

Albert Einstein

Eine neue Bestimmung der Moleküldimensionen

Diss. ETH No. ?

PESHO IVANOV

OPTIMAL SEQUENCE ALIGNMENT USING A*

DISS. ETH NO. ?

OPTIMAL SEQUENCE ALIGNMENT USING A*

A dissertation submitted to attain the degree of

DOCTOR OF SCIENCES of ETH ZURICH

(Dr. sc. ETH Zurich)

presented by

PESHO IVANOV

Dipl., Eidgenössisches Polytechnikum

born on 29 May 1989

citizen of Bulgaria

accepted on the recommendation of

Prof. Dr. Martin Vechev, examiner

Prof. Dr. Gunnar Rätsch, co-examiner

Prof. Dr. Veli Mäkinen, co-examiner

Prof. Dr. Paul Medvedev, co-examiner

2022

With gratitude to my high-school physics teacher Svilen Rusev. Thank you
for opening the world of science to me in an intuitive and visual way.

ABSTRACT

Sequence alignment is the process of detecting similarities between biological sequences, such as DNA, RNA and proteins. For the last half a century, sequence alignment has been of central importance for molecular biology. Applications include evolutionary biology, genome assembly, read mapping, variation detection and computational methods are crucial for the correctness of the analyses of the vast amounts of biological data. Can we use the the A* shortest path algorithm to find optimal alignments fast?

This thesis explores two variations of the alignment problem: *reads mapping* and *pairwise alignment*, also known as semi-global and global alignment. Biological sequences do not generally align perfectly due to biological differences and technical errors. Given two sequences, the desired alignment is a position-to-position correspondence between two sequences which minimizes the edit costs (substitutions, insertions or deletions). This task is closely related to calculating *edit distance*. Practical alignment algorithms are desired to ① find accurate alignments, ② apply to a wide range of data, and ③ use little time and memory.

Existing aligning algorithms are either optimal but quadratic or fast but approximate. Even though the resulting alignment is linear. This gap has motivated the development of various faster but approximate algorithms. Moreover, theoretical results suggest that it is that in general, alignment is not solvable in strongly subquadratic time.

We consider a principled alignment formulation based on shortest paths and demonstrate that the A* shortest path algorithm can be used to outperform current methods. Unlike existing methods, A* enables an *informed search* based on information from the unaligned sequence suffixes, thus radically improving the empirical runtime scaling (up to linear) in the average case while providing optimality guarantees. On real data, this approach reaches orders of magnitude of speedup compared to existing approaches.

An optimal alignment can naturally be represented as a shortest path in an alignment graph (equivalent to the DP table). In order to find such a shortest path with minimal exploration, we instantiate the A* algorithm with a novel problem-specific heuristic function based on the unaligned parts of the sequences. This additional information is a problem-specific heuristic function and it heavily determines the efficiency of the search. For any explored state by A*, this heuristic function should compute a lower

bound on the remaining path length, or more specifically, the minimal cost of edit operations needed to align the remaining sequences.

Seed heuristic. In practice, while achieving polynomial speed ups on real data. It ① provides optimality guarantees according to edit distance, ② scales to long and noisy sequences, and ③ scales subquadratically with sequence length. To scale to large reference sequences, we extend the graph with a trie index. To scale to long queries, we introduce design an admissible *seed heuristic*, which is provably-optimal also efficient to compute. To scale to high error rates, we design

Many tools do not have a well-stated problem they optimize.

Probabilistic approach. Focus on the metric. Extend to MSA, local, affine. Prototype implementations No asymptotics.

ZUSAMMENFASSUNG

Deutsche Zusammenfassung hier.

ACKNOWLEDGEMENTS

I would like to thank ...

CONTENTS

1	INTRODUCTION	1
1.1	Introduction	4
1.1.1	Related work on alignment	4
1.1.2	Contributions	6
1.2	Preliminaries	7
1.3	Methods	9
1.3.1	Overview	10
1.3.2	Seed heuristic and chaining seed heuristic	11
1.3.3	Seed heuristic as a maximization problem	12
1.3.4	Chaining seed heuristic as a maximization problem	13
1.3.5	Match pruning	15
1.4	Algorithm	16
1.4.1	A* algorithm	16
1.4.2	Computing the seed heuristic	17
1.4.3	Computing the chaining seed heuristic	18
1.4.4	Implementation	19
1.5	Results	20
1.5.1	Setup	21
1.5.2	Comparison on synthetic data	23
1.5.3	Comparison on real data	24
1.6	Discussion	25
1.7	Conclusion	27
A	Appendix	28
A.1	The seed heuristic and chaining seed heuristic are admissible	28
A.2	Partially admissible heuristic	29
A.3	Match pruning preserves finding a shortest path	30
A.4	Aligner versions and parameters	31
A.5	Effect of pruning, chaining, and inexact matches	32
A.6	Expanded states and equivalent band	33
A.7	Performance variation on human genome	33
A.8	Complex alignments	34
A.9	Notation	35
A	SUMMARY	37
A	APPENDIX	39

NOTATION

FREQUENTLY USED SYMBOLS

E energy

m rest mass

p impulse

PHYSICAL CONSTANTS

c speed of light in vacuum, $c = 299\,792\,458\,\text{m s}^{-1}$

(CODATA 2014 [**codata**])

Maxwell [1]

INTRODUCTION

It is better to be wrong than to be vague.

— Freeman Dyson

The number of possible alignments grow exponentially with length. The usual underlying question to finding “correct” alignments. Regarding the precision of alignment, one is usually interested in base-to-base (aka letter-to-letter) correspondence between the sequences, even though for some applications a less detailed solution is sufficient: only the similarity between sequences or the location where a read maps to a reference. Exact alignment is only useful for very short sequences (often kmers), and for all other cases the optimized metric may be hamming distance, edit distance (unit costs), Levenshtein distance, affine costs, convex and concave costs, general costs and others.

Depending on the the number of aligned sequences, there is pairwise alignment and multiple sequence alignment (MSA). Depending on the parts of the sequences that are aligned to each other, we differentiate global, local and various semi-global alignemnts. There are generalizations to sequence-to-sequence alignment, including aligning to nonlinear structures, such as directed acyclic graphs, DAGs, general graphs and others. These structures are nowadays becoming more common as a compressed form of representing a set of references to which a sequence can be aligned. Often, one best alignment is sufficient but finding several best (top-K) alignments. In the context of read mapping, a set of reads is aligned to the same reference sequence so an indexing procedure is often useful for the performance.

We specifically consider the mapping of a set of reads to a general graph, and the global pairwise alignment.

Existing optimal algorithms are based on dynamic programming (DP) and run in quadratic time (assuming that the number of errors is proportional to the length)

we employ the A* algorithm which is an *informed search* algorithm. TODO: a case for the informed algorithms

MOTIVATION Sequence alignment has been a core problem in computational biology for the last half-century. It is an open problem whether exact pairwise alignment is possible in linear time for related sequences [2].

METHODS We solve exact global pairwise alignment with respect to edit distance by using the A* shortest path algorithm on the edit graph. In order to efficiently align long sequences with high error rate, we extend the *seed heuristic* for A* [3] with *match chaining*, *inexact matches*, and the novel *match pruning* optimization. We prove the correctness of our algorithm and provide an efficient implementation in A*PA.

RESULTS We evaluate A*PA on synthetic data (random sequences of length n with uniform mutations with error rate e) and on real long ONT reads of human data. On the synthetic data with $e=5\%$ and $n \leq 10^7$ bp, A*PA exhibits a near-linear empirical runtime scaling of $n^{1.08}$ and achieves $>250\times$ speedup compared to the leading exact aligners EDLIB and BtWFA. Even for a high error rate of $e=15\%$, the empirical scaling is $n^{1.28}$ for $n \leq 10^7$ bp. On two real datasets, A*PA is the fastest aligner for 58% of the alignments when the reads contain only sequencing errors, and for 17% of the alignments when the reads also include biological variation.

AVAILABILITY github.com/RagnarGrootKoerkamp/astar-pairwise-aligner

CONTACT ragnar.grootkoerkamp@inf.ethz.ch, pesho@inf.ethz.ch

1.1 INTRODUCTION

The problem of aligning one biological sequence to another has been formulated over half a century ago [4] and is known as *global pairwise alignment* [5]. Pairwise alignment has numerous applications in computational biology, such as genome assembly, read mapping, variant detection, multiple sequence alignment, and differential expression [6]. Despite the centrality and age of pairwise alignment, “a major open problem is to implement an algorithm with linear-like empirical scaling on inputs where the edit distance is linear in n ” [2].

Alignment accuracy affects the subsequent analyses, so a common goal is to find a shortest sequence of edit operations (insertions, deletions, and substitutions of single letters) that transforms one sequence into the other. Finding such a sequence of operations is at least as hard as computing the *edit distance*, which has recently been proven to not be computable in strongly subquadratic time, unless SETH is false [7]. Given that the number of sequencing errors is proportional to the length, existing exact aligners are limited by quadratic scaling not only in the worst case but also in practice. This is a computational bottleneck given the growing amounts of biological data and the increasing sequence lengths [8].

1.1.1 *Related work on alignment*

We outline algorithms for exact pairwise alignment and their fastest implementations for biological sequences. Refer to Kucherov [8] for approximate, probabilistic, and non-edit distance algorithms and aligners.

DYNAMIC PROGRAMMING The standard approach to sequence alignment is by successively aligning prefixes of the first sequence to prefixes of the second. Vintsyuk [9] was the first to introduce this $O(nm)$ dynamic programming (DP) approach for comparing a pair speech signals with n and m elements. Independently, it was applied to compute edit distance for biological sequences [4, 10–12]. This well-known algorithm is implemented in modern aligners like SEQAN [13] and PARASAIL [14]. Improving the performance of these quadratic algorithms has been a central goal in later works.

BANDING AND BIT-PARALLELIZATION When similar sequences are being aligned, the whole DP table may not need to be computed. One such

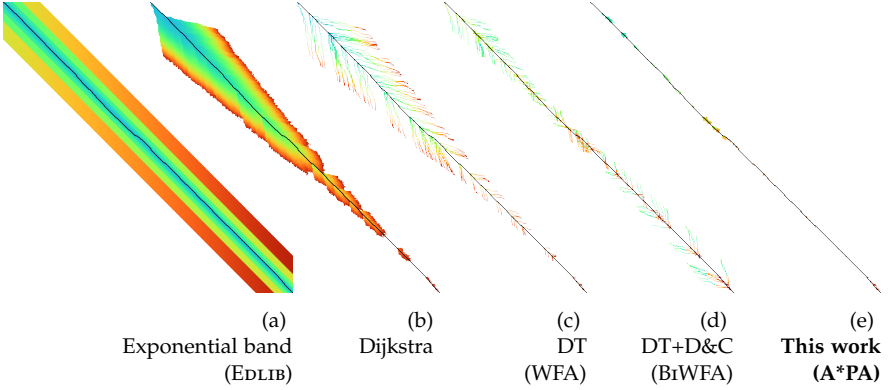


FIGURE 1.1: Demonstration of the computed states by various optimal alignment algorithms and corresponding aligners that implement them on synthetic data (length $n=500$ bp, error rate $e=20\%$). Blue-to-red coloring indicates the order of computation. a Exponential banding algorithm (EDLIB), b Dijkstra, c Diagonal transition/DT (WFA), d Diagonal transition with divide-and-conquer/D&C (BrWFA), e A^* with chaining seed heuristic and match pruning (seed length $k=5$ and exact matches).

output-sensitive algorithm is the *banded* algorithm of Ukkonen [15] (fig. 1.1a) which considers only states near the diagonal within an exponentially increasing *band*, and runs in $O(ns)$ time, where s is the edit distance between the sequences. This algorithm, combined with the *bit-parallel optimization* by Myers [16] is implemented by the EDLIB aligner [17] that runs in $O(ns/w)$ runtime, where w is the machine word size (nowadays 32 or 64).

