

Subquadratic optimal alignment using A*

Pesho Ivanov and Ragnar Groot Koerkamp

ETH Zurich

 **SRILAB**

<https://sri.inf.ethz.ch/>



[eth-sri/astarix](https://github.com/eth-sri/astarix)



[RagnarGrootKoerkamp/astar-pairwise-aligner](https://github.com/RagnarGrootKoerkamp/astar-pairwise-aligner)

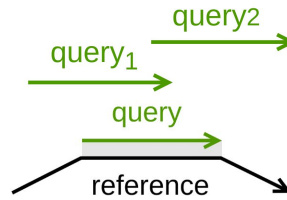
 **BIOMEDICAL
INFORMATICS**

<https://bmi.inf.ethz.ch/>

ETH zürich

Pairwise alignment variants

Semi-global /
Mapping



now

Global



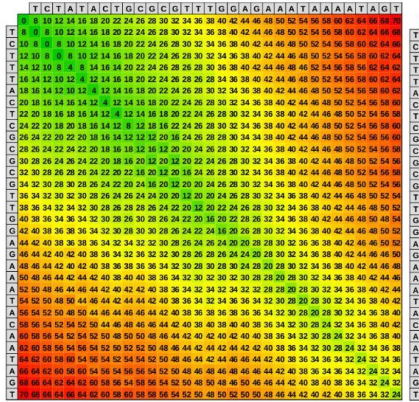
in 15min

Local



some day?

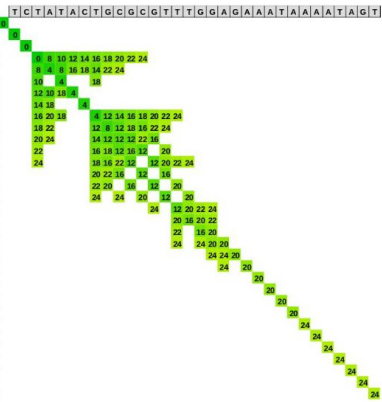
Optimal alignment algorithms



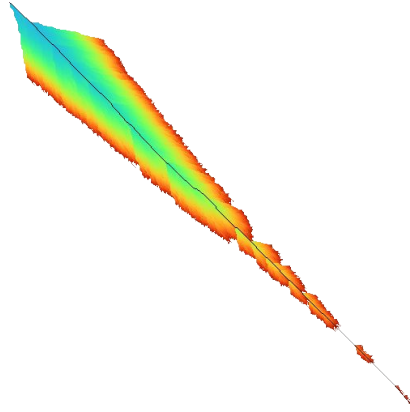
NW
 $O(n^2)$

WFA
 $O(sn) \equiv O(en^2)$

Uninformed
search

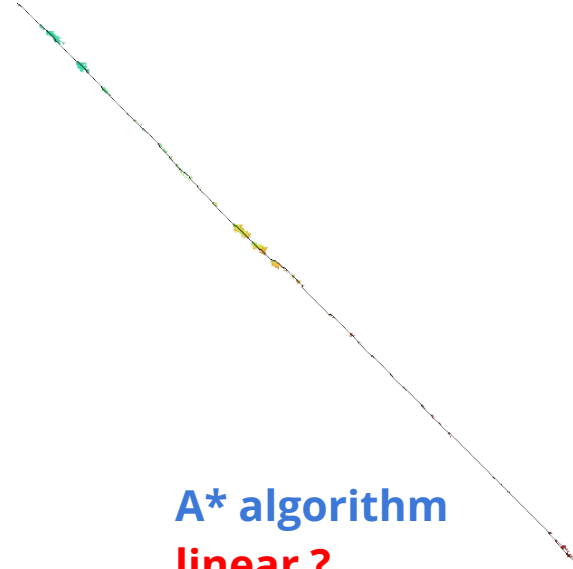


Dijkstra
 $O(n^2)$



A* algorithm
linear?

