JOURNAL OF COMPUTATIONAL BIOLOGY Volume 29, Number 11, 2022

© Mary Ann Liebert, Inc. Pp. 1237–1251

DOI: 10.1089/cmb.2022.0266

Open camera or QR reader and scan code to access this article and other resources online.



# Algorithms for Colinear Chaining with Overlaps and Gap Costs

CHIRAG JAIN, DANIEL GIBNEY, and SHARMA V. THANKACHAN3

#### **ABSTRACT**

Colinear chaining has proven to be a powerful heuristic for finding near-optimal alignments of long DNA sequences (e.g., long reads or a genome assembly) to a reference. It is used as an intermediate step in several alignment tools that employ a seed-chain-extend strategy. Despite this popularity, efficient subquadratic time algorithms for the general case where chains support anchor overlaps and gap costs are not currently known. We present algorithms to solve the colinear chaining problem with anchor overlaps and gap costs in O(n)time, where n denotes the count of anchors. The degree of the polylogarithmic factor depends on the type of anchors used (e.g., fixed-length anchors) and the type of precedence an optimal anchor chain is required to satisfy. We also establish the first theoretical connection between colinear chaining cost and edit distance. Specifically, we prove that for a fixed set of anchors under a carefully designed chaining cost function, the optimal "anchored" edit distance equals the optimal colinear chaining cost. The anchored edit distance for two sequences and a set of anchors is only a slight generalization of the standard edit distance. It adds an additional cost of one to an alignment of two matching symbols that are not supported by any anchor. Finally, we demonstrate experimentally that optimal colinear chaining cost under the proposed cost function can be computed orders of magnitude faster than edit distance, and achieves correlation coefficient >0.9 with edit distance for closely as well as distantly related sequences.

**Keywords:** alignment, colinear chaining, edit distance.

## 1. INTRODUCTION

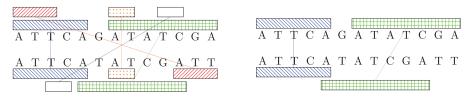
NE OF THE MOST FUNDAMENTAL PROBLEMS in computational biology is computing the optimal alignment of two sequences. Unfortunately, conditional lower bounds suggest that an algorithm for computing an optimal alignment, or edit distance, in strongly subquadratic time is unlikely (Backurs and

<sup>&</sup>lt;sup>1</sup>Department of Computational and Data Sciences, Indian Institute of Science, Bengaluru, India.

<sup>&</sup>lt;sup>2</sup>School of Computational Science and Engineering, Georgia Institute of Technology Atlanta, Georgia, USA.

<sup>&</sup>lt;sup>3</sup>Department of Computer Science, University of Central Florida, Orlando, Florida, USA.

A preliminary version of this article appeared in RECOMB 2022.



**FIG. 1.** (Left) Anchors representing a set of exact matches are shown as rectangles. The colinear chaining problem is to find an optimal ordered subset of anchors subject to some cost function. (Right) A chain of overlapping anchors.

Indyk, 2015; Hoppenworth et al, 2020). This lower bound indicates a challenge for scaling the computation of edit distance to high-throughput sequencing data. Instead, heuristics are often used to obtain an approximate solution in less time and space. One such popular heuristic is colinear chaining.

This technique involves precomputing fragments between the two sequences that closely agree (in this study, exact matches are called *anchors*), then determining which of these anchors should be kept within the alignment (Fig. 1). Techniques along these lines are used in long-read mappers (Chaisson and Tesler, 2012; Jain et al, 2022; Li, 2018; Li et al, 2020; Ren and Chaisson, 2021; Sahlin and Mäkinen, 2021; Sedlazeck et al, 2018) and generic sequence aligners (Abouelhoda et al, 2008; Bray et al, 2003; Kurtz et al, 2004; Marçais et al, 2018; Otto et al, 2011). We focus on the following problem (described formally in Section 2): Given a set of *n* anchors, determine an optimal ordered subset (or chain) of these anchors.

Several algorithms have been developed for colinear chaining (Abouelhoda and Ohlebusch, 2005; Mäkinen and Sahlin, 2020; Shibuya and Kurochkin, 2003; Uricaru et al, 2011) and even more in the context of sparse dynamic programming (Eppstein et al, 1992a; Eppstein et al, 1992b; Myers and Miller, 1995; Morgenstern, 2002; Mäkinen et al, 2019; Wilbur and Lipman, 1983). Solutions with different time complexities exist for different variations of this problem. These depend on the cost function assigned to a chain and the types of chains permitted. Solutions include an algorithm running in  $O(n \log n \log \log n)$  time for a simpler variant of the problem where anchors used in a solution must be nonoverlapping (Abouelhoda and Ohlebusch, 2005).

More recently, Mäkinen and Sahlin (2020) gave an algorithm running in  $O(n \log n)$  time where anchor overlaps are allowed, but gaps between anchors are not considered in the cost function. None of the solutions introduced thus far provide a subquadratic time algorithm for variations that use both overlap and gap costs. However, including overlaps and gaps into a cost function is a more realistic model for anchor chaining. For example, consider a simple scenario where minimizers (Schleimer et al, 2003) are used to identify anchors. Suppose query and reference sequences are identical, then adjacent minimizer anchors will likely overlap. Not allowing anchor overlaps during chaining will lead to a penalty cost associated with gaps between chained anchors despite the two strings being identical. Therefore, depending on the type of anchor, there may be no reason to assume that in an optimal alignment, the anchors would be nonoverlapping. At the same time, not penalizing long gaps between the anchors is unlikely to produce correct alignments. This is why both anchor overlaps and gap costs are supported during chaining in widely used aligners, for example, Minimap2 (Kalikar et al, 2022; Li, 2018) and Nucmer4 (Marçais et al, 2018).

This study's contribution is the following:

- We provide the first algorithm running in subquadratic  $\tilde{O}(n)$  time for chaining with overlap and gap costs.\* Refinements based on the specific type of anchor and chain under consideration are also given. These refinements include an  $O(n \log^2 n)$  time algorithm for the case where all anchors are of the same length, as is the case with k-mers.
- When n is not too large (less than the sequence lengths), we present an algorithm with  $O(n \cdot \text{OPT} + n \log n)$  average case time, where OPT is the optimal solution value. This provides a simple algorithm that is efficient in practice.
- Using a carefully designed cost function, we mathematically relate the optimal chaining cost with a generalized version of edit distance, which we call *anchored edit distance*. This is equivalent to the usual edit distance with the modification that matches performed without the support of an anchor have unit cost. A more formal definition appears in Section 2. With our cost function, we prove that the optimal chaining cost is equal to the anchored edit distance.

 $<sup>*\</sup>tilde{O}(\cdot)$  hides polylogarithmic factors.

• We empirically demonstrate that computing optimal chaining cost is orders of magnitude faster than computing edit distance, especially in semiglobal comparison mode. We also demonstrate a strong correlation between optimal chaining cost and edit distance. The correlation coefficients are favorable when compared with suboptimal chaining methods implemented in Minimap2 and Nucmer4.

#### 2. PRELIMINARIES

Let  $S_1$  and  $S_2$  be two strings of lengths  $|S_1|$  and  $|S_2|$ , respectively. An anchor interval pair ([a..b], [c..d]) signifies an exact match between  $S_1[a..b]$  and  $S_2[c..d]$ . For an anchor I, we denote these values as I.a, I.b, I.c, and I.d. Here b-a=d-c and  $S_1[a+j]=S_2[c+j]$  for all  $0 \le j \le b-a$ . We say that the character match  $S_1[a+j]=S_2[c+j]$ ,  $0 \le j \le b-a$ , is *supported* by the anchor ([a..b], [c..d]). Maximal exact matches (MEMs), maximal unique matches (MUMs), or k-mer matches are some of the common ways to define anchors. MUMs (Delcher et al, 1999) are a subset of MEMs, having the added constraint that the pattern involved occurs only once in both strings. If all intervals across all anchors have the same length (e.g., using k-mers), we say that the *fixed-length* property holds.

Our algorithms will make use of dynamic range minimum queries (RmQs). For a set of n d-dimensional points, each with an associated weight, a "query" consists of an orthogonal d-dimensional range. The query response is the point in that range with the smallest weight. Using known techniques in computational geometry, a data structure can be built in  $O(n \log^{d-1} n)$  time and space, which can both answer queries and modify a point's weight in  $O(\log^d n)$  time (de Berg et al, 2008).