DIAGONAL TRANSITION AND WFA The $O(ns)$ runtime complexity can be improved using the algorithm independently discovered by Ukkonen [15] and Myers [18] that is known as *diagonal transition* [5] (fig. 1.1c). It has an $O(ns)$ runtime in the worst-case but only takes expected $O(n + s^2)$ time under the assumption that the input sequences are random [18]. This algorithm has been extended to linear and affine costs in the *wavefront alignment* (WFA) algorithm [19] in a way similar to Gotoh [20], and has been improved to only require linear memory in BrWFA [21] by combining it with the *divide and conquer* approach of Hirschberg [22], similar to Myers [18] algorithm for unit edit costs (fig. 1.1d). Note that when each sequence

letter has an error with a constant probability, the total number of errors s is proportional to n , so that even these algorithms have a quadratic runtime.

SHORTEST PATHS AND A^* A pairwise alignment that minimizes edit distance corresponds to a shortest path in the *edit graph* [9, 15]. Assuming non-negative edit costs, a shortest path can be found using Dijkstra’s algorithm [15] (fig. 1.1b) or A^* [23]. A^* is an informed search algorithm which uses a task-specific heuristic function to direct its search. Depending on the heuristic function, a shortest path may be found significantly faster than by an uninformed search such as Dijkstra’s algorithm. In the context of semi-global sequence-to-graph alignment, A^* has been used to empirically scale sublinearly with the reference size for short reads [24].

SEEDS AND MATCHES Seed-and-extend is a commonly used paradigm for solving semi-global alignment approximately [8]. Seeds are also used to define and compute LCSk [25], a generalization of longest common subsequence (LCS). In contrast, the *seed heuristic* by Ivanov, Bichsel & Vechev [3] speeds up finding an optimal alignment by using seed matches to speed up the A^* search. The seed heuristic enables empirical near-linear scaling to long HiFi reads (up to 30 *kbp* with 0.3% errors) [3]. A limitation of the existing seed heuristic is the low tolerance to increasing error rates due to using only long exact matches without accounting for their order.

1.1.2 Contributions

Our algorithm exactly solves global pairwise alignment for edit distance costs, also known as *Levenshtein distance* [26]. It uses the A^* algorithm to find a shortest path in the edit graph.

SEED HEURISTIC In order to handle higher error rates, we extend the *seed heuristic* [3] to *inexact matches*, allowing up to 1 error in each match. To handle cases with a large number of seed matches we introduce *match chaining*, constraining the order in which seed matches can be linked [27, 28]. We prove that our *chaining seed heuristic* with inexact matches is admissible, which guarantees that A^* finds a shortest path.

MATCH PRUNING In order to reduce the number of states expanded by A^* we apply the *multiple-path pruning* observation of Poole & Mackworth [29]: once a shortest path to a vertex has been found, no other paths to this

vertex can improve the global shortest path. We prove that when a state at the start of a match is expanded, a shortest path to this state has been found. Since no other path to this state can be shorter, we show that we can *prune* (remove) the match, thus improving the seed heuristic. This incremental heuristic search has some similarities to Real-time Adaptive A* [30].

IMPLEMENTATION We efficiently implement our algorithm in the A*PA aligner. In particular, we use *contours* [31–33] to efficiently compute the chaining seed heuristic, and update them when pruning matches.

SCALING AND PERFORMANCE We compare the scaling and performance of our algorithm to other exact aligners on synthetic data, consisting of random genetic sequences with up to 15% uniform errors and up to 10^7 bases. We demonstrate that inexact matches and match chaining enable scaling to higher error rates, while match pruning enables near-linear scaling with length by reducing the number of expanded states to not much more than the best path (fig. 1.1e). Our empirical results show that for $e=5\%$ and $n=10^7$ bp, A*PA outperforms the leading aligners EDLIB [17] and BiWFA [21] by more than 250 times.

We demonstrate a limited applicability of our algorithm to long Oxford Nanopore (ONT) reads from human samples. A*PA is the fastest exact aligner on 58% of the alignments on a dataset with only sequencing errors, and on 17% of the alignments on a dataset with biological variation.

1.2 PRELIMINARIES

This section provides background definitions which are used throughout the paper.

SEQUENCES The input sequences $A = \overline{a_0 a_1 \dots a_i \dots a_{n-1}}$ and $B = \overline{b_0 b_1 \dots b_j \dots b_{m-1}}$ are over an alphabet Σ with 4 letters. We refer to substrings $\overline{a_i \dots a_{i'-1}}$ as $A_{i \dots i'}$, to prefixes $\overline{a_0 \dots a_{i-1}}$ as $A_{<i}$, and to suffixes $\overline{a_i \dots a_{n-1}}$ as $A_{\geq i}$. The *edit distance* $\text{ed}(A, B)$ is the minimum number of insertions, deletions, and substitutions of single letters needed to convert A into B .

EDIT GRAPH Let *state* $\langle i, j \rangle$ denote the subtask of aligning the prefix $A_{<i}$ to the prefix $B_{<j}$. The *edit graph* (also called *alignment graph*) $G(V, E)$ is a weighted directed graph with vertices $V = \{\langle i, j \rangle \mid 0 \leq i \leq n, 0 \leq j \leq m\}$ corresponding to all states, and edges connecting tasks to subtasks: edge

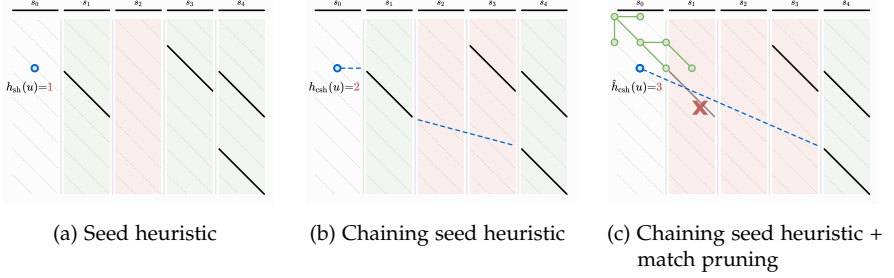


FIGURE 1.2: A demonstration of the seed heuristic, chaining seed heuristic, and match pruning. Sequence A is split into 5 seeds drawn as horizontal black segments (—) on top. Their exact matches in B are drawn as diagonal black segments (↖). The heuristic is evaluated at the blue state (u, \circ) , based on the 4 remaining seeds. The dashed blue paths show maximal length chains of seed matches (green columns, ■). The seeds that are not matched (red columns, ■) count towards the heuristic. a The seed heuristic $h_{sh}(u) = 1$ is the number of remaining seeds that do not have matches (s_2). b The chaining seed heuristic $h_{csh}(u) = 2$ is the number of not matched remaining seeds (s_2 and s_3) for a path going only down and to the right containing a maximal number of matches. c Once the start of a match is expanded (shown as green circles, \circ), the match is *pruned* (marked with red cross, \times), and future computations of the heuristic ignore it. This reduces the maximum chain of matches starting at u by 1 (s_1) so that $\hat{h}_{csh}(u)$ increases by 1.

$\langle i, j \rangle \rightarrow \langle i+1, j+1 \rangle$ has cost 0 if $a_i = b_j$ (match) and 1 otherwise (substitution), and edges $\langle i, j \rangle \rightarrow \langle i+1, j \rangle$ (deletion) and $\langle i, j \rangle \rightarrow \langle i, j+1 \rangle$ (insertion) have cost 1. We denote the root state $\langle 0, 0 \rangle$ by v_s and the target state $\langle n, m \rangle$ by v_t . For brevity we write $f(\langle i, j \rangle)$ as $f\langle i, j \rangle$. The edit graph is a natural representation of the alignment problem that provides a base for all alignment algorithms.

PATHS, ALIGNMENTS, SEEDS AND MATCHES Any path from $\langle i, j \rangle$ to $\langle i', j' \rangle$ in the edit graph G represents a *pairwise alignment* (or just *alignment*) of the substrings $A_{i\dots i'}$ and $B_{j\dots j'}$. We denote with $d(u, v)$ the distance between states u and v . A shortest path π^* corresponds to an optimal alignment, thus $\text{cost}(\pi^*) = d(v_s, v_t) = \text{ed}(A, B)$. For a state u we write $g^*(u) := d(v_s, u)$ and $h^*(u) := d(u, v_t)$ for the distance from the start to u and from u to the

target v_t , respectively. We define a *seed* as a substring $A_{i\dots i'}$ of A , and an *exact match* of a seed as an alignment of the seed of cost 0.

DIJKSTRA AND A^* Dijkstra's algorithm [34] finds a shortest path from v_s to v_t by *expanding* vertices in order of increasing distance $g^*(u)$ from the start. The A^* algorithm [35, 36], instead, directs the search towards a target by expanding vertices in order of increasing $f(u) := g(u) + h(u)$, where $h(u)$ is a heuristic function that estimates the distance $h^*(u)$ to the end and $g(u)$ is the shortest length of a path from v_s to u found so far. A heuristic is *admissible* if it is a lower bound on the remaining distance, $h(u) \leq h^*(u)$, which guarantees that A^* has found a shortest path as soon as it expands v_t . Heuristic h_1 *dominates* another heuristic h_2 when $h_1(u) \geq h_2(u)$ for all vertices u . A dominant heuristic will usually, but not always [37], expand less vertices. Note that Dijkstra's algorithm is equivalent to A^* using a heuristic that is always 0, and that both algorithms require non-negative edge costs. Our variant of the A^* algorithm is provided in section 1.4.1.

CHAINS A state $u = \langle i, j \rangle \in V$ *precedes* a state $v = \langle i', j' \rangle \in V$, denoted $u \preceq v$, when $i \leq i'$ and $j \leq j'$. Similarly, a match m precedes a match m' , denoted $m \preceq m'$, when the end of m precedes the start of m' . This makes the set of matches a partially ordered set. A state u precedes a match m (denoted $u \preceq m$) when it precedes the start of the match. A *chain* of matches is a (possibly empty) sequence of matches $m_1 \preceq \dots \preceq m_l$.

CONTOURS To efficiently calculate maximal chains of matches, *contours* are used. Given a set of matches \mathcal{M} , $S(u)$ is the number of matches in the longest chain $u \preceq m_0 \preceq \dots$, starting at u . The function $S\langle i, j \rangle$ is non-increasing in both i and j . *Contours* are the boundaries between regions of states with $S(u) = \ell$ and $S(u) < \ell$ (see fig. 1.3). Note that contour ℓ is completely determined by the set of matches $m \in \mathcal{M}$ for which $S(\text{start}(m)) = \ell$ [31]. Hunt & Szymanski [32] give an algorithm to efficiently compute S when \mathcal{M} is the set of single-letter matches between A and B , and Deorowicz & Grabowski [38] give an algorithm when \mathcal{M} is the set of k -mer exact matches.

1.3 METHODS

In this section we define the *seed heuristic* with *inexact matches* and *chaining*, and *match pruning* (see fig. 1.2). We motivate each of these extensions

intuitively, define them formally, and prove that A^* is guaranteed to find a shortest path. We end with a reformulation of our heuristic definitions that allows for more efficient computation. appendix A.9 contains a summary of the notation we use.