Informed
search



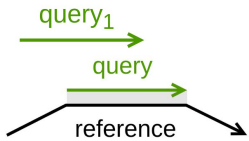
Results

Core approach

Empirical runtime

Tool

Sequence-to-graph



trie / suffix tree Ivanov et al. (2020)
seed heuristic Ivanov et al. (2022)

$O(mN) \rightarrow m^{1.2}N^{0.4} ?$

 AStarix

Global aligning



seed heuristic pruning Groot Koerkamp & Ivanov (draft)

$O(en^2) \rightarrow n ?$

 A*PA

m - query length

N - reference size

n - sequences length

e - error rate

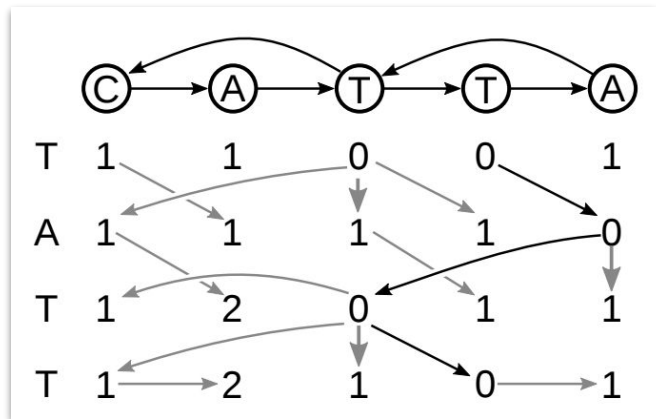
A* algorithm and admissible heuristic



A* heuristic:

- Admissible
- Precise
- Efficient to compute

Sequence-to-graph mappers



Edit graph

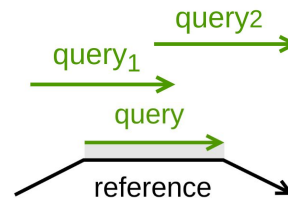
$dp(u,i)$ — Minimal cost to align $q[1..i]$
from node u in the reference

Mikko Rautiainen and Tobias Marschall (2017) —
Aligning sequences to general graphs in $O(V+mE)$ time



Except of **AStarix** other **optimal** mappers are **quadratic**:
GraphAligner, **Vargas**, **PaSGAL**

A* for sequence-to-graph



Scaling with

Reference size

A* + trie

$N^{0.11-0.46}$

Query length

seed heuristic

$m^{1.21}$

Ivanov P, Bichsel B, Mustafa H, Kahles A, Rätsch G, Vechev M (RECOMB 2020) – AStarix: Fast and Optimal Sequence-to-Graph Alignment

Ivanov P, Bichsel B, Vechev M (RECOMB 2022) – Fast and Optimal Sequence-to-Graph Alignment Guided by Seeds

Fastest (consistently >60x) for:

- Illumina reads (200bp, 2-4% error rate)
- HiFi reads (up to 25kbp, 0.3% error rate)