## 2.1. Colinear chaining with overlap and gap costs

Given a set of n anchors  $\mathcal{A}$  for strings  $S_1$  and  $S_2$ , we assume that  $\mathcal{A}$  already contains two *end-point* anchors  $\mathcal{A}_{left} = ([0,0], [0,0])$  and  $\mathcal{A}_{right} = ([|S_1|+1,|S_1|+1], [|S_2|+1, |S_2|+1])$ . We define the strict precedence relationship  $\prec$  between two anchors  $I' := \mathcal{A}[j]$  and  $I := \mathcal{A}[i]$  as  $I' \prec I$  if and only if  $I' \cdot a \leq I \cdot a$ ,  $I' \cdot b \leq I \cdot b$ ,  $I' \cdot c \leq I \cdot c$ ,  $I' \cdot d \leq I \cdot d$ , and strict inequality holds for at least one of the four inequalities. In other words, the interval belonging to I' for  $S_1$  (respectively  $S_2$ ) should start before or at the starting position of the interval belonging to I for  $S_1$  (respectively  $S_2$ ) and should not extend past it.

We also define the weak precedence relation  $\prec_w$  as  $I' \prec_w I$  if and only if  $I'.a \le I.a$ ,  $I'.c \le I.c$ , and strict inequality holds for at least one of the two inequalities, that is, intervals belonging to I' should start before or at the starting position of intervals belonging to I, but now intervals belonging to I' can be extended past the intervals belonging to I. The aim of the problem is to find a totally ordered subset (a chain) of A that achieves the minimum cost under the cost function presented next. We specify whether we mean a chain under strict precedence or under weak precedence when necessary.

2.1.1. Cost function. For  $I' \prec I$ , the function connect(I', I) is designed to indicate the cost of connecting anchor I' to anchor I in a chain. The chaining problem asks for a chain of  $m \le n$  anchors,  $\mathcal{A}'[1], \mathcal{A}'[2], \ldots, \mathcal{A}'[m]$ , such that the following properties hold: (i)  $\mathcal{A}'[1] = \mathcal{A}_{left}$ , (ii)  $\mathcal{A}'[m] = \mathcal{A}_{right}$ , (iii)  $\mathcal{A}'[1] \prec \mathcal{A}'[2] \prec \ldots \prec \mathcal{A}'[m]$ , and (iv) the cost  $\sum_{i=1}^{m-1} connect(\mathcal{A}'[i], \mathcal{A}'[i+1])$  is minimized.

We next define the function *connect*. In Section 4, we see that this definition is well motivated by the relationship with anchored edit distance. For a pair of anchors I', I such that  $I' \prec I$ :

- The gap in string  $S_1$  between anchors I' and I is  $g_1 = \max(0, I.a-I'.b-1)$ . Similarly, the gap between the anchors in string  $S_2$  is  $g_2 = \max(0, I.c-I'.d-1)$ . Define the gap cost  $g(I', I) = \max(g_1, g_2)$ .
- The overlap  $o_1$  is defined such that  $I'.b o_1$  reflects the nonoverlapping prefix of anchor I' in string  $S_1$ . Specifically,  $o_1 = \max(0, I'.b - I.a + 1)$ . Similarly, define  $o_2 = \max(0, I'.d - I.c + 1)$ . We define the overlap cost as  $o(I', I) = |o_1 - o_2|$ .
- Lastly, define connect(I', I) = g(I', I) + o(I', I).

The same definitions are used for weak precedence, only using  $\prec_w$  in the place of  $\prec$ . Regardless of the definition of *connect*, the mentioned problem can be trivially solved in  $O(n^2)$  time and O(n) space. First sort the anchors by the component  $\mathcal{A}[\cdot]$ . a and let  $\mathcal{A}'$  be the sorted array. The chaining problem then has a direct dynamic programming solution by filling an n-sized array C from left to right, such that C[i] reflects the cost of an optimal chain that ends at anchor  $\mathcal{A}'[i]$ . The value C[i] is computed using the recursion:

 $C[i] = \min_{\mathcal{A}'[k] \prec \mathcal{A}'[i]} (C[k] + connect(\mathcal{A}'[k], \mathcal{A}'[i]))$ , where the base case associated with anchor  $\mathcal{A}_{left}$  is C[1] = 0. The optimal chaining cost will be stored in C[n] after spending  $O(n^2)$  time. We provide an  $O(n \log^4 n)$  time algorithm for this problem.

#### 2.2. Anchored edit distance

The edit distance problem is to identify the minimum number of operations (substitutions, insertions, or deletions) that must be applied to string  $S_2$  to transform it to  $S_1$ . Edit operations can be equivalently represented as an alignment (a.k.a. edit transcript) that specifies the associated matches, mismatches, and gaps while placing one string on top of another. The *anchored edit distance problem* is as follows: given strings  $S_1$  and  $S_2$  and a set of n anchors  $\mathcal{A}$ , compute the optimal edit distance subject to the condition that a match supported by an anchor has edit cost 0, and a match that is not supported by an anchor has edit cost 1.

The mentioned problem is solvable in  $O(|S_1||S_2|)$  time and space. We can assume that input does not contain redundant anchors, therefore, the count of anchors is  $\leq |S_1||S_2|$ . Next, the standard dynamic programming recursion for solving the edit distance problem can be revised. Let D[i,j] denote anchored edit distance between  $S_1[1,i]$  and  $S_2[1,j]$ , then  $D[i,j] = \min(D[i-1,j-1] + x, D[i-1,j] + 1, D[i,j-1] + 1)$ , where x=0 if  $S_1[i]=S_2[j]$  and the match is supported by some anchor, and x=1 otherwise.

# 2.3. Representing sequence alignment as a graph

It is useful to consider the following representation of an alignment of two strings  $S_1$  and  $S_2$ . As illustrated in Figure 2, we have a set of  $|S_1|$  top vertices and  $|S_2|$  bottom vertices. There are two types of edges between the top and bottom vertices: (1) A solid edge from *i*th top vertex to the *j*th bottom vertex. This represents an anchor supported character match between the *i*th character in  $S_1$  and the *j*th character in  $S_2$ . (2) A dashed edge from the *i*th top vertex to the *j*th bottom vertex. This represents a character being substituted to form a match between  $S_1[i]$  and  $S_2[j]$  or a character match not supported by an anchor. All unmatched vertices are labeled with an "x" to indicate that the corresponding character is deleted. An important observation is that no two edges cross.

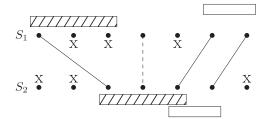
In a solution to the anchored edit distance problem, every solid edge must be "supported" by an anchor. By "supported" here we mean that the match between the corresponding characters in  $S_1$  and  $S_2$  is supported by an anchor. In Figure 2, these anchors are represented with rectangles above and below the vertices. We use  $\mathcal{M}$  to denote a particular alignment. We also associate an edit cost with the alignment, denoted as  $EDIT(\mathcal{M})$ . This is equal to the number of vertices marked with x in  $\mathcal{M}$  plus the number of dashed edges in  $\mathcal{M}$ .

## 3. ALGORITHMS FOR COLINEAR CHAINING

**Theorem 1.** The colinear chaining problem with overlap and gap costs can be solved in time  $\tilde{O}(n)$ . In particular, in time  $O(n \log^2 n)$  for chains with fixed-length anchors, in time  $O(n \log^3 n)$  for chains under weak precedence, and in time  $O(n \log^4 n)$  for chains under strict precedence.

The proposed algorithm still uses the recursive formula given in Section 2.1. However, it uses RmQ data structures to avoid having to check every anchor A[k], where A[k].a < A[i].a. We achieve this by considering six cases concerning the optimal choice of the prior anchor. We use the best of the six distinct possibilities to determine the optimal C[i] value. This C[i] value is then used to update the RmQ data

**FIG. 2.** The graph representation of an alignment. Solid edges represent anchor-supported character matches, dashed edges represent substitutions and unsupported matches, and x's represent deletions. We use  $\mathcal{M}$  to denote an alignment. Here  $EDIT(\mathcal{M}) = 7$ , the total number of x's and dashed edges.



structures. For the strict precedence case, some of the six cases require up to four dimensions for the RmQs. When only weak precedence is required, we reduce this to at most three dimensions. When the fixed-length property holds (e.g., *k*-mers), we reduce this to two dimensions.

# 3.1. Algorithm for chains under strict precedence

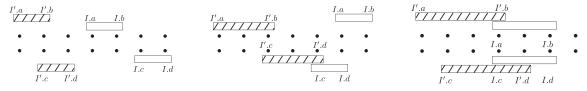
The first step is to sort the set of anchors  $\mathcal{A}$  using the key  $\mathcal{A}[\cdot].a$ . Let  $\mathcal{A}'$  be the sorted array. We will next use six RmQ data structures labeled  $\mathcal{T}_{1a}$ ,  $\mathcal{T}_{1b}$ ,  $\mathcal{T}_{2a}$ ,  $\mathcal{T}_{2b}$ ,  $\mathcal{T}_{3a}$ , and  $\mathcal{T}_{3b}$ . These RmQ data structures are initialized with the following points for every anchor: For anchor  $I \in \mathcal{A}'$ ,  $\mathcal{T}_{1a}$  is initialized with the point (I.b, I.d-I.b),  $\mathcal{T}_{1b}$  with (I.d, I.d-I.b),  $\mathcal{T}_{2a}$  with (I.b, I.c, I.d),  $\mathcal{T}_{2b}$  with (I.b, I.d),  $\mathcal{T}_{3a}$  with (I.b, I.c, I.d, I.d-I.b), and  $\mathcal{T}_{3b}$  with (I.b, I.d, I.d-I.b).