1.3.1 Overview

In order to find a minimal cost alignment of A and B , we use the A^* algorithm to find a shortest path from the starting state $v_s = \langle 0, 0 \rangle$ to the target state $v_t = \langle n, m \rangle$ in the edit graph. We present two heuristic functions, the seed heuristic and the chaining seed heuristic, and prove that they are admissible, so that A^* always finds a shortest path.

To define the seed heuristic h_{sh} , we split A into short, non-overlapping substrings (*seeds*) of fixed length k . Since the whole of sequence A has to be aligned, each of the seeds also has to be aligned somewhere. If a seed cannot be *exactly* matched in B without mistakes, then at least one edit has to be made to align it. We first compute all positions in B where each seed from A matches exactly. Then, a lower bound on the edit distance between the remaining suffixes $A_{\geq i}$ and $B_{\geq j}$ is given by the number of seeds entirely contained in $A_{\geq i}$ that do not match anywhere in B . An example is shown in fig. 1.2a.

We improve the seed heuristic by enforcing that the seed matches occur in the same order as their seeds occur in A , i.e., they form a *chain*. Now, the number of upcoming errors is at least the minimal number of remaining seeds that cannot be aligned on a single path to the target. The chaining seed heuristic h_{csh} equals the number of remaining seeds minus the maximum length of a chain of matches, see fig. 1.2b.

To scale to larger error rates, we generalize both heuristics to use *inexact matches*. For each seed from A , our algorithm finds all its inexact matches in B with cost at most 1. Then, not using any match of a seed will require at least $r=2$ edits for the seed to be eventually aligned.

Finally, we use *match pruning*: once we find a shortest path to a state v , no shorter path to that state is possible. Since we are only looking for a single shortest path, we may ignore any other path going to v . In particular, when we are evaluating the heuristic h in some state u preceding v , and v is at the start of a match, all paths containing the match are excluded, and thus we may simply ignore (*prune*) this match for future computations of the heuristic. This increases the value of the heuristic, as shown in fig. 1.2c, and

leads to a significant reduction of the number of states being expanded by the A^* .

1.3.2 Seed heuristic and chaining seed heuristic

We start with the formal definitions of seeds and matches that are the basis of our heuristic functions and algorithms.

SEEDS We split the sequence A into a set of consecutive non-overlapping substrings (*seeds*) $\mathcal{S} = \{s_0, s_1, s_2, \dots, s_{\lfloor n/k \rfloor - 1}\}$, such that each seed $s_l = A_{lk \dots lk+k}$ has length k . After aligning the first i letters of A , the information for the heuristic will come from the *remaining* seeds $\mathcal{S}_{\geq i} := \{s_l \in \mathcal{S} \mid lk \geq i\}$ contained in the suffix $A_{\geq i}$.

SEED ALIGNMENT An *alignment* of a seed $s = A_{i \dots i+k}$ is a path π_s in the edit graph G from a state $u = \langle i, j \rangle$ to a state $v = \langle i+k, j' \rangle$, where i is the start and $i+k$ is the end of s . We refer to the cost of the alignment π_s as $\text{cost}(\pi_s)$.

MATCHES In order to only consider good alignments of seeds, we fix a threshold cost r called the *seed potential*. We define a *match* m as an alignment of a seed with cost $\text{cost}(m) < r$. For each seed s , the set of all of its matches is \mathcal{M}_s . The inequality is strict so that $\mathcal{M}_s = \emptyset$ implies that aligning the seed will incur cost at least r . Let $\mathcal{M} = \bigcup_s \mathcal{M}_s$ denote the set of all matches. With $r=1$ we allow only *exact* matches, while for $r=2$ we allow both exact and *inexact* matches with one edit. In this paper we do not consider larger r .

Next, we define two heuristic functions:

Definition 1 (Seed heuristic). *Given a set of matches \mathcal{M} with costs less than r , the seed heuristic $h_{\text{sh}}^{\mathcal{M}}$ in a state $u = \langle i, j \rangle$ is the minimal total cost to independently align all remaining seeds: each seed is either matched or incurs cost r if it cannot be matched in \mathcal{M} :*

$$h_{\text{sh}}^{\mathcal{M}}(u) := \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \min_{m \in \mathcal{M}_s} \text{cost}(m) & \text{if } \mathcal{M}_s \neq \emptyset, \\ r & \text{otherwise.} \end{cases} \quad (1.1)$$

Definition 2 (Chaining seed heuristic). *Given a set of matches \mathcal{M} of costs less than r , the chaining seed heuristic $h_{\text{csh}}^{\mathcal{M}}$ in a state $u = \langle i, j \rangle$ is the minimal cost to jointly align each remaining seed on a single path in the edit graph:*

$$h_{\text{csh}}^{\mathcal{M}}(u) := \min_{\pi \in G(u \rightsquigarrow v_t)} \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \text{cost}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ r & \text{otherwise.} \end{cases} \quad (1.2)$$

where π runs over all paths from u to v_t , and π_s is the shortest part of the path π that is an alignment of the seed s .

We abbreviate $h_{\text{sh}} := h_{\text{sh}}^{\mathcal{M}}$ and $h_{\text{csh}} := h_{\text{csh}}^{\mathcal{M}}$ when \mathcal{M} is the set of all matches of cost strictly less than r . For such \mathcal{M} , both heuristics guarantee finding an optimal alignment.

Theorem 1. *The seed heuristic h_{sh} and the chaining seed heuristic h_{csh} are admissible.*

The proof follows directly from the definitions and can be found in appendix A.1. As a consequence of the proof, $h_{\text{csh}}(u)$ dominates $h_{\text{sh}}(u)$, that is, $h_{\text{csh}}(u) \geq h_{\text{sh}}(u)$ for all u . Hence $h_{\text{csh}}(u)$ usually expands fewer states, at the cost of being more complex to compute.

1.3.3 Seed heuristic as a maximization problem

In order to efficiently evaluate the seed heuristic, eq. (1.1), we rewrite it from a minimization of match costs to a maximization of match scores. We first define a few new concepts, and then rewrite the equation into form that will be simpler to compute.

score We define the *score of a match m* as $\text{score}(m) := r - \text{cost}(m)$. This is always positive since match costs are defined to be strictly less than r . The *score of a seed s* is the maximal score of a match of s :

$$\text{score}(s) := \begin{cases} \max_{m \in \mathcal{M}_s} \text{score}(m) & \text{if } \mathcal{M}_s \neq \emptyset, \\ 0 & \text{otherwise.} \end{cases} \quad (1.3)$$

The *score of a state* is the sum of the scores of the remaining seeds,

$$S_{\text{sh}} \langle i, j \rangle := \sum_{s \in \mathcal{S}_{\geq i}} \text{score}(s). \quad (1.4)$$

POTENTIAL The *potential* $P\langle i, j \rangle$ is the value of the heuristic when there are no matches, and is the maximal value the heuristic can take in a given state:

$$P\langle i, j \rangle := r \cdot |\mathcal{S}_{\geq i}|. \quad (1.5)$$

OBJECTIVE eq. (1.1) can now be rewritten in terms of potential and score,

$$\begin{aligned} h_{\text{sh}}^{\mathcal{M}}(u) &\stackrel{(1.1)}{=} \sum_{s \in \mathcal{S}_{\geq i}} \left(r - \begin{cases} \max_{m \in \mathcal{M}_s} r - \text{cost}(m) & \text{if } \mathcal{M}_s \neq \emptyset, \\ 0 & \text{otherwise} \end{cases} \right) \\ &= r \cdot |\mathcal{S}_{\geq i}| - \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \max_{m \in \mathcal{M}_s} \text{score}(m) & \text{if } \mathcal{M}_s \neq \emptyset, \\ 0 & \text{otherwise} \end{cases} \\ &\stackrel{(1.3)}{=} r \cdot |\mathcal{S}_{\geq i}| - \sum_{s \in \mathcal{S}_{\geq i}} \text{score}(s) \stackrel{(1.4),(1.5)}{=} P(u) - S_{\text{sh}}(u). \end{aligned} \quad (1.6)$$

The potential $P(u)$ of a state is simple to compute, and the score $S_{\text{sh}}(u)$ can be computed by using *layers*.

LAYER Let *layer* \mathcal{L}_{ℓ} be the set of states u with score $S_{\text{sh}}(u) \geq \ell$. Since S_{sh} is non-increasing in i and independent of j , and only changes value at the start of seeds, \mathcal{L}_{ℓ} is fully determined by largest i_{ℓ} such that $S_{\text{sh}}\langle i_{\ell}, \cdot \rangle \geq \ell$. The score $S_{\text{sh}}\langle i, j \rangle$ is then the largest ℓ such that $i_{\ell} \geq i$.

1.3.4 Chaining seed heuristic as a maximization problem

Similar to the seed heuristic, we rewrite eq. (1.2) into a maximization form that will be simpler to compute.

CHAIN SCORE The *score of a chain* is the sum of the scores of the matches in the chain. Let $S_{\text{chain}}(m)$ be the maximum score of a chain starting with match m . This satisfies the recursion

$$S_{\text{chain}}(m) := \text{score}(m) + \max_{m' \preceq m'} S_{\text{chain}}(m'). \quad (1.7)$$

where this and the following maxima are taken as 0 when they are over an empty set. For the chaining seed heuristic, the score at a state is

$$S_{\text{csh}}(u) := \max_{u \preceq m} S_{\text{chain}}(m). \quad (1.8)$$

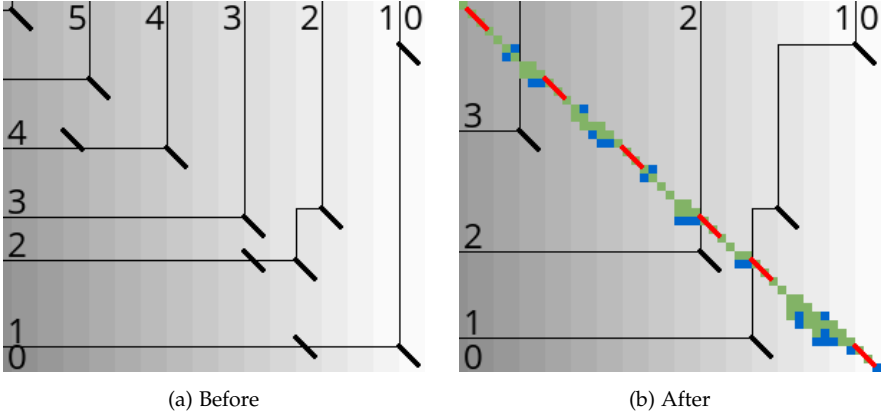


FIGURE 1.3: An example of the layers and contours used by the chaining seed heuristic before and after the A* execution with match pruning for $r=1$ and seed length $k=3$. Exact matches with $\text{score}(m)=1$ are shown as black diagonal segments (\blacktriangledown). Contours are shown as horizontal and vertical black lines and indicate the bottom-right boundaries of layers \mathcal{L}_ℓ consisting of the states above/left of the contour marked with ℓ . States u between two contours have the same maximal number of matches $S_{\text{csh}}(u)$ on a chain to the end. The value of the heuristic is shown in the background, from white ($h=0$) to darker grey. Expanded states are shown in green (\blacksquare), open states in blue (\blacksquare), and pruned matches in red (\blacklozenge). Note that pruning matches during the A* execution shifts the contours and changes the layers.

