy, A. A large neighborhood search heuristic for the longest common subsequence problem. *J. Heuristics* **2008**, *14*, 271–283. [[CrossRef](#)]
38. Blum, C.; Blesa, M.J. Probabilistic beam search for the longest common subsequence problem. In Proceedings of the International Workshop on Engineering Stochastic Local Search Algorithms, Brussels, Belgium, 6–8 September 2007; pp. 150–161.
39. Djukanovic, M.; Kartelj, A.; Matic, D.; Grbic, M.; Blum, C.; Raidl, G. Solving the Generalized Constrained Longest Common Subsequence Problem with Many Pattern Strings; Technical Report AC-TR-21-008. 2021. Available online: <https://www.ac.tuwien.ac.at/files/tr/ac-tr-21-008.pdf> (accessed on 25 June 2021).
40. Kesten, H.; Morse, N. A property of the multinomial distribution. *Ann. Math. Stat.* **1959**, *30*, 120–127. [[CrossRef](#)]
41. Calvo, B.; Santafé Rodrigo, G. scmamp: Statistical Comparison of Multiple Algorithms in Multiple Problems. 2016. Available online: <https://journal.r-project.org/archive/2016/RJ-2016-017/RJ-2016-017.pdf> (accessed on 25 June 2021).
42. Pohlert, T. The pairwise multiple comparison of mean ranks package (PMCMR). *R Package* **2014**, *27*, 9.
43. Magara, M.B.; Ojo, S.O.; Zuva, T. A comparative analysis of text similarity measures and algorithms in research paper recommender systems. In Proceedings of the 2018 Conference on Information Communications Technology and Society (ICTAS), Durban, South Africa, 8–9 March 2018; pp. 1–5.
44. Gotthilf, Z.; Hermelin, D.; Lewenstein, M. Constrained LCS: Hardness and Approximation. In Proceedings of the CPM 2008—The 19th Annual Symposium on Combinatorial Pattern Matching, Pisa, Italy, 18–20 June 2008; Volume 5029, pp. 255–262.
45. Adi, S.S.; Braga, M.D.; Fernandes, C.G.; Ferreira, C.E.; Martinez, F.V.; Sagot, M.F.; Stefanos, M.A.; Tjandraatmadja, C.; Wakabayashi, Y. Repetition-free longest common subsequence. *Discret. Appl. Math.* **2010**, *158*, 1315–1324. [[CrossRef](#)]
46. Zhu, D.; Wang, X. A Simple Algorithm for Solving for the Generalized Longest Common Subsequence (LCS) Problem with a Substring Exclusion Constraint. *Algorithms* **2013**, *6*, 485–493. [[CrossRef](#)]
47. Chowdhury, S.R.; Hasan, M.M.; Iqbal, S.; Rahman, M.S. Computing a Longest Common Palindromic Subsequence. *Fundam. Informaticae* **2014**, *129*, 329–340. [[CrossRef](#)]