All weights are initially set to  $\infty$  except for  $I = \mathcal{A}_{left}$ , where the corresponding points are given weight 0. We then process the anchors in sorted order and update the RmQ data structures after each iteration. On the ith iteration, for j < i, we let C[j] be the optimal colinear chaining cost of any ordered subset of  $\mathcal{A}'[1]$ ,  $\mathcal{A}'[2]$ , ...,  $\mathcal{A}'[j]$  that ends with  $\mathcal{A}'[j]$ . For i > 1, RmQs are used to find the optimal j < i by considering six different cases. We let  $I = \mathcal{A}'[i]$ ,  $I' = \mathcal{A}'[j]$ , and C[I'] = C[j].

The query for each RmQ structure is determined by the different inequalities relating I.a, I.b, I.c, and I.d to previous anchors in the case considered. For example, in Case 1.a (Fig. 3), it can be seen that I'.b < I.a and I.a-I'.b < I.c-I'.d, making  $I'.b \in [0, I.a-1]$  and  $I'.d-I'.b \in [-\infty, I.c-I.a]$ , motivating the query input  $[0, I.a-1] \times [-\infty, I.c-I.a]$ . At the same time, the values stored in these RmQ structures are determined by the expression for the colinear chaining cost in that case, C[I'] + I.c-I'.d-1.

Note that the values stored in each RmQ structure depend only on previously processed anchors and are combined with the values I.a, I.b, I.c, and I.d for the current anchor I being processed to obtain the appropriate cost. Hence, for  $\mathcal{T}_{1a}$  we store values of the form C[I'] - I'.d and combine this with I.c to obtain the cost. The other cases can be similarly analyzed.

- 1. Case: I' disjoint from I.
- (a) Case: The gap in  $S_1$  is less or equal to gap in  $S_2$  (Fig. 3 Left). The RmQ (query input) is of the form:  $[0, I.a 1] \times [-\infty, I.c I.a]$ . Let the query response (weight) from  $\mathcal{T}_{1a}$  be  $v_{1a} = \min\{C[I'] I'.d: (I'.b, I'.d I'.b) \in [0, I.a 1] \times [-\infty, I.c I.a]\}$  and let  $C_{1a} = v_{1a} + I.c 1$ .
- (b) Case: The gap in  $S_2$  is less than gap in  $S_1$ . The RmQ is of the form  $[0, I.c-1] \times [I.c-I.a+1, \infty]$ . Let the query response from  $\mathcal{T}_{1b}$  be  $v_{1b} = \min\{C[I'] I'.b : (I'.d, I'.d-I'.b) \in [0, I.c-1] \times [I.c-I.a+1, \infty]\}$  and let  $C_{1b} = v_{1b} + I.a-1$ .
  - 2. Case: I' and I overlap in only one dimension.
- (a) Case: I' and I overlap only in  $S_2$  (Fig. 3 Middle). The RmQ is of the form  $[0, I.a-1] \times [0, I.c] \times [I.c, I.d]$ . Let the query response from  $\mathcal{T}_{2a}$  be  $v_{2a} = \min\{C[I'] I'.b + I'.d : (I'.b, I'.c, I'.d) \in [0, I.a-1] \times [0, I.c] \times [I.c, I.d]\}$  and let  $C_{2a} = v_{2a} + I.a I.c$ .
- (b) Case: I' and I overlap only in  $S_1$ . Since the anchors are sorted on  $\mathcal{A}[\cdot].a$ , this can be done with a two-dimensional RmQ structure. The RmQ is of the form  $[I.a, I.b] \times [0, I.c 1]$ . Let the query response from  $\mathcal{T}_{2b}$  be  $v_{2b} = \min\{C[I'] + I'.b I'.d : (I'.b, I'.d) \in [I.a, I.b] \times [0, I.c 1]\}$  and let  $C_{2b} = v_{2b} + I.c I.a$ .
  - 3. Case: I' and I overlap in both dimensions.
- (a) Case: Greater overlap in  $S_2$  (Fig. 3 Right). Here,  $|o_1 o_2| = o_2 o_1 = l'.d l.c (l'.b l.a)$ . The RmQ is of the form  $[I.a, I.b] \times [0, I.c] \times [I.c, I.d] \times [I.c I.a + 1, \infty]$ . Let the query response from  $\mathcal{T}_{3a}$  be  $v_{3a} = \min\{C[l'] l'.b + l'.d : (l'.b, l'c, l'.d, l'.d l'.b) \in [I.a, I.b] \times [0, I.c] \times [I.c, I.d] \times [I.c I.a + 1, \infty]\}$  and let  $C_{3a} = v_{3a} + I.a I.c$ .



**FIG. 3.** (Left) Case 1.a. Colinear chaining cost is  $C[I'] + g_2 = C[I'] + I.c - I'.d - 1$ . (Middle) Case 2.a. Chaining cost is  $C[I'] + g_1 + o_2 = C[I'] + I.a - I'.b + I'.d - I.c$ . (Right) Case 3.a. Chaining cost is  $C[I'] + o_2 - o_1 = C[j] + I'.d - I.c - (I'.b - I.a)$ .

(b) Case: Greater or equal overlap in  $S_1$ . Here,  $|o_1 - o_2| = o_1 - o_2 = I'.b - I.a - (I'.d - I.c)$ . If  $o_1 \ge o_2 > 0$ ,  $I'.b \in [I.a, I.b]$ , and  $I'.a \in [0, I.a]$ , then  $I'.c \in [0, I.c]$ . Hence, the RmQ is of the form  $[I.a, I.b] \times [I.c, I.d] \times [-\infty, I.c - I.a]$ . Let the query response from  $\mathcal{T}_{3b}$  be  $v_{3b} = \min\{C[I'] + I'.b - I'.d : (I'.b, I'.d, I'.d - I'.b) \in [I.a, I.b] \times [I.c, I.d] \times [-\infty, I.c - I.a]\}$  and let  $C_{3b} = v_{3b} - I.a + I.c$ .

Finally, let  $C[i] = \min(C_{1a}, C_{1b}, C_{2a}, C_{2b}, C_{3a}, C_{3b})$  and update the RmQ structures as shown in the pseudo code in Algorithm 1. In the pseudo code, every RmQ structure  $\mathcal{T}$  has the query method  $\mathcal{T}.RmQ()$  that takes as arguments an interval for each dimension. It also has the method  $\mathcal{T}.update()$ , which takes a point and a weight and updates the point to have the new weight. The four-dimensional RmQ structures for Case 3.a require  $O(\log^4 n)$  time per query and update, causing an overtime complexity that is  $O(n \log^4 n)$ .

# 3.2. Modifications for weak precedence and fixed-length anchors

We first consider the case of weak precedence. In Case 3.a, the anchor end I'.d can be positioned arbitrarily to the right of I.c. Moreover, since by the first dimension of the RmQ there is positive overlap in  $S_1$  and by the fourth dimension there is greater overlap in  $S_2$ , we know that  $I'.d \ge I.c$ . Hence, we can then remove the third dimension from the RmQ. The query will then be of the form  $(I'.b, I'.c, I'.d - I'.b) \in [I.a, \infty] \times [0, I.c] \times [I.c - I.a + 1, \infty]$ .

In Case 3.b, where there is greater or equal overlap in  $S_1$ , we can similarly ignore I'.b, but to match our definition of weak precedence we must also ensure  $I'.c \in [0, I.c - 1]$  (this is unnecessary for I'.a in Case 3.a as the strictly greater overlap in  $S_2$  ensures I'.a < I.a). We modify the query to be of the form  $(I'.c, I'.d, I'.d - I'.b) \in [0, I.c - 1] \times [I.c, \infty] \times [-\infty, I.c - I.a]$ . Since each RmQ has at most three dimensions, the total time complexity can be brought down to  $O(n \log^3 n)$ .