OBJECTIVE Similar to the seed heuristic, the chaining seed heuristic is computed via the following equality,

$$\begin{aligned}
 h_{\text{csh}}^{\mathcal{M}}(u) &\stackrel{(1.2)}{=} \min_{\pi \in G(u \rightsquigarrow v_t)} \sum_{s \in \mathcal{S}_{\geq i}} \left(r - \begin{cases} \text{score}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ 0 & \text{otherwise} \end{cases} \right) \\
 &= r \cdot |\mathcal{S}_{\geq i}| - \max_{\pi \in G(u \rightsquigarrow v_t)} \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \text{score}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ 0 & \text{otherwise} \end{cases} \\
 &= r \cdot |\mathcal{S}_{\geq i}| - \max_{u \preceq m_1 \preceq \dots \preceq m_l} \{ \text{score}(m_1) + \dots + \text{score}(m_l) \} \\
 &\stackrel{(1.5), (1.8)}{=} P(u) - S_{\text{csh}}(u).
 \end{aligned} \tag{1.9}$$

CONTOURS Like fore the seed heuristic, we define layer \mathcal{L}_ℓ as those states u with score $S_{\text{csh}}(u)$ at least ℓ . The ℓ th *contour* is the bottom-right boundary of \mathcal{L}_ℓ (see fig. 1.3). A state $u \in \mathcal{L}_\ell$ is *dominant* when there is no state $v \in \mathcal{L}_\ell$ with $v \neq u$ and $u \preceq v$. Both \mathcal{L}_ℓ and the corresponding contour are completely determined by the matches starting in the dominant states. The following lemma makes this precise:

Lemma 1. *For $\ell > 0$ the matches m with $\ell \leq S_{\text{chain}}(m) < \ell + r$ fully determine layer \mathcal{L}_ℓ :*

$$\mathcal{L}_\ell = \{u \mid \exists m \in \mathcal{M} : u \preceq m \text{ and } \ell \leq S_{\text{chain}}(m) < \ell + r\}. \quad (1.10)$$

Proof. It follows directly from eq. (1.8) that $\mathcal{L}_\ell = \{u \mid \exists m \in \mathcal{M} : u \preceq m \text{ and } S_{\text{chain}}(m) \geq \ell\}$. Suppose that m has score $S_{\text{chain}}(m) \geq \ell + r$. By definition, $\text{score}(m) \leq r$, so by eq. (1.7) there must be a match $m \preceq m'$ with $S_{\text{chain}}(m') = S_{\text{chain}}(m) - \text{score}(m) \geq (\ell + r) - r = \ell$. This implies that $u \preceq m \preceq m'$, which allows m , and hence any match with $S_{\text{chain}}(m) \geq \ell + r$, to be omitted from consideration. \square \square

We will use this lemma to efficiently find the largest ℓ such that layer \mathcal{L}_ℓ contains a state u , which gives $S_{\text{csh}}(u)$.

Note that the formulas for the chaining seed heuristic become equivalent to those for the seed heuristic when the partial order $\langle i, j \rangle \preceq \langle i', j' \rangle$ is redefined to mean $i \leq i'$, ignoring the j -coordinate.

1.3.5 Match pruning

In order to reduce the number of states expanded by the A^* algorithm, we apply the *multiple-path pruning* observation: once a shortest path to a vertex has been found, no other path to this vertex could possibly improve the global shortest path [29]. When A^* expands the start of a match we *prune* this match, so that the heuristic values of preceding states do no longer benefit from it, thus getting deprioritized by the A^* . We define *pruned* variants of the seed heuristic and the chaining seed heuristic that ignore pruned matches:

Definition 3 (Match pruning). *Let E be the set of expanded states during the A^* search, and let $\mathcal{M} \setminus E$ be the set of unpruned matches, i.e. those matches not starting in an expanded vertex. Then the pruning seed heuristic is $\hat{h}_{\text{sh}} := h_{\text{sh}}^{\mathcal{M} \setminus E}$ and the pruning chaining seed heuristic is $\hat{h}_{\text{csh}} := h_{\text{csh}}^{\mathcal{M} \setminus E}$.*

The hat (\hat{h}) denotes the implicit dependency on the progress of the A^* . Even though match pruning breaks the admissibility of our heuristics for some vertices, we prove that A^* is still guaranteed to find a shortest path:

Theorem 2. *A^* with match pruning heuristic \hat{h}_{sh} or \hat{h}_{csh} finds a shortest path.*

The proof can be found in appendix A.2. Since the heuristic may increase over time, our algorithm ensures that f is up-to-date when expanding a state by *reordering* nodes with outdated f values (see remark 1 in appendix A.2).

1.4 ALGORITHM

Our algorithm computes the shortest path in the edit graph using a variant of A^* that handles pruning (section 1.4.1). This depends on the efficient computation of the heuristics (sections 1.4.2 and 1.4.3). section 1.4.4 contains further implementation notes.

At a high level, we first initialize the heuristic by finding all seeds and matches and precomputing the potential $P[i]$ and layers \mathcal{L}_ℓ . Then we run the A^* search that evaluates the heuristic in many states, and updates the heuristic whenever a match is pruned.

1.4.1 A^* algorithm

We give our variant of the A^* algorithm [35] that adds steps A4 and A6a (marked as **bold**) to handle the pruning of matches. All computed values of g are stored in a map, and all states in the front are stored in a priority queue of tuples $(v, g(v), f(v))$ ordered by increasing f .

- A1. Set $g(v_s) = 0$ in the map and initialize the priority queue with $(v_s, g(v_s)=0, f(v_s))$.
- A2. Pop the tuple (u, g, f) with minimal f .
- A3. If $g > g(u)$, i.e. the value of $g(u)$ in the map has decreased since pushing the tuple, go to step 2. This is impossible for a consistent heuristic.
- A4.** If $f \neq f(u)$, *reorder u* : push $(u, g(u), f(u))$ and go to step 2. This can only happen when f changes value, i.e. after pruning.
- A5. If $u = v_t$, terminate and report a best path.

A6. Otherwise, *expand* u :

a. Prune matches starting at u .

b) For each successor v of u , *open* v when either $g(v)$ has not yet been computed or when $g(u) + d(u, v) < g(v)$. In this case, update the value of $g(v)$ in the map and insert $(v, g(v), f(v))$ in the priority queue. Go to step 2. This makes u the *parent* of v .

Note that the heuristic is evaluated in steps A4 and A6b for the computation of $f(u)$ and $f(v)$ respectively.

A vertex is called *closed* after it has been expanded for the first time, and called *open* when it is present in the priority queue.

1.4.2 Computing the seed heuristic

We compute the seed heuristic using an array LS that contains for each layer \mathcal{L}_ℓ the *layer start* $LS[\ell] := i_\ell = \max\{i : S_{\text{sh}}\langle i, \cdot \rangle \geq \ell\}$.

We give the algorithms to precompute LS , to evaluate the heuristic using it, and to update it after pruning.

PRECOMPUTATION

- R1. Compute the set of seeds \mathcal{S} by splitting A into consecutive kmers.
- R2. Build a hashmap containing all kmers of B .
- R3. For each seed, find all matches via a lookup in the hashmap. In case of inexact matches ($r = 2$), all sequences at distance 1 from each seed are looked up independently. The matches \mathcal{M} are stored in a hashmap keyed by the start of each match.
- R4. Initialize the array of potentials P by iterating over the seeds backwards.
- R5. Initialize the array of layer starts LS by setting $LS[0] = n$ and iterating over the seeds backwards. For each seed $s = A_{i\dots i'}$ that has matches, push $\text{score}(s)$ copies of i to the end of LS .

EVALUATING THE HEURISTIC IN $u = \langle i, j \rangle$

- E1. Look up the potential $P[i]$.
- E2. $S_{\text{sh}}(u) = \max\{\ell : LS[\ell] \geq i\}$ is found by a binary search over ℓ .
- E3. Return $h(u) = P[i] - S_{\text{sh}}(u)$.

PRUNING WHEN EXPANDING MATCH START $u = \langle i, j \rangle$

- P1. Compute $seedscore_{old} := \text{score}(s)$ and $\ell_{old} := S_{sh}(u)$.
- P2. Remove all matches from \mathcal{M} that start at u .
- P3. Compute $seedscore_{new} := \text{score}(s)$ and set $\ell_{new} := \ell_{old} - seedscore_{old} + seedscore_{new}$. If $\ell_{new} < \ell_{old}$, remove layers $LS[\ell_{new} + 1]$ to $LS[\ell_{old}]$ from LS and shift down larger elements $LS[\ell]$ with $\ell > \ell_{old}$ correspondingly.

1.4.3 Computing the chaining seed heuristic

The computation of the chaining seed heuristic is similar to that of the seed heuristic. It involves a slightly more complicated array LM of *layer matches* that associates to each layer \mathcal{L}_ℓ a list of matches with score ℓ : $LM[\ell] = \{m \in \mathcal{M} \mid S_{chain}(m) = \ell\}$. The score $S_{csh}(u)$ is then the largest ℓ such that $LM[\ell]$ contains a match m preceded by u .

Our algorithm for computing the chaining seed heuristic is similar to the one for the seed heuristic (section 1.4.2). Hence we highlight only the steps which differ (e.g. step 2' is used instead of 2 for the seed heuristic).

PRECOMPUTATION

- R5'. Initialize LM by setting $LM[0] = \emptyset$. Iterate over all matches in order of decreasing i of the match start, compute $\ell = S_{chain}(m)$ (see point 2' below), and add m to $LM[\ell]$.

EVALUATING THE HEURISTIC IN $u = \langle i, j \rangle$

- E2'. Because of ?? 1, the score $S_{csh}(u)$ can be computed using binary search: $S_{csh}(u) \geq \ell$ holds if and only if one of the layers $LM[\ell']$ with $\ell \leq \ell' < \ell + r$ contains a match m with $u \preceq m$. This is checked by simply iterating over all the matches in these layers.

PRUNING WHEN EXPANDING MATCH START $u = \langle i, j \rangle$

- P2'. Compute $\ell_u = S_{csh}(u)$, and remove all matches that start at u from layers $LM[\ell_u - r + 1]$ to $LM[\ell_u]$.
- P3'. Iterate over increasing ℓ starting at $\ell = \ell_u + 1$ and recompute $\ell' := S_{chain}(m) \leq \ell$ for all matches m in $LM[\ell]$. Move m from $LM[\ell]$ to layer $LM[\ell']$ whenever $\ell' \neq \ell$. Stop when either r consecutive layers

are unchanged, in which case no further changes of $S_{\text{chain}}(m)$ can happen because of eq. (1.7) and $\text{score}(m) \leq r$, or when all matches in r consecutive layers have shifted down by the same amount, say $\Delta := \ell - \ell'$. In the latter case, $S_{\text{chain}}(m)$ decreases by Δ for all matches with score at least ℓ . We remove the emptied layers $LM[\ell - \Delta + 1]$ to $LM[\ell]$ so that all higher layers shift down by Δ .

1.4.4 Implementation

This section covers some implementation details which are necessary for a good performance.

BUCKET QUEUE We use a hashmap to store all computed values of g in the A^* algorithm. Since the edit costs are bounded integers, we implement the priority queue using a *bucket queue* [39]. Unlike heaps, this data structure has amortized constant time push and pop operations since the difference between the value of consecutive pop operations is bounded.

GREEDY MATCHING OF LETTERS From a state $\langle i, j \rangle$ where $a_i = b_j$, it is sufficient to only consider the matching edge to $\langle i + 1, j + 1 \rangle$ [24, 40], and ignore the insertion and deletion edges to $\langle i, j + 1 \rangle$ and $\langle i + 1, j \rangle$. During alignment, we greedily match as many letters as possible within the current seed before inserting only the last opened state in the priority queue. We do not cross seed boundaries in order to not interfere with match pruning. We include greedily matched states in the reported number of expanded states.

PRIORITY QUEUE OFFSET Pruning the last match in a layer may cause an increase of the heuristic in all states preceding the start u of the match. This invalidates f values in the priority queue and causes reordering (step A4). We skip most of the update operations by keeping a global offset to the f -values in the priority queue, which is updated it when all states in the priority queue precede u .

SPLIT VECTOR FOR LAYERS Pruning a match may lead to the removal of one or more layers of the array of layer starts LS or the array of layer matches LM , after which all higher layers are shifted down. To support this efficiently, we use a *split vector*: a data structure that internally consists of two stacks, one containing the layers before the last deletion, and one containing the layers after the last deletion in reverse order. To delete a

layer, we move layers from the top of one of the stacks to the top of the other, until the layer to be deleted is at the top of one of the stacks, and then remove it. When layers are removed in decreasing order of ℓ , this takes linear total time.

BINARY SEARCH HINTS When we open v from its parent u in step A6b of the A* algorithm, the value of $h(v)$ is close to the value of $h(u)$. In particular, $S_{\text{sh}}(v)$ (resp. $S_{\text{csh}}(v)$) is close to and at most the score in u . Instead of a binary search (step E2), we do a linear search starting at the value of $\ell = S_{\text{sh}}(u)$ (resp. $\ell = S_{\text{csh}}(u)$).

We also speed up the binary search in the recomputation of $f(u)$ in the reordering check (step A4) by storing the last computed value of $\text{hint}(u) := S_{\text{sh}}(v_s) - S_{\text{sh}}(u)$ (resp. $S_{\text{csh}}(v_s) - S_{\text{csh}}(u)$) in the hashmap. When evaluating $S_{\text{sh}}(u)$, we start the linear search at $S_{\text{sh}}(v_s) - \text{hint}(u)$ with a fallback to binary search in case the linear search needs more than 5 iterations. Since $\text{hint}(u)$ remains constant when matches starting after u are pruned, it is a good starting point for the search when only matches near the tip of the A* search are pruned.

CODE CORRECTNESS Our implementation A*PA is written in Rust, contains many assertions testing e.g. the correctness of our A* and layers data structure implementation, and is tested for correctness and performance on randomly-generated sequences. Correctness is checked against simpler algorithms (Needleman-Wunsch) and other aligners (EDLIB, BiWFA).

1.5 RESULTS

We implemented the algorithms from section 1.2 in the aligner A*PA and refer to the two of our algorithms as SH for seed heuristic, and CSH for chaining seed heuristic. Here we empirically compare the runtime and memory usage to the exact optimal aligners EDLIB and BiWFA on both synthetic and real data. We demonstrate the benefit of match pruning on the runtime scaling with sequence length, and the benefits from match chaining and inexact matches on scaling to high error rates. All code, evaluation scripts and data are available in the A*PA repository.

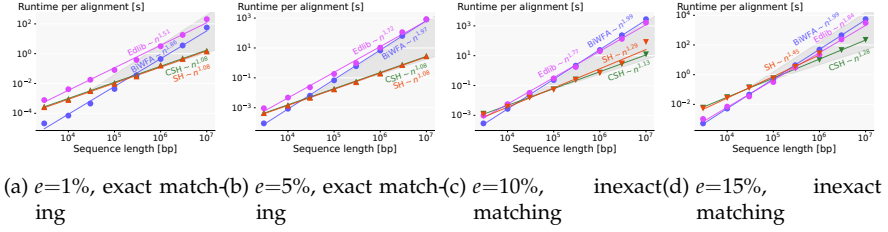


FIGURE 1.4: Log-log plots of the runtime for aligning **synthetic sequences** of increasing length for A*PA seed heuristic (SH), A*PA chaining seed heuristic (CSH), EDLIB (●) and BiWFA (●). The slopes of the bottom (top) of the dark-grey cones correspond to linear (quadratic) growth. For $e \leq 5\%$, SH (▼) and CSH (▼) use $k=15$, $r=1$. For $e \geq 10\%$, SH (▲) and CSH (▲) use $k=15$, $r=2$. The missing data points for SH at $e=15\%$ are due to exceeding the memory limit (30 GB). Each runtime is the average over $\lfloor 10^7/n \rfloor$ alignments.

1.5.1 Setup

SYNTHETIC DATA We measure the performance of aligners on a set of randomly-generated sequences. Each set is parametrized by the number of sequence pairs, the sequence length n , and the error rate e . The first sequence in a pair is generated by concatenating n i.i.d. letters from Σ . The second sequence is generated by sequentially applying $\lfloor e \cdot n \rfloor$ edit operations (insertions, deletions, and substitutions with equal probability) to the first sequence. Note that errors can cancel each other, so the final distance between the sequences is usually less than $\lfloor e \cdot n \rfloor$. In order to minimize the measurement errors, each test consists of a number of sequence pairs.

HUMAN DATA We consider two datasets¹ of Oxford Nanopore Technologies (ONT) reads that are aligned to regions of the telomere-to-telomere assembly of a human genome CHM13 (v1.1) [41]. Statistics are shown in table 1.1.

- *CHM13: Human reads without biological variation.* We randomly sampled 50 reads of length at least 500 *kbp* from the first 12 GB of the ultra-Long ONT reads that were used to assemble CHM13. We removed soft clipped

¹ <https://github.com/RagnarGrootKoerkamp/astar-pairwise-aligner/releases/tag/datasets>

Dataset	Biological mutations	Length [kbp]			Error rate [%]		
		min	mean	max	min	mean	max
CHM13	No	500	590	840	2.7	6.3	18.0
NA12878	Yes	502	624	1053	4.4	7.4	19.8

TABLE 1.1: Statistics on the **real data**: ONT reads from human samples.

regions and paired each read to its corresponding reference region in CHM13².

- *NA12878: Human reads with biological variation.* We consider the long ONT MinION reads from a different reference sample NA12878 [42] which was used to evaluate BiWFA [21]. This dataset contains 48 reads longer than 500 *kbp* that were mapped to CHM13.

COMPARED ALGORITHMS AND ALIGNERS We compare the seed heuristic and chaining seed heuristic, implemented in A*PA, to the state-of-the-art exact pairwise aligners BiWFA and EDLIB. In order to study the performance of the A* heuristic functions and the pruning optimization, we also compare to Dijkstra’s algorithm (which is equivalent to A* with a zero heuristic) and to a no-pruning variant of A*, all implemented in A*PA. The exact aligners SEQAN and PARASAIL are not included in this evaluation since they have been outperformed by BiWFA and EDLIB [19]. We execute all aligners with unit edit costs and compute not only the edit distance but also an optimal alignment. See appendix A.4 for aligner versions and parameters.

EXECUTION All evaluations are executed on Arch Linux on a single thread on an Intel Core i7-10750H CPU @ 2.6GHz processor with 64 GB of memory and 6 cores, with hyper-threading and performance mode disabled. We fix the CPU frequency to 2.6GHz and limit the available memory per execution to 30 GB using `ulimit -v 30000000`. To speed up the evaluations, we run 3 jobs in parallel, pinned to cores 0, 2, and 4.

MEASUREMENTS The runtime (wall-clock time) and memory usage (resident set size) of each run are measured using `ttime`. The runtime in all plots

² <https://github.com/RagnarGrootKoerkamp/bam2seq>

Aligner	Algorithm	Runtime [s]				
		$e =$	1%	5%	10%	15%
●	EDLIB	Banding, exponential search, bit-parallel	206	839	1653	3073
●	BIWFA	Diagonal transition, divide & conquer	56.9	806	2799	5492
▲▼	A*PA	A*, match-pruning, seed heuristic	1.41	2.79	82.5	ML
▲▼	A*PA	A*, match-pruning, chaining seed heuristic	1.60	2.98	12.3	220

TABLE 1.2: Runtime and memory comparison on **synthetic sequences** of length $n=10^7$ bp for various error rates. ML stands for exceeding the memory limit of 30 GB. Note that we run our A* heuristics with exact matches (▲▲) when $e \leq 5\%$, and with inexact matches (▼▼) when $e \geq 10\%$.

and tables refers to the average alignment time per pair of sequences. To reduce startup overhead, we average faster alignments over more pairs of sequences, as specified in the figure captions. Memory usage is measured as the maximum used memory during the processing of the whole input file. To estimate how algorithms scale with sequence length, we calculate best fit polynomials via a linear fit in the log-log domain using the least squares method.

PARAMETERS FOR THE A* HEURISTICS Longer sequences contain more potential locations for off-path seed matches (i.e. lying outside of the resulting optimal alignment). Each match takes time to be located and stored, while also potentially worsening the heuristic. To keep the number of matches low, longer seeds are to be preferred. On the other hand, to handle higher error rates, a higher number of shorter reads is preferable. For simplicity, we fix $k=15$ throughout our evaluations as a reasonable trade-off between scaling to large n and scaling to high e . We use exact matches ($r=1$) for low error rates ($e \leq 5\%$), and inexact matches ($r=2$) for high error rates ($e \geq 10\%$). For human datasets we use $r=2$ since they include sequences with error rates higher than 10%.

1.5.2 Comparison on synthetic data

Here we compare the runtime scaling and performance of the exact optimal aligners on synthetic data. appendix A.5 compares the effects of pruning

and inexact matches, whereas appendix A.6 compares SH and CSH in terms of the number of expanded states.

SCALING WITH SEQUENCE LENGTH First, we compare the aligners by runtime scaling with sequence length n for various error rates e (fig. 1.4). As expected from the theoretical analysis, EDLIB and BiWFA scale approximately quadratically regardless of the error rate. Unlike them, SH and CSH scale subquadratically which explains why they are faster for long enough sequences.

More specifically, for low error rates ($e=1\%$ and $e=5\%$), both SH and CSH with exact matches scale near-linearly with sequence length (best fit $n^{1.08}$ for $n \leq 10^7$), and the benefit from chaining is negligible.

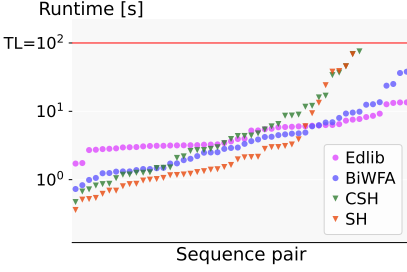
For high error rates ($e \geq 10\%$) and large n , the need to match the seeds inexactly causes a significant number of off-path matches. These off-path matches lower the value of the heuristics and prevent the punishment of suboptimal paths. This causes SH to expand a super-linear number of states (appendix A.6). Chaining the matches enforces them to follow the order of the seeds, which greatly reduces the negative effects of the off-path matches, leading to only linearly many expanded states. Nevertheless, the runtime scaling of CSH is super-linear as a result of the increasing fraction of time (up to 93% for $n=10^7$) spent on reordering states because of outdated f values after pruning.

PERFORMANCE table 1.2 compares the runtime and memory of A*PA to EDLIB and BiWFA for aligning a single pair of sequences of length $n=10^7$. A* with CSH is more than 250 times faster than EDLIB and BiWFA at $e=5\%$. For low error rate ($e \leq 5\%$), there is no significant performance difference between SH and CSH. The memory usage of A*PA for $e \leq 5\%$ is less than 600 MB for any n . At $e=15\%$, CSH uses at most 4.7 GB while SH goes out of memory (≥ 30 GB) because there are too many expanded states to store.