**Algorithm 1:**  $O(n \log^4 n)$  time algorithm for strict precedence.

```
Input: n anchors A[1, n] including A_{\text{left}} = A[1] and A_{\text{right}} = A[n].
Output: Array C[1, n] such that C[i] is the optimal colinear chaining cost for any ordered subset of A[1, i] ending
Let \mathcal{A}'[1], \ldots, \mathcal{A}'[n] be the anchors \mathcal{A} sorted on \mathcal{A}[\cdot].a;
Construct RmQ structures with weights set to \infty;
Initialize array C of size n to 0;
for i \leftarrow 1 to n do
   I \leftarrow \mathcal{A}'[i];
   if i \geq 2 then
      C_{1a} \leftarrow \mathcal{T}_{1a}. RmQ([0, I.a - 1], [-\infty, I.c - I.a]) + I.c - 1;
      C_{1b} \leftarrow \mathcal{T}_{1b}.RmQ([0,I.c-1],[I.c-I.a+1,\infty]) + I.a-1;
      C_{2a} \leftarrow T_{2a}.RmQ([0, I.a-1], [0, I.c], [I.c, I.d]) + I.a-I.c;
      C_{2b} \leftarrow \mathcal{T}_{2b}.RmQ([I.a, I.b], [0, I.c-1]) - I.a + I.c;
      C_{3a} \leftarrow T_{3a}.RmQ([I.a, I.b], [0, I.c], [I.c, I.d], [I.c-I.a+1, \infty]) + I.a-I.c;
      C_{3b} \leftarrow T_{3b}.RmQ([I.a, I.b], [I.c, I.d], [-\infty, I.c-I.a]) - I.a + I.c;
      /* Take optimal choice */
     C[i] \leftarrow \min(C_{1a}, C_{1b}, C_{2a}, C_{2b}, C_{3a}, C_{3b});
   /* Update RmQ structures */
   \mathcal{T}_{1a}.update((I.b, I.d-I.b), C[i]-I.d);
   \mathcal{T}_{1b}.update((I.d, I.d-I.b), C[i]-I.b);
   T_{2a}.update((I.b, I.c, I.d), C[i] - I.b + I.d);
   \mathcal{T}_{2b}.update((I.b, I.d), C[i] + I.b - I.d);
   T_{3a}.update((I.b, I.c, I.d, I.d-I.b), C[i]-I.b+I.d);
   T_{3b}.update((I.b, I.d, I.d-I.b), C[i]+I.b-I.d);
end
return C[1, n]
```

In the case of fixed-length anchors, the RmQ for Case 2.a. can be made  $(I'.b, I'.d) \in [0, I.a-1] \times [I.c, I.d]$ . The modifications for Cases 3.a and 3.b are more involved. We keep a pointer  $p_a$  to indicate the current a value of the interval, initially setting  $p_a = \mathcal{A}[1] \cdot a$ . Conceptually, before processing anchor I we increment  $p_a$  from its previous position to I.a. If for some anchor I' the end I'.b is passed by  $p_a$ , we

update the points associated with I' in  $\mathcal{T}_{3a}$  and  $\mathcal{T}_{3b}$  to have the weight  $\infty$ . This eliminates the need to use a range query to check  $I'.b \in [I.a, I.b]$ , since any points not within that range are effectively removed from consideration.

Hence, we can reduce the query for Case 3.a. (overlap in  $S_2$  greater than overlap in  $S_1$ ) to  $(I'.c, I'.d - I'.b) \in [0, I.c] \times [I.c - I.a + 1, \infty]$ . and the query for Case 3.b (overlap in  $S_1$  greater or equal to overlap in  $S_2$ ) to  $(I'.d, I'.d - I'.b) \in [I.c, I.d] \times [-\infty, I.c - I.a]$ . To avoid the  $|S_1|$  time complexity, the anchors that would be encountered while incrementing  $p_a$  can be found by looking at which anchors have b values between the previous  $p_a$  value and I.a. Because each update requires  $O(\log^2 n)$  time, these updates cost time  $O(n \log^2 n)$  in total.

# 4. PROOF OF EQUIVALENCE

**Theorem 2.** For a fixed set of anchors A, the following quantities are equal: the anchored edit distance, the optimal colinear chaining cost under strict precedence, and the optimal colinear chaining cost under weak precedence.

The optimal colinear chaining cost is defined using the cost function described in Section 2.1. An implication of Theorems 1 and 2 is that if only the anchored edit distance is desired (and not an optimal strictly ordered anchor chain), there exists a  $O(n \log^3 n)$  for computing this value.

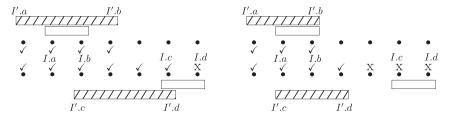
Theorem 2 will follow from Lemmas 1 and 2.

**Lemma 1.** Anchored edit distance  $\leq$  optimal colinear chaining cost under weak precedence  $\leq$  optimal colinear chaining cost under strict precedence.

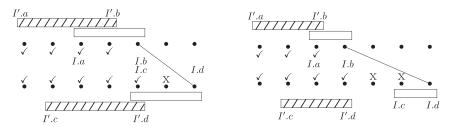
*Proof.* The second inequality follows from the observation that every set of anchors ordered under strict precedence is also ordered under weak precedence. We now focus on the inequality between anchored edit distance and colinear chaining cost under weak precedence. Starting with an anchor chain under weak precedence, A[1], A[2], ... with associated colinear chaining cost x, we provide an alignment with an anchored edit distance that is at most x. This alignment is obtained using a greedy algorithm that works from left to right, always taking the closest exact match when possible, and when not possible, a character substitution or unsupported exact match, or if none of these are possible, a deletion. We now present the details.

**The Greedy Algorithm.** Assume inductively that all symbols in  $S_1[1, A[i].b]$  and  $S_2[1, A[i].d]$  have been processed, that is, either matched, substituted, or deleted (represented by check marks in Figs. 4–6). The base case of this induction holds trivially for  $A_{left}$ . We consider the anchor A[i+1] and the possible cases regarding its position relative to A[i]. Symmetric cases that only swap the roles of  $S_1$  and  $S_2$  are ignored. To ease notation, let I' = A[i] and I = A[i+1].

- 1. Case  $I'.b \ge I.b$  and  $I'.d \ge I.c$  (Fig. 4): To continue the alignment, delete the substring  $S_2[I'.d+1, I.d]$  from  $S_2$ . This has edit cost I.d-I'.d. We can assume both intervals of I' are not nested in intervals of I, hence  $connect(I', I) = o_1 o_2 = I'.b I.a I'.d + I.c \ge I.c + I.b I.a I'.d = I.d I'.d$ .
- 2. **Case**  $I'.b \ge I.b$  and I'.d < I.c (Fig. 4): Delete the substring  $S_2[I'.d+1, I.d]$  from  $S_2$ , with edit cost I.d-I'.d. Also  $connect(I', I) = o_1 + g_2 = I'.b I.a + I.c I'.d \ge I.c + I.b I.a I'.d = I.d I'.d$ .
- 3. **Case** I.b > I'.b,  $I.a \le I'.b$ ,  $I.c \le I'.d$  (Fig. 5): Supposing wlog that  $o_1 > o_2$ , delete  $S_2[I'.d+1, I'.d+o_1-o_2]$ , and match  $S_1[I'.b+1, I.b]$  and  $S_2[I'.d+o_1-o_2+1, I.d]$ . This has edit cost  $o_1-o_2$  and  $connect(I', I) = o_1 o_2$ .
- 4. **Case** I.b > I'.b,  $I.a \le I'.b$ , I.c > I'.d (Fig. 5): We delete  $S_2[I'.d+1, I'.d+o_1+g_2]$  and match  $S_1[I'.b+1, I.b]$  with  $S_2[I'.d+o_1+g_2+1, I.d]$ . This has edit cost  $o_1+g_2$  and  $connect(I', I)=o_1+g_2$ .



**FIG. 4.** Cases in Proof of Lemma 1. The  $\checkmark$  symbol indicates symbols processed before considering *I*. (Left) Case  $I', b \ge I.b$  and  $I'.d \ge I.c$   $I'.b \ge I.b$  and I'.d < I.c.



**FIG. 5.** Cases in proof of Lemma 1. (Left) Case I.b > I'.b,  $I.a \le I'.b$ ,  $I.c \le I'.d$ . (Right) Case I.b > I'.b,  $I.a \le I'.b$ , I.c > I'.d.

5. Case I.a > I'.b, I.c > I'.d (Fig. 6): Supposing wlog  $g_2 \ge g_1$ , match with substitutions or unsupported exact matches  $S_1[I'.b+1, I'.b+g_1]$  and  $S_2[I'.d+1, I'.d+g_1]$ . Delete the substring  $S_2[I'.d+g_1+1, I.c-1]$ . Finally, match  $S_1[I.a, I.b]$  and  $S_2[I.c, I.d]$ . The edits consist of  $g_1$  of substitutions or unsupported exact matches and  $g_2 - g_1$  deletions, which is  $g_2$  edits in total. Also,  $connect(I', I) = max\{g_1, g_2\} = g_2$ .

Continuing this process until  $A_{\text{right}}$ , all symbols in  $S_1$  and  $S_2$  become included in the alignment.  $\square$  We delay the details of Lemma 2's proof to Section 4.1.

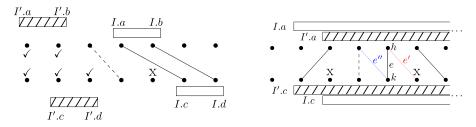
**Lemma 2.** For a set of anchors A, optimal chaining cost under strict precedence  $\leq$  anchored edit distance.