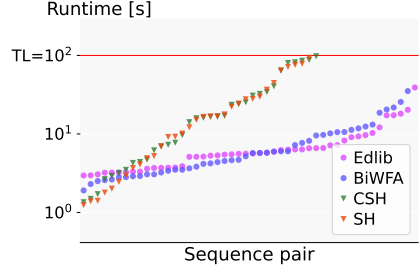
1.5.3 Comparison on real data

In this section we compare the exact optimal aligners on long ONT reads from our human datasets (fig. 1.5). With the presented minimalistic features, A*PA aligns some sequences faster than BiWFA and EDLIB, but the high runtime variance makes it slower overall (appendix A.7).

On the dataset without biological variation CHM13, SH is faster than BiWFA and EDLIB on 58% of the alignments (29 of 50). On the dataset with



(a) CHM13: ONT read errors



(b) NA12878: ONT read errors + biological mutations

FIGURE 1.5: Log plot comparison of the aligners' runtime on **real data**. The data points for each individual aligner are sorted by alignment time. Alignments that timed out after 100 seconds are not shown.

biological variation NA12878, SH outperforms BiWFA and EDLIB on 17% of the alignments (8 of 48) and in other cases is over an order of magnitude slower. In both datasets, SH and CSH time out for the sequences with the highest edit distances, because they have an error rate larger than the heuristic can handle efficiently ($e \geq r/k = 2/15 = 13.3\%$).

CSH usually explores fewer states than SH since chaining seed heuristic dominates the seed heuristic. However, in certain cases CSH is slower than SH since needs more time to update the heuristic after pruning (Step 3' in section 1.4.3).

1.6 DISCUSSION

Our graph-based approach to alignment differs considerably from dynamic programming approaches, mainly because of the ability to use information from the entire sequences. This additional information enables radically more focused path-finding at the cost of more complex algorithms.

LIMITATIONS Our presented method has several limitations:

1. *Complex regions trigger quadratic search.* Since it is unlikely that edit distance in general can be solved in strongly subquadratic time, it is inevitable that there are inputs for which our algorithm requires quadratic time. In particular, regions with high error rate, long indels, and too many matches (appendix A.8) are challenging and trigger quadratic exploration.

2. *High constant in runtime complexity.* Despite the near-linear scaling of the number of expanded states (appendix A.6), A*PA only outperforms EDLIB and BrWFA for sufficiently long sequences (fig. 1.4) due to the relatively high computational constant that the A* search induces.
3. *Complex parameter tuning.* The performance of our algorithm depends heavily on the sequences to be aligned and the corresponding choice of parameters (whether to use chaining, the seed length k , and whether to use inexact matches r). The parameter tuning (currently very simple (section 1.5.1) may require a more comprehensive framework when introducing additional optimizations.
4. *Real data.* The efficiency of the presented algorithm has high variability on real data (section 1.5.3) due to high error rates, long indels, and multiple repeats (demonstrated in fig. .8). Further optimizations are needed to align complex data.

FUTURE WORK We foresee a multitude of extensions and optimizations that may lead to efficient global aligning for production usage.

1. *Performance.* The practical performance of our A* approach could be improved using multiple existing ideas from the alignment domain: diagonal transition method, variable seed lengths, overlapping seeds, combining heuristics with different seed lengths, gap costs between matches in a chain [15, 27], more aggressive pruning, and better parameter tuning. More efficient implementations may be possible by using computational domains [23], bit-parallelization [16], and SIMD [19].
2. *Generalizations.* Our method can be generalized to more expressive cost models (non-unit costs, affine costs) and different alignment types (semi-global, ends-free, and possibly local alignment).
3. *Relaxations.* Abandoning the optimality guarantee enables various performance optimizations. Another relaxation of our algorithm would be to validate the optimality of a given alignment more efficiently than finding an optimal alignment from scratch.
4. *Analysis.* The near-linear scaling behaviour requires a thorough theoretical analysis [43]. The fundamental question that remains to be answered is: *What sequences and what errors can be tolerated while still scaling near-linearly with the sequence length?* We expect both theoretical and practical contributions to this question.

1.7 CONCLUSION

We presented an algorithm with an implementation in A*PA solving pairwise alignment between two sequences. The algorithm is based on A* with a seed heuristic, inexact matching, match chaining, and match pruning, which we proved to find an exact solution according to edit distance. For random sequences with up to 15% uniform errors, the runtime of A*PA scales near-linearly to very long sequences (10^7 bp) and outperforms other exact aligners. We demonstrate that on real ONT reads from a human genome, A*PA is faster than other aligners on only a limited portion of the reads.

ACKNOWLEDGEMENTS

The authors are grateful to Mykola Akulov for his help with producing figs. 1.1 and 1.3, and to Benjamin Bichsel, Maximilian Mordig and André Kahles for valuable feedback on various draft versions. We thank Sergey Nurk for his help with the preparation of the human data. Ragnar Groot Koerkamp is financed by ETH Grant ETH-17 21-1 to Gunnar Rätsch.

A APPENDIX

A.1 The seed heuristic and chaining seed heuristic are admissible

Theorem 3. *The seed heuristic h_{sh} and the chaining seed heuristic h_{csh} are admissible.*

Proof. Let $u = \langle i, j \rangle$ be any state in the edit graph $G(V, E)$, and let π be any path from u to the target t . A heuristic h is admissible if $h(u) \leq \text{cost}(\pi)$ for all $u \in V$ and for all paths π .

We first show that h_{csh} is admissible. Each remaining seed in $\mathcal{S}_{\geq i}$ is aligned somewhere by the path π . Since the seeds do not overlap, their shortest alignments in π do not have overlapping edges. Each edge in π has a non-negative cost, and hence $\sum_{s \in \mathcal{S}_i} \text{cost}(\pi_s) \leq \text{cost}(\pi)$. The alignment π_s of each seed either has cost less than r and equals a match in \mathcal{M}_s , or it has cost at least r otherwise. Together this implies that $h_{\text{csh}}^{\mathcal{M}}$ is a lower bound on the remaining cost $h^*(u)$:

$$\begin{aligned} h_{\text{csh}}^{\mathcal{M}}(u) &= \min_{\pi \in G(u \rightsquigarrow t)} \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \text{cost}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ r & \text{otherwise.} \end{cases} \\ &\leq \min_{\pi \in G(u \rightsquigarrow t)} \sum_{s \in \mathcal{S}_{\geq i}} \text{cost}(\pi_s) \leq \min_{\pi \in G(u \rightsquigarrow t)} \text{cost}(\pi) = h^*(u). \end{aligned} \quad (.11)$$

We now show that the seed heuristic $h_{\text{sh}}^{\mathcal{M}}(u)$ is bounded above by $h_{\text{csh}}^{\mathcal{M}}(u)$, so that it is also admissible:

$$\begin{aligned} h_{\text{sh}}^{\mathcal{M}}(u) &= \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \min_{m \in \mathcal{M}_s} \text{cost}(m) & \text{if } \mathcal{M}_s \neq \emptyset, \\ r & \text{otherwise.} \end{cases} \\ &\leq \sum_{s \in \mathcal{S}_{\geq i}} \min_{\pi \in G(u \rightsquigarrow t)} \begin{cases} \text{cost}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ r & \text{otherwise.} \end{cases} \\ &\leq \min_{\pi \in G(u \rightsquigarrow t)} \sum_{s \in \mathcal{S}_{\geq i}} \begin{cases} \text{cost}(\pi_s) & \text{if } \pi_s \in \mathcal{M}_s, \\ r & \text{otherwise.} \end{cases} = h_{\text{csh}}^{\mathcal{M}}(u). \end{aligned}$$

□

A.2 Partially admissible heuristic

Here, we generalize the concept of an *admissible* heuristic to that of a *partially admissible heuristic* that may change value as the A^* progresses, and even become inadmissible in some vertices. In ?? 4 we show that A^* still finds a shortest path when used with a partial admissible heuristic.

Definition 4 (Fixed vertex). *A fixed vertex is a closed vertex u to which A^* has found a shortest path, that is, $g(u) = g^*(u)$.*

A fixed vertex is never opened, and hence remains fixed.

Definition 5 (Partial admissible). *A heuristic \hat{h} is partially admissible when there exists a shortest path π^* from v_s to v_t for which at any point during the A^* and for any non-fixed vertex $u \in \pi^*$, the heuristic in u is a lower bound on the remaining distance $h^*(u)$: $\hat{h}(u) \leq h^*(u)$.*

We follow the structure and notation of Hart, Nilsson & Raphael [35] to prove that A^* finds a shortest path when used with a partially admissible heuristic.

Lemma 2 (Generalization of Hart, Nilsson & Raphael [35], Lemma 1). *Assuming v_t is not fixed, there exists an open vertex n' on π^* with $g(n') = g^*(n')$ such that all vertices in π^* following n' are not fixed.*

Proof. Let n^* be the last fixed vertex in π^* , and let n' be its successor on π^* . We have $g(n') = g^*(n')$ since π^* is a shortest path and n^* was expanded, and n' is open by our choice of n^* . □ □

Corollary 1 (Generalization of Hart, Nilsson & Raphael [35], Corollary to Lemma 1). *Suppose that \hat{h} is partially admissible, and suppose A^* has not terminated. Then n' as in ?? 2 has $f(n') \leq g^*(v_t)$.*

Proof. By definition of f we have $f(n') = g(n') + \hat{h}(n')$. Since $n' \in \pi^*$ and π^* does not contain any fixed vertices after n' , the partial admissibility of \hat{h} implies $\hat{h}(n') \leq h^*(n')$. Thus, $f(n') = g(n') + \hat{h}(n') = g^*(n') + \hat{h}(n') \leq g^*(n') + h^*(n') = g^*(v_t)$. □ □

Theorem 4. *A^* with a partially admissible heuristic finds a shortest path.*

Proof. The proof of Theorem 1 in Hart, Nilsson & Raphael [35] applies, with the remark that we use the specific path π^* and the specific vertex n' from ?? 1, instead of an arbitrary shortest path P . □ □

Remark 1 (Monotone heuristic). We call $\hat{h}(u)$ monotone when it does not decrease as the A^* progresses. To find the vertex u with minimal $f(u)$, we add step A4 to the A^* algorithm in section 1.4.1 to check whether the value of f in the priority queue is still up to date. If $f = f(u)$, we know that for each vertex (v, g_v, f_v) in the queue we have $f(u) \leq f_v \leq f(v)$ for all other open vertices v in the queue with stored value f_v , ensuring that u indeed has the minimal value of f value of all open vertices.

A.3 Match pruning preserves finding a shortest path

Lemma 3. For both \hat{h}_{sh} and \hat{h}_{csh} , any expanded state at the start of a seed is fixed.

Proof. Let \hat{h} be either \hat{h}_{sh} or \hat{h}_{csh} . We use a proof by contradiction. Thus, suppose that a state u at the start of a seed has minimal f among all open states, but the shortest path π^* from v_s to u has length $g^*(u) < g(u)$.