*Proof.* We start with an arbitrary alignment  $\mathcal{M}$  supported by  $\mathcal{A}$ . We show in Lemma 3 how to obtain a subset  $\mathcal{B} \subseteq \mathcal{A}$  totally ordered under strict precedence and supporting an alignment  $\mathcal{M}'$  where  $EDIT(\mathcal{M}') \leq EDIT(\mathcal{M})$ . We then show in Lemma 4 that the edit cost of  $\mathcal{M}'$  is greater or equal to the edit cost of the alignment  $\mathcal{M}_G$  given by the greedy algorithm on  $\mathcal{B}$ . Finally, in Lemma 5 we show that the colinear chaining cost of  $\mathcal{B}$  is equal to the edit cost of  $\mathcal{M}_G$ . Combining, we have  $EDIT(\mathcal{M}) \geq EDIT(\mathcal{M}') \geq EDIT(\mathcal{M}_G) =$  the colinear chaining cost on  $\mathcal{B} \geq$  optimal colinear chaining cost under strict precedence for  $\mathcal{A}$ . The result follows from the fact that  $EDIT(\mathcal{M})$  equals the anchored edit distance when  $\mathcal{M}$  is an optimal alignment for  $\mathcal{A}$ .

#### 4.1. Details of Lemma 2 proof

We apply Algorithm (i) followed by Algorithm (ii) to convert a supporting set of anchors  $\mathcal{A}$  for  $\mathcal{M}$  into the totally ordered subset of anchors  $\mathcal{B}$  supporting  $\mathcal{M}'$ . Note that these algorithms are only for the purpose of the proof. Moving forward, we call an edge  $e = (S_1[h], S_2[k])$  contained but not supported by I if  $h \in [I.a, I.b]$  or  $k \in [I.c, I.d]$  and  $h-I.a \neq k-I.c$ . We define for e the two edges  $e' = (S_1[h], S_2[I.c + h - I.a])$  and  $e'' = (S_1[I.a + k - I.c], S_2[k])$ , which are supported by I.

Algorithm (i). Algorithm for removing incomparable anchors. Let I and I' be two incomparable anchors under weak precedence (Fig. 6). The anchor that has the rightmost supported solid edge will be the anchor we keep. Suppose wlog is I. Working from right to left, starting with that rightmost edge, for any edge e that is contained but not supported by I, we replace e with the rightmost of e' and e''. Note that at least one side of every edge supported by I' is within an interval of I. Hence, all edges supported by I' are eventually replaced. We then remove I'. This algorithm is repeated until a total ordering under weak precedence is possible.



**FIG. 6.** (Left) Case I.a > l'.b, I.c > l'.d. (Right) Anchors I and I' are incomparable. The current alignment is shown with black solid and dashed edges. To remove I' we sweep from right to left, replacing edges not supported by I with edges supported by I. Here,  $e = (S_1[h], S_2[k])$  is not supported by I and will be replaced with  $e' = (S_1[h], S_2[I.c + h - I.a])$  (in red), which is supported by I.

Algorithm (ii). Algorithm for removing anchors with nested intervals. Consider two anchors I and I' where wlog I' has an interval nested in one of the intervals belonging to I. Let  $e_R$  be the rightmost edge supported by I. Working from right to left, we replace any edge e to the left of  $e_R$  that is contained but not supported by I with the rightmost of e' and e''. Next, working from left to right, we replace any edge e to the right of  $e_R$  that is contained but not supported by I with the leftmost of e' and e''. These procedures combined will replace all edges supported by I' with those supported by I. We repeat this until there are no two nested intervals among all remaining anchors. Finally, remove all anchors that do not support any edge. We call such an anchor chain where every anchor supports at least one edge minimal.

#### **Lemma 3.** $EDIT(\mathcal{M}') \leq EDIT(\mathcal{M})$ .

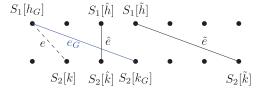
*Proof.* For Algorithm (i), suppose we are replacing an edge e not supported by anchor I, the anchor we wish to keep. Suppose wlog that e' is the rightmost of e' and e'', so we replace e with e'. Because the edge immediately to the right of e is also aligned with I, deleting  $S_2[k]$  and matching  $S_2[I.c + h - I.a]$  does not require modifying any additional edges. If e was a solid edge, the edit cost is unaltered, since the total number of deletions and matches is unaltered. If e was a dashed edge, replacing e with e' converts a substitution or unsupported exact match at  $S_2[k]$  to a deletion, and removes a deletion at  $S_2[I.c + h - I.a]$ , decreasing the edit cost by 1.

The same arguments hold for Algorithm (ii) when we replace edges from right to left. In Algorithm (ii) when we process edges from left to right, since any edges to left of the edge e being replaced are supported by I, replacing e with the leftmost of e' and e'' does not require modifying any additional edges. Again, if e is solid, the edit cost is unaltered, and if e is dashed, the edit cost is decreased by 1.

**Lemma 4.** The greedy algorithm described in the proof of Lemma 1 produces an optimal alignment for a "minimal" anchor chain under strict precedence.

*Proof.* This follows from an exchange argument. Let  $\mathcal{M}$  be an optimal alignment on the anchor chain  $\mathcal{B}$  Further, suppose that  $\mathcal{M}$  is not the same as the alignment produced by the greedy algorithm  $\mathcal{M}_G$ . As we process the edges from left to right, consider when the first discrepancy in the edges is found, the leftmost edges  $e = (S_1[h], S_2[k])$  in  $\mathcal{M}$  and  $e_G = (S_1[h_G], S_2[k_G])$  in  $\mathcal{M}_G$  that are not equal. Let  $e_{\text{prev}}$  be the previous edge on which the  $\mathcal{M}$  and  $\mathcal{M}_G$  coincided. We claim e can be replaced with  $e_G$  without increasing the edit cost.

- Case:  $e_G$  is solid and left of e. Then e can clearly be exchanged for  $e_G$  with no increase in edit cost.
- Case:  $e_G$  is solid and not left of e. We assume WLOG  $k < k_G$  (Fig. 7). Suppose  $e_G$  is supported by an anchor I. We claim that (i)  $h \ge h_G$ , and (ii) there exists an edge  $\tilde{e} = (S_1[\tilde{h}], S_2[\tilde{k}])$  in  $\mathcal{M}$  to the right of e and supported by I. Note that  $\tilde{h} h_G + 1$  solid edges can be obtained by matching  $S_1[h_G, \tilde{h}]$  to  $S_2[k_G, \tilde{k}]$ , and that this alignment has at most the edit cost incurred within the intervals  $[h_G, \tilde{h}]$  in  $S_1$  and  $[k, \tilde{k}]$  in  $S_2$  by  $\mathcal{M}$ . Thanks to (i) and (ii), swapping all of the edges in  $\mathcal{M}$  inside these intervals with the  $\tilde{h} h_G + 1$  edges supported by I already mentioned can be done without effecting edges outside those intervals. The proofs for claims (i) and (ii) are shown hereunder.
- (i)  $h \ge h_G$ : Suppose to the contrary that  $h < h_G$ . Let  $e_{\text{prev}} = (S_1[h_{\text{prev}}], S_2[k_{\text{prev}}])$ . Since  $h_{\text{prev}} < h < h_G$  and  $k_{\text{prev}} < k < k_G$ , this would imply that there are deletions (x's) at both  $S_1[h]$  and  $S_2[k]$  in  $\mathcal{M}_G$ . However, the greedy algorithm would instead create an edge including  $S_1[h]$  or  $S_2[k]$  (or both), causing there to be an edge to the left of  $e_G$  and to the right of  $e_{\text{prev}}$ . This contradicts our assumption that  $e_G$  is the first edge in  $\mathcal{M}_G$  to the right of  $e_{\text{prev}}$ .
- (ii) An edge  $\tilde{e} = (S_1[\tilde{h}], S_2[\tilde{k}])$  exists in  $\mathcal{M}$  to the right of e and supported by I: We claim  $e_G$  is the leftmost edge supported by I in  $\mathcal{M}_G$ . Otherwise, since only consecutive edges supported by an anchor are taken in the greedy algorithm (Cases 3, 4, and 5 in proof of Lemma 1),  $e_{\text{prev}} = (S_1[h_G 1], S_2[k_G 1])$ . Then  $e_{\text{prev}}$  was also in  $\mathcal{M}$  by our assumption that  $\mathcal{M}_G$  and  $\mathcal{M}$  agreed



**FIG. 7.** Consider a greedy alignment  $\mathcal{M}_G$  and optimal alignment  $\mathcal{M}$  where the leftmost difference is e and  $e_G$ . Here,  $\hat{e}$  prevents swapping e and  $e_G$  without modifying additional edges in  $\mathcal{M}$ 

until  $e_G$ . However,  $e_{\text{prev}} = (S_1[h_G - 1], S_2[k_G - 1])$  would cross or share a vertex with  $e = (S_1[h], S_2[k])$  in  $\mathcal{M}$  since  $h_G - 1 < h$  and  $k \le k_G - 1$ , a contradiction. Now, if no other edges supported by I are in  $\mathcal{M}$ , then  $\mathcal{B}$  is not minimal. Hence, we can assume that an edge  $\tilde{e} = (S_1[\tilde{h}], S_2[\tilde{k}])$  exists in  $\mathcal{M}$  to the right of e and supported by I.

• Case:  $e_G$  is a dashed edge. Since dashed edges are only created in Case 5 of the proof of Lemma 1 (Fig. 6),  $e_{prev} = (S_1[h_G - 1], S_2[k_G - 1])$  must have been either the rightmost supported edge for some anchor, or another dashed edge. In addition, no other supported edges are possible that include either  $S_1[h]$  or  $S_2[k]$ . Here the dashed edge  $e_G$  is optimal and swapping e with  $e_G$  will reduce the edit cost.

**Lemma 5.** For an anchor chain under strict precedence, the edit cost of the alignment produced by the greedy algorithm described in the proof of Lemma 1 is equal to the chaining cost.

*Proof.* This follows from induction on the number of anchors processed, using the same arguments used in the proof of Lemma 1. However, only I'.b=I.b needs to be considered in Cases 1 and 2 leading to equality in these cases.

#### 5. IMPLEMENTATION

In multidimensional RmQs,  $O(n \log^{d-1} n)$  storage requirement and irregular memory access during a query can limit their efficacy in practice (de Berg et al, 2008). We can take advantage of two observations to design a more practical algorithm. First, if sequences are highly similar, their edit distance will be relatively small. Hence the anchored edit distance, denoted in this section as OPT, will be relatively small for MUM or MEM anchors. Second, if the sequences are dissimilar, then the number of MUM or MEM anchors, n, will likely be small.

These observations allow us to design an alternative algorithm (Algorithm 1) that requires O(n) worst-case space and  $O(n \cdot OPT + n \log n)$  average case time over all possible inputs where  $n \leq \max(|S_1|, |S_2|)$ , that is, the number of anchors is less than the longer sequence length. This is stated formally in Lemma 6. This property always holds when the anchors are MUMs and is typically true for MEMs as well. This makes the algorithm presented here a practical alternative.

As before, let  $\mathcal{A}$  be the initial (possibly unsorted) set of anchors, but with  $\mathcal{A}_{left} = \mathcal{A}[1]$  and  $\mathcal{A}_{right} = \mathcal{A}[n]$ . We assume wlog  $|S_1| \geq |S_2|$ . We begin by sorting anchor set  $\mathcal{A}$  by the component  $\mathcal{A}[\cdot]$ . a and making a guess for the optimal solution, B (Algorithm 2). The value B is used at every step to bound the range of  $\mathcal{A}[\cdot]$ . a values that need to be examined. This bounds the number of anchors that need to be considered (on average). If C[n] is greater than our current guess B after processing all B anchors, we update our guess to  $B_2 \cdot B$ .

**Algorithm 2:**  $O(OPT \cdot n + n \log n)$  average case algorithm.

```
Input: n anchors \mathcal{A} and parameters B_1 and B_2.

Output: C[1, n] such that C[i] is optimal colinear chaining cost for any ordered subset of \mathcal{A}[1, i] ending at \mathcal{A}[i]. Let \mathcal{A}'[1], ... \mathcal{A}'[n] be the set of anchors \mathcal{A} sorted on \mathcal{A}[\cdot].a; Initialize array C of size n to 0 and B \leftarrow B_1;

do

\begin{vmatrix}
j \leftarrow 1; \\
\text{for } i \leftarrow 1 \text{ to } n \text{ do} \\
& | j \leftarrow j + 1; \\
& \text{end} \\
& C[i] \leftarrow \min\{C[k] + connect(\mathcal{A}'[k], \mathcal{A}'[i])|j \leq k < i \text{ and } \mathcal{A}'[k] \prec \mathcal{A}'[i]\} \\
& \text{end} \\
& B_{\text{last}} \leftarrow B; \\
& B \leftarrow B_2 \cdot B; \\
& \text{while } C[n] > B_{\text{last}}; \\
& \text{return } C[1, n]
```

Extending the mentioned pseudo code to enable semiglobal chaining, that is, free anchor gap on both ends of reference sequences, is also simple. In each i-loop, the connection to anchor  $\mathcal{A}_{left}$  must be always considered, and for last iteration when i=n, j must be set to 1. Second, a revised cost function must be used when connecting to either  $\mathcal{A}_{left}$  or  $\mathcal{A}_{right}$ , where a gap penalty is used only for anchor gap over the query sequence. The experiments in the next section use an implementation of this algorithm.

**Lemma 6.** Algorithm 2 runs in  $O(n \cdot OPT + n \log n)$  average case time over all inputs where  $n \le \max(|S_1|, |S_2|)$ .

*Proof.* The  $n \log n$  term is from the sorting the anchors. To analyze the second portion of the algorithm, we first let  $X_{h,j}$  be 1 if  $\mathcal{A}[h].a$  is located at index j in  $S_1$ . Under the assumption of a random placement of anchors,  $\mathbb{E}[X_{h,j}] = 1/|S_1|$ . Let  $X_i$  be the number of anchors,  $\mathcal{A}[h]$ , where  $\mathcal{A}[h].a \in [\mathcal{A}[i].a - B, \mathcal{A}[i].a - 1]$ . We have that  $X_i = \sum_{h=1}^n \sum_{j=\mathcal{A}[i].a-B}^{\mathcal{A}[i].a-1} X_{h,j}$ . Letting X be the total number of anchors processed,  $X = \sum_{i=1}^n X_i$  and

$$\mathbb{E}[X] = \sum_{i=1}^{n} \sum_{h=1}^{n} \sum_{i=A[i]}^{A[i]} \sum_{a=B}^{a-1} \mathbb{E}[X_{h,j}] = \frac{n^2 \cdot B}{|S_1|} \leq nB.$$

The total expected time is a constant factor from  $B_1 n(1 + B_2 + ... + B_2^{\lceil \log_{B_2} OPT \rceil}) = O(n \cdot OPT)$ .  $\square$ 

## 6. EXPERIMENTAL EVALUATION

There are multiple open-source libraries/tools that implement edit distance computation. Edlib (v1.2.7) (Šošić and Šikić, 2017) uses Myers's bit-vector algorithm (Myers, 1999) and Ukkonen's banded algorithm (Ukkonen, 1985), and is known to be the fastest implementation currently. In this section, we aim to show that (1) the proposed algorithm as well as existing chaining methods achieves significant speedup compared with computing exact edit distance using Edlib, and (2) in contrast to existing chaining methods, our implementation consistently achieves high Pearson correlation (> 0.90) with edit distance while requiring modest time and memory resources.

We implemented Algorithm 2 in C++, and refer to it as ChainX. The code is available at (https://github.com/at-cg/ChainX). Inputs are a target string, query strings, comparison mode (global or semi-global), anchor type preferred, that is, MUMs or MEMs, and a minimum match length. We include a preprocessing step to index target string using the same suffix array-based algorithm (Vyverman et al, 2013) used in Nucmer4 (Marçais et al, 2018). Chaining costs computed using ChainX for each query-target pair are provably optimal.

# Existing colinear chaining implementations

Colinear chaining has been implemented previously as a stand-alone utility (Abouelhoda et al, 2008; Otto et al, 2011) and also used as a heuristic inside widely used sequence aligners (Bray et al, 2003; Li, 2018; Marçais et al, 2018). Out of these, Clasp (v1.1), Nucmer4 (v4.0.0rc1), and Minimap2 (v2.22-r1101) tools are available as open source, and used here for comparison purpose. Unlike our algorithm where the optimization problem involves minimizing a cost function, these tools execute their respective chaining algorithms using a score maximization objective function. Clasp, being a stand-alone chaining method, returns chaining scores in its output, whereas we modified Minimap2 and Nucmer4 to print the maximum chaining score for each query–target string pair, and skip subsequent steps.

To enable a fair comparison, all methods were run with single thread and same minimum anchor size 20. Accordingly, ChainX, Clasp, and Nucmer4 were run with MUMs of length  $\geq$  20, and Minimap2 was configured to use minimizer k-mers of length 20. For these tests, we made use of an Intel Xeon E5-2698 v3 processor with 32 cores and 128 GB RAM. All tools were required to match only the forward strand of each query string. ChainX and Clasp are both exact solvers of colinear chaining problem, but use different gap cost functions. Clasp only permits nonoverlapping anchors in a chain, and supports two cost functions that were referred to as *sum-of-pair* and *linear* gap cost functions in their article (Otto et al, 2011).