Let $n \neq u$ be the last expanded state on π^* , and let v be its successor. Let the unexpanded states of π^* following n that are at the start of a seed be $v \preceq w_0 \prec \dots \prec w_l = u$, in this order. We will show that $f(v) < f(u)$, resulting in a contradiction with the fact that A^* always expands the open state with minimal f :

$$\begin{aligned}
 f(v) &= g(v) + \hat{h}(v) = g^*(v) + \hat{h}(v) && \pi^* \text{ is a shortest path to } v, \\
 &\leq g^*(w_0) + \hat{h}(w_0) && \text{proven in part 1,} \\
 &\leq \dots \leq g^*(w_l) + \hat{h}(w_l) && \text{proven in part 2,} \\
 &= g^*(u) + \hat{h}(u) < g(u) + \hat{h}(u) = f(u) && \text{by contradiction.}
 \end{aligned}$$

Part 1: Proof of $g^*(v) + h(v) \leq g^*(w_0) + h(w_0)$. Since w_0 is the first start of a seed at or after v , the set of seeds following v is the same as the set of seeds following u . Thus, the summation in $\hat{h}(v)$ and $\hat{h}(w_0)$ in eq. (1.1) or eq. (1.2) is over the same seeds. Since $v \preceq w_0$, this implies $\hat{h}(v) \leq \hat{h}(w_0)$ in both cases, and $g^*(v) \leq g^*(w_0)$ since $v \preceq w_0$.

Part 2: Proof of $g^*(w_i) + \hat{h}(w_i) \leq g^*(w_{i+1}) + \hat{h}(w_{i+1})$. Note that w_i and w_{i+1} are at the start of the same or consecutive seeds. In the first case, $\hat{h}(w_i) \leq \hat{h}(w_{i+1})$ follows by definition since $w_i \preceq w_{i+1}$. Else, we have $\hat{h}(w_i) \leq \hat{h}(w_{i+1}) + r$. When $d := d(w_i, w_{i+1}) < r$, there is an unpruned

match m from w_i to w_{i+1} in $M - E$ of cost d , since w_i is not expanded. This implies $\hat{h}(w_i) \leq \hat{h}(w_{i+1}) + d$, and we obtain

$$\begin{aligned} g^*(w_i) + \hat{h}(w_i) &= (g^*(w_{i+1}) - d) + \hat{h}(w_i) \\ &\leq (g^*(w_{i+1}) - \min(d, r)) + (\hat{h}(w_{i+1}) + \min(d, r)) \\ &= g^*(w_{i+1}) + \hat{h}(w_{i+1}). \end{aligned} \quad \square$$

□

Lemma 4. *The heuristics \hat{h}_{csh} and \hat{h}_{sh} are partially-admissible.*

Proof. Let π^* be a shortest path from s to v_t . Let n be the last expanded state on π^* at the start of a seed. By ?? 3 n is fixed. By choice of n , no states on π^* following n are expanded, so no matches on π^* following n are pruned. The proof of ?? 3 applies to the path π^* without changes, implying that $h(u) \leq h^*(u)$ for all u on π^* following n , for both $h = \hat{h}_{\text{sh}}$ and $h = \hat{h}_{\text{csh}}$. □

Theorem 5. *A^* with match pruning heuristic \hat{h}_{sh} or \hat{h}_{csh} finds a shortest path.*

Proof. Follows from ?? 3 in appendix A.2 and ?? 4. □

A.4 Aligner versions and parameters

A*PA github.com/RagnarGrootKoerkamp/astar-pairwise-aligner (commit 550bbe5, and for human data commit 394b629), run as

```
astar-pairwise-aligner -i {input} --silent2
-k {k} -r {r} --algorithm {Dijkstra,SH,CSH}
[--no-prune]
```

BIWFA github.com/smarco/WFA2-lib (commit ffa2439), run as

```
align_benchmark -i {input} -a edit-wfa
--wfa-memory-mode ultralow
```

EDLIB We forked v1.2.7 at github.com/RagnarGrootKoerkamp/edlib (commit b34805c) to support WFA's input format. Run as

```
edlib-aligner -p -s {input}
```

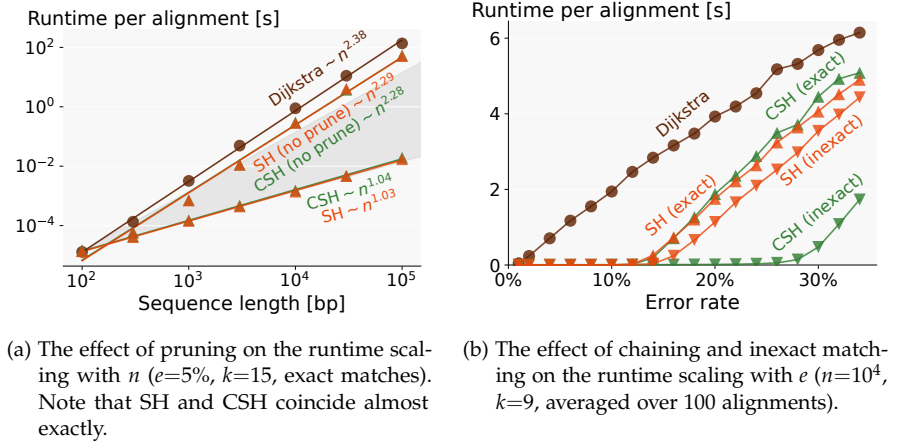


FIGURE .6: The effect of optimizations on runtime scaling on **synthetic data**.

A.5 Effect of pruning, chaining, and inexact matches

The optimizations and generalizations of the seed heuristic (section 1.2) impact the performance in a complex way. Here we aim to provide intuitive explanations (fig. .6).

PRUNING ENABLES NEAR-LINEAR SCALING WITH LENGTH fig. .6a shows that match pruning has a crucial effect on the runtime scaling with length for both SH and CSH. Essentially, this optimization changes the quadratic runtime to near-linear runtime. The pruned variants of SH and CSH are averaged over $\lfloor 10^7/n \rfloor$ sequence pairs, while the no-prune variants and Dijkstra are averaged over $\lfloor 10^5/n \rfloor$ pairs.

INEXACT MATCHING AND MATCH CHAINING ENABLE SCALING TO HIGH ERROR RATES fig. .6b shows that inexact matches can tolerate higher error rates. Because of the larger number of matches, chaining is needed to preserve the near-linear runtime. There are two distinctive modes of operation: the runtime is close to constant up to a certain error rate, after which the runtime grows linearly in e . Thus, our heuristics can direct the search up to a certain fraction of errors, after which does a Dijkstra-like exploration step for each additional error. A reasonable quantification of the effect of different optimizations is to mark the error rate at which the heuristic transfers to the second (slow) mode of operation. For $n=10^4$ and

$k=9$, Dijkstra starts a linear exploration at $e=0\%$, SH and CSH with exact matches start at around 12%, SH with inexact matches start at around 14%, and CSH with inexact matches start at around 27%.

A.6 Expanded states and equivalent band

e	n	SH				CSH			
		10^4	10^5	10^6	10^7	10^4	10^5	10^6	10^7
1%		1.08	1.08	1.08	1.08	1.07	1.07	1.07	1.07
5%		1.90	1.92	1.92	1.95	1.89	1.91	1.90	1.91
10%		2.85	2.88	3.16	16.6	2.79	2.80	2.81	2.82
15%		18.9	21.7	43.4	ML	18.5	20.1	20.4	20.3

TABLE .3: Equivalent band for aligning **synthetic sequences** of a given length and error rate. Each cell averages over $\lfloor 10^7/n \rfloor$ alignments. ML stands for exceeding the memory limit of 30 GB.

table .3 shows the *equivalent band* of each alignment: the number of expanded states divided by n . An equivalent band of 1 is the theoretical optimum for equal sequences, resulting from expanding only the states on the main diagonal. The equivalent band for CSH is always lower than the equivalent band for SH because the chaining seed heuristic dominates over seed heuristic. In practice, this difference becomes significant for large enough e and n ($e > 5\%$ and $n > 10^5$). The equivalent band of SH and CSH is constant for $e=1\%$, indicating a linear scaling with n . In other cases ($e \leq 5\%$ for SH and $e \leq 15\%$ for CSH), the band increases slowly with n , indicating a near-linear scaling.

A.7 Performance variation on human genome

Note that the BiWFA data points lie on a line, corresponding to the expected runtime of BiWFA of $O(s^2)$. Furthermore, the runtime of EDLIB can be seen to jump up around powers of two ($2^{14}=16\,384$, $2^{15}=32\,768$, $2^{16}=65\,536$), corresponding to the exponential search of the edit distance. EDLIB has a runtime complexity of $O(ns)$, and hence is slower than BiWFA for small edit distances, but faster than BiWFA for large edit distances.

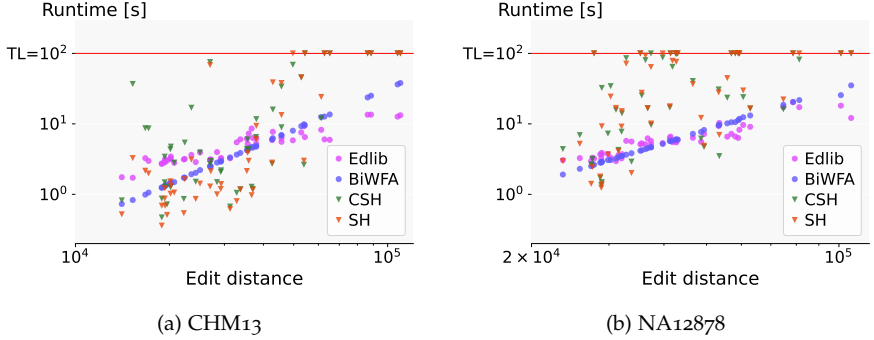


FIGURE .7: Log-log plots of the runtime for aligning **real sequences** from two datasets of varying error distance. The sequence length varies between 500 *kbp* and 840 *kbp* for the CHM13 reads, and between 502 *kbp* and 1 053 *kbp* for the NA12878 reads. Alignments that timed out after 100 seconds are shown at 100 s. Parameters used for SH and CSH are $k=15$ and $r=2$.

A.8 Complex alignments

Our algorithm finds optimal alignments very efficiently when both the sequence and the errors are uniformly random and the error rate is limited (section 1.5). Nevertheless, since the alignment problem is fundamentally unsolvable in strictly subquadratic time [7], there are sequences which cannot be aligned fast. We give three cases (fig. .8) where the complexity of our algorithm degrades.

1. *High error rate.* When the error rate becomes too high (larger than r/k), seeds can not increase the heuristic enough penalize all errors. Each unpenalized error increases f and needs more states to be searched, similar to Dijkstra. This can be mitigated by using shorter seeds and/or inexact matches, at the cost of introducing more matches.
2. *Long indel.* If there is a long insertion or deletion, the search has to accumulate the high cost for the long indel. Our heuristics do not account for long indels, and hence the search needs to expand states until the end of the gap is reached. Again, this causes a big increase of f and a search similar to Dijkstra. This may be improved in future work by introducing a *gap-cost* term to the chaining seed heuristic that penalizes gaps between consecutive matches in a chain.