We tested Clasp with both of its gap cost functions, and refer to these two versions as Clasp-sop and Clasp-linear, respectively. Both the versions solve colinear chaining using RmQ data structures,

Table 1. Runtime and Memory Usage Comparison of Edit Distance Solver Edlib and Colinear Chaining Methods ChainX, Clasp, Nucmer4, and Minimap2

Similarity	No. of MUMs	Edlib	ChainX	Clasp-sop	Clasp-linear	Nucmer4	Minimap2
		Time (Mem)	Time (Mem)	Time (Mem)	Time (Mem)	Time (Mem)	Time (Mem)
Semigloba	l pairwise	e sequence com	parisons, seque	ence sizes $10^4 \times 5$	* 10 <sup>6</sup>		
99%	67	190 ( <b>17</b> )	2.0 (57)	1.8 (57)	0.9 (57)	1.8 (60)	1.9 (75)
97%	160	642 ( <b>17</b> )	2.9 (57)	4.8 (57)	1.8 (57)	4.1 (60)	2.3 (75)
94%	176	1165 ( <b>17</b> )	3.0 (57)	5.9 (57)	2.1 (57)	3.2 (60)	1.6 (75)
90%	135	2168 ( <b>17</b> )	5.6 (57)	4.7 (57)	2.0 (57)	5.5 (60)	<b>1.9</b> (75)
80%	28	2360 (17)	4.2 (57)	2.5 (57)	<b>2.2</b> (57)	3.4 (60)	4.3 (75)
70%	3	4297 ( <b>17</b> )	3.7 (57)	2.2 (57)	2.3 (57)	5.5 (60)	1.1 (75)
Global pai	rwise seq	uence comparis	sons, sequence	sizes $10^6 \times 10^6$			
99%	7012	949 (8)	<b>47.2</b> (24)	1236.8 (1800)	182.8 (257)	68.7 (26)	193.5 (35)
97%	15, 862	1308 ( <b>8</b> )	490.4 (24)	5363.7 (8742)	765.4 (1278)	<b>87.8</b> (26)	179.0 (36)
94%	18, 389	2613 <b>(8</b> )	677.9 (24)	11737.1 (20, 501)	1021.0 (1694)	113.5 (27)	116.9 (34)
90%	14, 472	6233 <b>(8</b> )	851.5 (24)	5110.3 (8277)	115.3 (27)	121.8 (26)	<b>94.8</b> (33)
80%	2964	12, 506 ( <b>8</b> )	158.8 (24)	504.8 (572)	133.7 (24)	148.9 (26)	<b>69.5</b> (32)
70%	195	29, 602 ( <b>8</b> )	136.5 (23)	140.6 (23)	139.6 (23)	167.3 (26)	<b>55.6</b> (32)

The best numbers are highlighted in bold.

Runtime is measured in milliseconds across the columns, and memory usage (Mem) is noted in MBs. In this experiment, ChainX, Clasp-sop, Clasp-linear, and Nucmer4 used MUMs of length  $\geq 20$  as input anchors, while Minimap2 used fixed-length minimizer k-mers of size 20.

MUM, maximal unique match.

requiring  $O(n \log^2 n)$  and  $O(n \log n)$  time, respectively. Both require a set of anchors as input, therefore, we supplied the same set of anchors, that is, MUMs of length  $\geq 20$  as used by ChainX. Minimap2 and Nucmer4 use colinear chaining as part of their seed-chain-extend pipelines. Both Minimap2 and Nucmer2 support anchor overlaps in a chain, as well as penalize gaps using custom functions. However, both these tools employ heuristics (e.g., enforce a maximum gap between adjacent chained anchors) for faster processing that can result in suboptimal chaining scores.

# Runtime and memory comparison

We downloaded the same set of query and target strings that were used for benchmarking in Edlib article (Šošić and Šikić, 2017) (https://github.com/Martinsos/edlib/tree/master/test\_data). These test strings, ranging from 10 to 5000 kbp in length, allowed us to compare tools for end-to-end global sequence comparisons as well as semiglobal comparisons at various degrees of similarity levels. To test end-to-end comparisons, the target string had been artificially mutated at various rates using mutatrix (https://github.com/ekg/mutatrix), whereas for the semiglobal comparisons, a substring of the target string had been sampled and mutated.

Table 2. Absolute Pearson Correlation Coefficients of Chaining Costs (or Scores) Computed by Various Methods with the Corresponding Edit Distances

	Similarity	Correlation coefficient					
Sequence sizes		ChainX	Clasp-sop	Clasp-linear	Nucmer4	Minimap2	
Semiglobal sequen	ice comparisons						
$10^{4} \times 5 * 10^{6}$	90%-100%	0.996	0.994	0.986	0.968	0.995	
$10^4 \times 5 * 10^6$	80%-90%	0.975	0.976	0.786	0.864	0.958	
$10^4 \times 5 * 10^6$	75%-80%	0.927	0.915	0.732	0.733	0.808	
Global sequence co	omparisons						
$10^6 \times 10^6$	90%-100%	0.999	0.997	0.994	0.991	0.999	
$10^6 \times 10^6$	80%-90%	0.998	0.998	0.922	0.955	0.996	
$10^6 \times 10^6$	75%-80%	0.992	0.993	0.871	0.907	0.952	

The best numbers are highlighted in bold.

<sup>100</sup> query strings were simulated and matched to the target string within each similarity range.

Using MUMs Using maximal exact matches  $Length \ge 7$ Length ≥20 Length ≥10 Length ≥7 Length ≥20 Length ≥10 (corr.) Similarity Time Time (corr.) Time Time (corr.) Time Time (corr.) (corr.) (corr.) 7.2 90%-100% (0.996)2.9 (0.997)3.5 (0.997)5.1 (0.996)8.1 (0.997)2652 (0.998)80%-90% 4.5 (0.975)(0.992)3.2 (0.992)(0.975)7.4 (0.993)5413 (0.995)5.6 4.5 75%-80% 5.3 (0.927)5.9 (0.977)1.9 (0.977)5.0 (0.927)10.9 (0.987)9221 (0.992)

Table 3. Effect of anchor precomputation method on the performance of ChainX

The best numbers are highlighted in bold.

Total runtime to do 100 pairwise semiglobal sequence comparisons (sequence size:  $10^4 \times 5 * 10^6$ ) is measured in seconds, and correlation (corr.) with the corresponding edit distances is computed using Pearson correlation coefficient.

Table 1 presents runtime and memory comparison of all tools. Columns of the table are organized to show tools in three categories: edit distance solver (Edlib); optimal colinear chaining solvers (ChainX, Clasp-sop, and Clasp-linear); and heuristic implementations (Nucmer4 and Minimap2). We make the following observations here. First, chaining methods (both optimal and heuristic based) are significantly faster than Edlib in most cases, and we see up to three orders of magnitude speedup.

Second, within optimal chaining methods, Clasp-sop's time and memory consumption increases quickly with increase in count of anchors, which is likely due to irregular memory access and storage overhead of its algorithm that uses a 2d-RmQ data structure. Finally, we note that Minimap2 and Nucmer4 are often faster than exact algorithms during global string comparisons due to their fast heuristics.

All tools (except Edlib) use an indexing step such as building a *k*-mer hash table (Minimap2) or computing suffix array (ChainX, Clasp-sop, Clasp-linear, and Nucmer4). Indexing time was excluded from reported results, and was found to be comparable. For instance, in the case of semiglobal comparisons, ChainX, Nucmer4, and Minimap2 required 590, 736, and 236 ms for index computation, respectively.

#### Correlation with edit distance

We checked how well the chaining cost (or score) correlates with edit distance. We use absolute value of Pearson correlation coefficients for a comparison. In this experiment, we simulated 100 query strings within three similarity ranges: 90% - 100%, 80% - 90%, and 75% - 80%. Table 2 gives the correlation achieved by all the tools. Here we observe that ChainX and Clasp-sop are more consistent in terms of maintaining high correlation across all similarity ranges. Between the two, ChainX was shown to offer superior scalability in terms of runtime and memory usage (Table 1). Hence, ChainX can be useful in practice when good performance and accuracy are desired across a wide similarity range.

## Effect of anchor type and minimum match length

How many anchors are given as input will naturally affect the performance and output quality of a chaining algorithm. We tested runtime and correlation with edit distance achieved by ChainX while varying the anchor type (MUMs/MEMs) and minimum match length  $l_{\min}$  parameter (Table 3). When MUMs are used as anchors, we observe good scalability, and lowering  $l_{\min}$  from 20 to 10 improves the correlation, but the correlation saturates afterward. This is because very short exact matches will unlikely be unique and will not be selected as MUMs. However, when MEMs are used as anchors, correlation continues to improve with decreasing minimum length parameter, however, runtime grows exponentially. Excessive count of anchors improves the correlation but then anchor chaining becomes computationally demanding.

#### 7. CONCLUSIONS

This study provides new algorithms for colinear chaining, a fundamental problem in bioinformatics. Variants of this technique have been regularly used in alignment tools since four decades (Wilbur and Lipman, 1983). We addressed an open problem pertaining to the general case of this problem that allows anchor overlaps and penalizes gap cost between adjacent chained anchors. The proposed algorithms for multiple versions of this problem are provably optimal and efficient, and can be incorporated in read

mappers. We also discussed a new cost function for the colinear chaining problem that enabled us to establish the first mathematical link between colinear chaining and the edit distance problem.

This result is a useful addition to a prior result (Mäkinen and Sahlin, 2020) where a connection between the colinear chaining problem and the longest common subsequence problem was made. Although we focused on a single cost function, it may be possible to derive a family of chaining cost functions by linking from weighted edit distance or linear and affine gap alignment scoring functions.

# **AUTHORS' CONTRIBUTIONS**

The authors confirm that all authors contributed equally to the conception and analysis of the ideas presented in this study, as well as the drafting and finalizing of this document.

# **AUTHOR DISCLOSURE STATEMENT**

The authors declare they have no conflicting financial interests.

#### **FUNDING INFORMATION**

This research is supported in part by the U.S. National Science Foundation (NSF) grants CCF-1704552, CCF-1816027, CCF-2112643, and CCF-2146003, and funding from the Indian Institute of Science. A preliminary version of this study appeared in RECOMB (Jain et al, 2022).

#### REFERENCES

- Abouelhoda M, Ohlebusch E. Chaining algorithms for multiple genome comparison. J Discrete Algorithms 2005; 3(2–4):321–341; doi: 10.1016/j.jda.2004.08.011
- Abouelhoda MI, Kurtz S, Ohlebusch E. Coconut: An efficient system for the comparison and analysis of genomes. BMC Bioinf 2008;9(1):476; doi: 10.1186/1471-2105-9-476
- Backurs A, Indyk P. Edit distance cannot be computed in strongly subquadratic time (unless SETH is false). In: Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC, Association for Computing Machinery, New York, NY, United States 2015; pp. 51–58.
- Bray N, Dubchak I, Pachter L. Avid: A global alignment program. Genome Res 2003;13(1):97–102; doi: 10.1101/gr.789803
- Chaisson MJ, Tesler G. Mapping single molecule sequencing reads using basic local alignment with successive refinement (blasr): Application and theory. BMC Bioinf 2012;13(328):1–17; doi: 10.1186/1471-2105-13-238
- de Berg M, Cheong O, van Kreveld MJ, et al. Computational geometry: Algorithms and applications, 3rd Edition. Springer-Verlag, Berlin, Heidelberg, 2008.
- Delcher AL, Kasif S, Fleischmann RD, et al. Alignment of whole genomes. Nucleic Acids Res 1999;27(11):2369–2376; doi: 10.1093/nar/27.11.2369
- Eppstein D, Galil Z, Giancarlo R, et al. Sparse dynamic programming i: Linear cost functions. JACM 1992a;39(3):519–545; doi: 10.1145/146637.146650
- Eppstein D, Galil Z, Giancarlo R, et al. Sparse dynamic programming ii: Convex and concave cost functions. JACM 1992b;39(3):546–567; doi: 10.1145/146637.146650.
- Hoppenworth G, Bentley JW, Gibney D, et al. The fine-grained complexity of median and center string problems under edit distance. In: 28th Annual European Symposium on Algorithms, ESA 2020, September 7–9, 2020, Vol. 173. Schloss Dagstuhl Leibniz-Zentrum für Informatik: Pisa, Italy, 2020; pp. 61:1–61:19.
- Jain C, Gibney D, Thankachan SV. Co-linear chaining with overlaps and gap costs. Research in Computational Biology In: 26th International Conference on Research in Computational Molecular Biology, RECOMB 2022, 22–25 May 2022, Springer, San Diego, CA, USA, pp. 246–262. bioRxiv 2021; doi: 10.1007/978-3-031-04749-7\\_15
- Jain C, Rhie A, Hansen NF, et al. Long-read mapping to repetitive reference sequences using winnowmap2. Nat Methods 2022;19(6):705–710; doi: 10.1038/s41592-022-01457-8
- Kalikar S, Jain C, Vasimuddin M, et al. Accelerating minimap2 for long-read sequencing applications on modern cpus. Nat Comput Sci 2022;2(2):78–83; doi: 10.1038/s43588-022-00201-8

- Kurtz S, et al. Versatile and open software for comparing large genomes. Stefan Kurtz. Adam M. Phillippy (eds)-Genome Biol 2004;5(2):R12; doi: 10.1186/gb-2004-5-2-r12
- Li H. Minimap2: Pairwise alignment for nucleotide sequences. Bioinformatics 2018;34(18):3094–3100; doi: 10.1093/bioinformatics/bty191
- Li H, Feng X, Chu C. The design and construction of reference pangenome graphs with minigraph. Genome Biol 2020;21(1):265; doi: 10.1186/s13059-020-02168-z
- Mäkinen V, Sahlin K. Chaining with overlaps revisited. In: 31st Annual Symposium on Combinatorial Pattern Matching, CPM 2020, June 17–19, 2020, Vol. 161. Schloss Dagstuhl Leibniz-Zentrum für Informatik: Copenhagen, Denmark, 2020; pp. 25:1–25:12.
- Mäkinen V, Tomescu AI, Kuosmanen A, et al. Sparse dynamic programming on dags with small width. ACM Trans Algorithms 2019;15(2):29:1–29:21; doi: 10.1145/3301312
- Marçais G, Delcher AL, Phillippy AM, et al. Mummer4: A fast and versatile genome alignment system. PLoS Comput Biol 2018;14(1):e1005944; doi: 10.1371/journal.pcbi.1005944
- Morgenstern B. A simple and space-efficient fragment-chaining algorithm for alignment of DNA and protein sequences. Appl Math Lett 2002;15(1):11–16; doi: 10.1016/S0893-9659(01)00085-4
- Myers G. A fast bit-vector algorithm for approximate string matching based on dynamic programming. JACM 1999;46(3):395–415; doi: 10.1145/316542.316550
- Myers G, Miller W. Chaining multiple-alignment fragments in sub-quadratic time. SODA '95: Proceedings of the Sixth Annual ACM-SIAM Symposium on Discrete Algorithms 1995;95:38–47; doi: 10.5555/313651
- Otto C, Hoffmann S, Gorodkin J, et al. Fast local fragment chaining using sum-of-pair gap costs. Algorithms Mol Biol 2011;6(1):4; doi: 10.1186/1748-7188-6-4
- Ren J, Chaisson MJP. Ira: A long read aligner for sequences and contigs. PLoS Comput Biol 2021;17(6):e1009078; doi: 10.1371/journal.pcbi.1009078
- Sahlin K, Mäkinen V. Accurate spliced alignment of long RNA sequencing reads. Bioinformatics 2021;37(24):4643–4651; doi: doi: 10.1093/bioinformatics/btab540
- Schleimer S, Wilkerson DS, Aiken A. Winnowing: Local algorithms for document fingerprinting. In: Proceedings of the 2003 ACM SIGMOD International Conference on Management of data. Association for Computing Machinery, New York, NY, USA, 2003; pp. 76–85.
- Sedlazeck FJ, Rescheneder P, Smolka M, et al. Accurate detection of complex structural variations using single-molecule sequencing. Nat Methods 2018;15(6):461–468; doi: 10.1038/s41592-018-0001-7
- Shibuya T, Kurochkin I. Match chaining algorithms for cDNA mapping. In: Algorithms in Bioinformatics, Third International Workshop, WABI 2003, Budapest, Hungary, September 15–20, 2003, Proceedings. Springer, Berlin, Heidelberg, Germany, 2003; pp. 462–475.
- Šošić M, Šikić M. Edlib: A C/C++ library for fast, exact sequence alignment using edit distance. Bioinformatics 2017;33(9):1394–1395; doi: 10.1093/bioinformatics/btw753.
- Ukkonen E. Algorithms for approximate string matching. Inf Control 1985;64(1-3):100-118; doi: 10.1016/S0019-9958(85)80046-2
- Uricaru R, Mancheron A, Rivals E, et al. Novel definition and algorithm for chaining fragments with proportional overlaps. J Comput Biol 2011;18(9):1141–1154; doi: 10.1089/cmb.2011.0126
- Vyverman M, De Baets B, et al. essamem: Finding maximal exact matches using enhanced sparse suffix arrays. Bioinformatics 2013;29(6):802–804; doi: 10.1093/bioinformatics/btt042
- Wilbur WJ, Lipman DJ. Rapid similarity searches of nucleic acid and protein data banks. Proc Natl Acad Sci U S A 1983;80(3):726–730; doi: 10.1073/pnas.80.3.726

E-mail: daniel.j.gibney@gmail.com