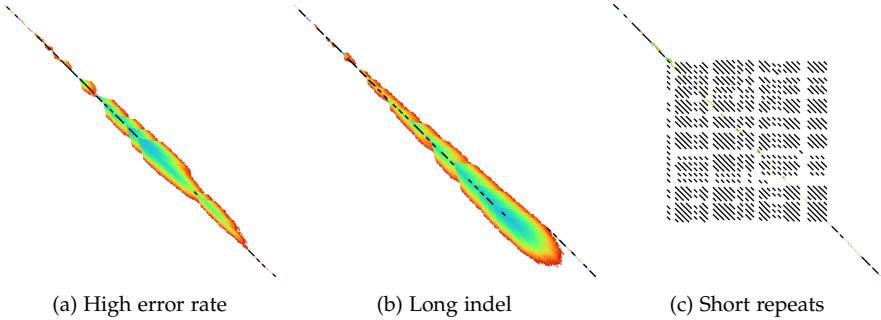


FIGURE .8: Expanded states by CSH ($r=2$, $k=10$) when aligning pairs of **synthetic sequences** ($n=1000$) with 8% random mutations containing (a) a region of length 200 with larger error rate than the heuristic can efficiently account for (50%), (b) a deletion of length 100, and (c) 60 mutated copies of a pattern of length 10. Seed matches are shown in black and occur on the best paths and in the repeated region. The order of expansion is shown by the color gradient from blue to red.

3. *Short repeats.* When a short pattern repeats many times, this can result in a quadratic number of matches. This makes the corresponding seeds ineffective for the seed heuristic, and slows down the computation and pruning of the chaining seed heuristic. This can be partially mitigated by increasing the seed length and/or exact matches, at the cost of reducing the potential of the heuristic. Alternatively, seeds with too many matches could be completely ignored.

A.9 Notation

table .4 summarizes the notation used in this work.

Object	Notation
Sequences	
Alphabet	$\Sigma = \{A, C, G, T\}$
Sequences	$A = \overline{a_0 a_1 \dots a_i \dots a_{n-1}} \in \Sigma^*$ $B = \overline{b_0 b_1 \dots b_j \dots b_{m-1}} \in \Sigma^*$
Subsequence	$a_{i..i'} := a_i a_{i+1} \dots a_{i'-1}$
Edit distance	$ed(A, B)$
Edit graph	
Graph	$G = (V, E)$
Vertices (states)	$u, v \in V = \{\langle i, j \rangle \mid 0 \leq i \leq n, 0 \leq j \leq m\}$
Edges	match/substitution $\langle i, j \rangle \rightarrow \langle i+1, j+1 \rangle$ deletion $\langle i, j \rangle \rightarrow \langle i+1, j \rangle$ insertion $\langle i, j \rangle \rightarrow \langle i, j+1 \rangle$
Distance	$d(u, v)$
Path, alignment	$\pi : (u \rightsquigarrow v)$
Shortest path	π^*
Cost	$\text{cost}(\pi)$
Preceding	$u \preceq v, m_1 \preceq m_2$
A*	
Start and target state	$v_s, v_t \in V$
Distance from v_s	$g^* = d(v_s, \cdot)$
Distance to v_t	$h^* = d(\cdot, v_t)$
Heuristic	h
Best distance from start	g
Estimated distance	$f = g + h$
Admissible heuristic	$h \leq h^*$
Consistent heuristic	$h(u) \leq d(u, v) + h(v)$
Seeds and matches	
Seed length and potential	k, r
Seeds	$s \in \mathcal{S}, s_l = A_{l..k..(l+1)..k}$
Seeds in suffix	$\mathcal{S}_{\geq i}$
Matches (per seed)	$m \in \mathcal{M}, \mathcal{M}_s$
Cost of match	$0 \leq \text{cost}(m) < r$
Score of match	$0 < \text{score}(m) = r - \text{cost}(m) \leq r$
Score of seed	$\text{score}(s) = \max_{m \in \mathcal{M}_s} \text{score}(m)$



SUMMARY

I dream my painting and I paint my dream.

— Vincent van Gogh

Summary here.



APPENDIX

Here be dragons.

BIBLIOGRAPHY

1. Maxwell, J. C. A Dynamical Theory of the Electromagnetic Field. *Philosophical Transactions of the Royal Society of London* **155**, 459 (1865).
2. Medvedev, P. Theoretical analysis of edit distance algorithms: an applied perspective. *arXiv preprint:2204.09535* (2022).
3. Ivanov, P., Bichsel, B. & Vechev, M. *Fast and Optimal Sequence-to-Graph Alignment Guided by Seeds* in *RECOMB 2022* (2022).
4. Needleman, S. B. & Wunsch, C. D. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology* **48**, 443 (1970).
5. Navarro, G. A guided tour to approximate string matching. *ACM computing surveys (CSUR)* **33**, 31 (2001).
6. Prijbelski, A. D., Korobeynikov, A. I. & Lapidus, A. L. in *Encyclopedia of Bioinformatics and Computational Biology: ABC of Bioinformatics* 292 (2018).
7. 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 symposium on Theory of computing* (2015), 51.
8. Kucherov, G. Evolution of biosequence search algorithms: a brief survey. *Bioinformatics* **35**, 3547 (2019).
9. Vintsyuk, T. K. Speech discrimination by dynamic programming. *Cybernetics* **4**, 52 (1968).
10. Sankoff, D. Matching sequences under deletion/insertion constraints. *Proceedings of the National Academy of Sciences* **69**, 4 (1972).
11. Sellers, P. H. On the theory and computation of evolutionary distances. *SIAM Journal on Applied Mathematics* **26**, 787 (1974).
12. Wagner, R. A. & Fischer, M. J. The string-to-string correction problem. *Journal of the ACM (JACM)* **21**, 168 (1974).
13. Reinert, K., Dadi, T. H., Ehrhardt, M., Hauswedell, H., Mehringer, S., Rahn, R., Kim, J., Pockrandt, C., Winkler, J., Siragusa, E., *et al.* The SeqAn C++ template library for efficient sequence analysis: a resource for programmers. *Journal of biotechnology* **261**, 157 (2017).

14. Daily, J. Parasail: SIMD C library for global, semi-global, and local pairwise sequence alignments. *BMC Bioinformatics* **17**, 1 (2016).
15. Ukkonen, E. Algorithms for approximate string matching. *Information and control* **64**, 100 (1985).
16. Myers, G. A fast bit-vector algorithm for approximate string matching based on dynamic programming. *Journal of the ACM (JACM)* **46**, 395 (1999).
17. Šošić, M. & Šikić, M. Edlib: a C/C++ library for fast, exact sequence alignment using edit distance. *Bioinformatics* **33**, 1394 (2017).
18. Myers, E. W. An $O(ND)$ difference algorithm and its variations. *Algorithmica* **1**, 251 (1986).
19. Marco-Sola, S., Moure, J. C., Moreto, M. & Espinosa, A. Fast gap-affine pairwise alignment using the wavefront algorithm. *Bioinformatics* **37**, 456 (2021).
20. Gotoh, O. An improved algorithm for matching biological sequences. *Journal of molecular biology* **162**, 705 (1982).
21. Marco-Sola, S., Eizenga, J. M., Guarracino, A., Paten, B., Garrison, E. & Moreto, M. Optimal gap-affine alignment in $O(s)$ space. *bioRxiv* (2022).
22. Hirschberg, D. S. A linear space algorithm for computing maximal common subsequences. *Communications of the ACM* **18**, 341 (1975).
23. Spouge, J. L. Speeding up dynamic programming algorithms for finding optimal lattice paths. *SIAM Journal on Applied Mathematics* **49**, 1552 (1989).
24. Ivanov, P., Bichsel, B., Mustafa, H., Kahles, A., Rätsch, G. & Vechev, M. T. *AStarix: Fast and Optimal Sequence-to-Graph Alignment* in RECOMB 2020 (2020).
25. Benson, G., Levy, A. & Shalom, R. *Longest Common Subsequence in k -length substrings* 2014.
26. Levenshtein, V. I. *et al. Binary codes capable of correcting deletions, insertions, and reversals* in *Soviet physics doklady* **10** (1966), 707.
27. Wilbur, W. J. & Lipman, D. J. The context dependent comparison of biological sequences. *SIAM Journal on Applied Mathematics* **44**, 557 (1984).
28. Benson, G., Levy, A., Maimoni, S., Noifeld, D. & Shalom, B. R. LCSk: a refined similarity measure. *Theoretical Computer Science* **638**, 11 (2016).

29. Poole, D. L. & Mackworth, A. K. *Artificial Intelligence: Foundations of Computational Agents* Second (Cambridge University Press, 2017).
30. Koenig, S. & Likhachev, M. *Real-time adaptive A* in Proceedings of the fifth international joint conference on Autonomous agents and multiagent systems* (2006), 281.
31. Hirschberg, D. S. Algorithms for the longest common subsequence problem. *Journal of the ACM (JACM)* **24**, 664 (1977).
32. Hunt, J. W. & Szymanski, T. G. A fast algorithm for computing longest common subsequences. *Communications of the ACM* **20**, 350 (1977).
33. Pavetić, F., Katanić, I., Matula, G., Žužić, G. & Šikić, M. Fast and simple algorithms for computing both LCSk and LCSk+. *arXiv preprint:1705.07279* (2017).
34. Dijkstra, E. W. A note on two problems in connexion with graphs. *Numerische mathematik* **1**, 269 (1959).
35. Hart, P. E., Nilsson, N. J. & Raphael, B. A formal basis for the heuristic determination of minimum cost paths. *IEEE transactions on Systems Science and Cybernetics* **4**, 100 (1968).
36. Pearl, J. *Heuristics: intelligent search strategies for computer problem solving* (Addison-Wesley Longman Publishing Co., Inc., 1984).
37. Holte, R. C. *Common misconceptions concerning heuristic search in Third Annual Symposium on Combinatorial Search* (2010).
38. Deorowicz, S. & Grabowski, S. Efficient algorithms for the longest common subsequence in k-length substrings. *Information Processing Letters* **114**, 634 (2014).
39. Bertsekas, D. P. *Linear network optimization: algorithms and codes* (MIT Press, 1991).
40. Allison, L. Lazy dynamic-programming can be eager. *Information Processing Letters* (1992).
41. Nurk, S., Koren, S., Rhie, A., Rautiainen, M., *et al.* The complete sequence of a human genome. *Science* **376**, 44 (2022).
42. Bowden, R., Davies, R. W., Heger, A., Pagnamenta, A. T., de Cesare, M., Oikkonen, L. E., Parkes, D., Freeman, C., Dhalla, F., Patel, S. Y., *et al.* Sequencing of human genomes with nanopore technology. *Nature communications* **10**, 1 (2019).
43. Medvedev, P. The limitations of the theoretical analysis of applied algorithms. *arXiv preprint:2205.01785* (2022).

CURRICULUM VITAE

PERSONAL DATA

Name	Albert Einstein
Date of Birth	March 14, 1879
Place of Birth	Ulm, Germany
Citizen of	Switzerland

EDUCATION

1896 – 1900	Eidgenössisches Polytechnikum, Zürich, Switzerland <i>Final degree:</i> Diploma
1895 – 1896	Aargauische Kantonsschule (grammar school) Aarau, Switzerland <i>Final degree:</i> Matura (university entrance diploma)
– July 1894	Luitpold-Gymnasium (grammar school) Munich, Germany

EMPLOYMENT

June 1902 –	Technical Expert, III Class <i>Federal Office for Intellectual Property,</i> Bern, Switzerland
-------------	--

PUBLICATIONS

Articles in peer-reviewed journals:

1. Einstein, A. Über die von der molekularkinetischen Theorie der Wärme geforderte Bewegung von in ruhenden Flüssigkeiten suspendierten Teilchen. *Annalen der Physik* **322**, 549 (1905).
2. Einstein, A. Zur Elektrodynamik bewegter Körper. *Annalen der Physik* **322**, 891 (1905).

Conference contributions:

3. Einstein, A. *Implications of a fixed vacuum speed of light in Relativity* Oct. 2–6, 1905 (1st Conference on Special Relativity, Zurich, Switzerland).