Index the reference



Prepare for a query

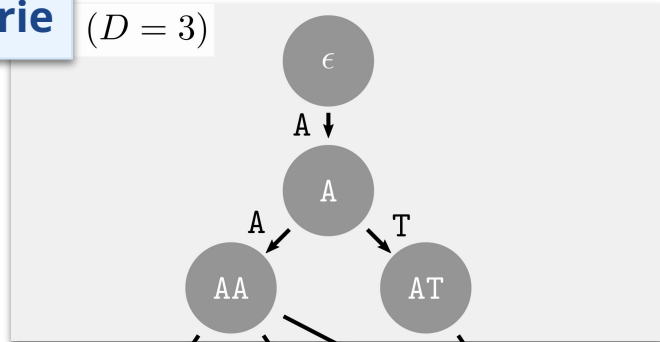


Align the query

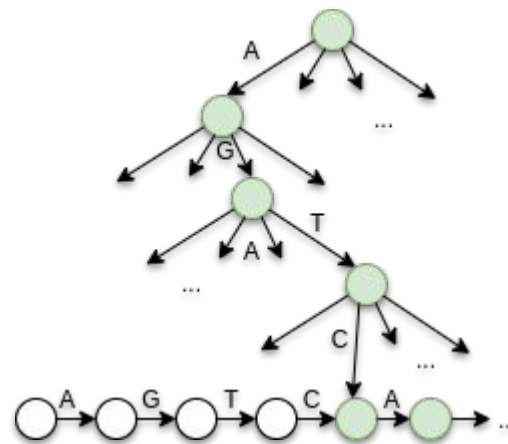
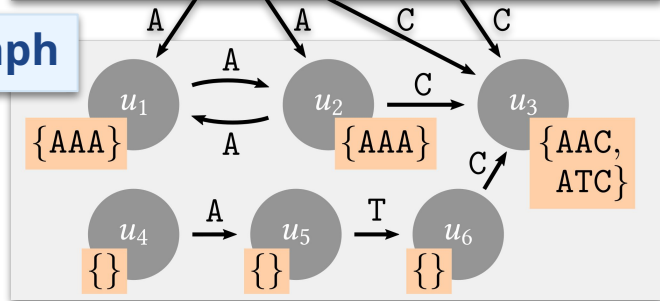


1. Index the reference

Trie ($D = 3$)



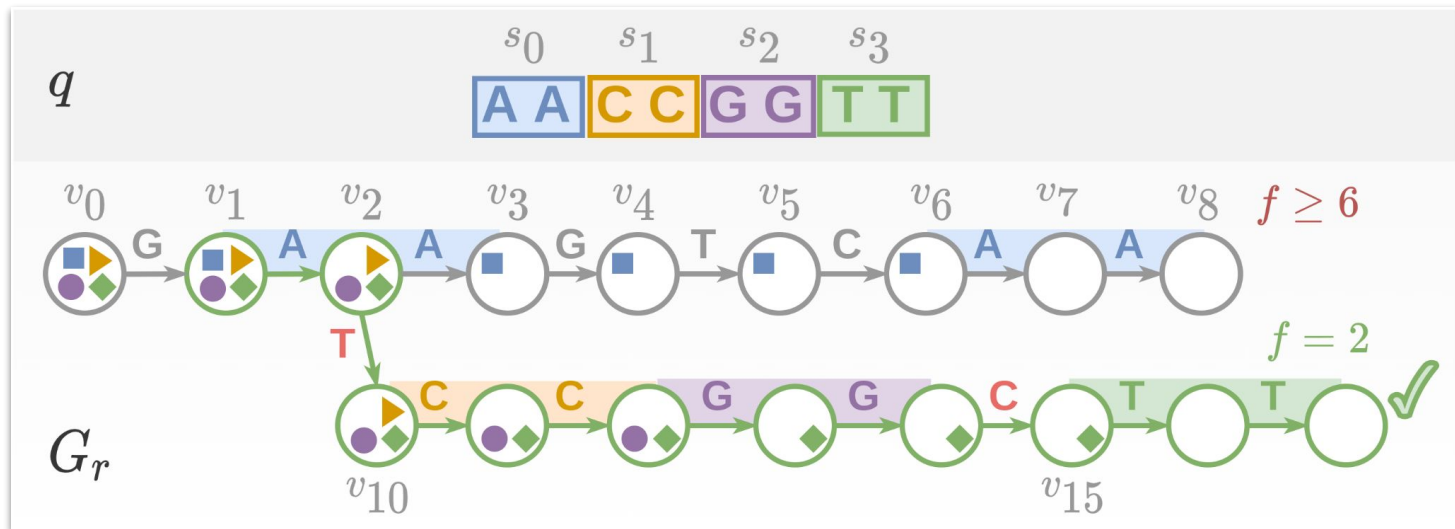
Reference graph



Aligning with a good heuristic

Dox and Fostier (Master's thesis, 2018), Efficient algorithms for pairwise sequence alignment on graphs
Ivanov et al. (RECOMB 2020), AStarix: Fast and Optimal Sequence-to-Graph Alignment

Seed heuristic: preparation

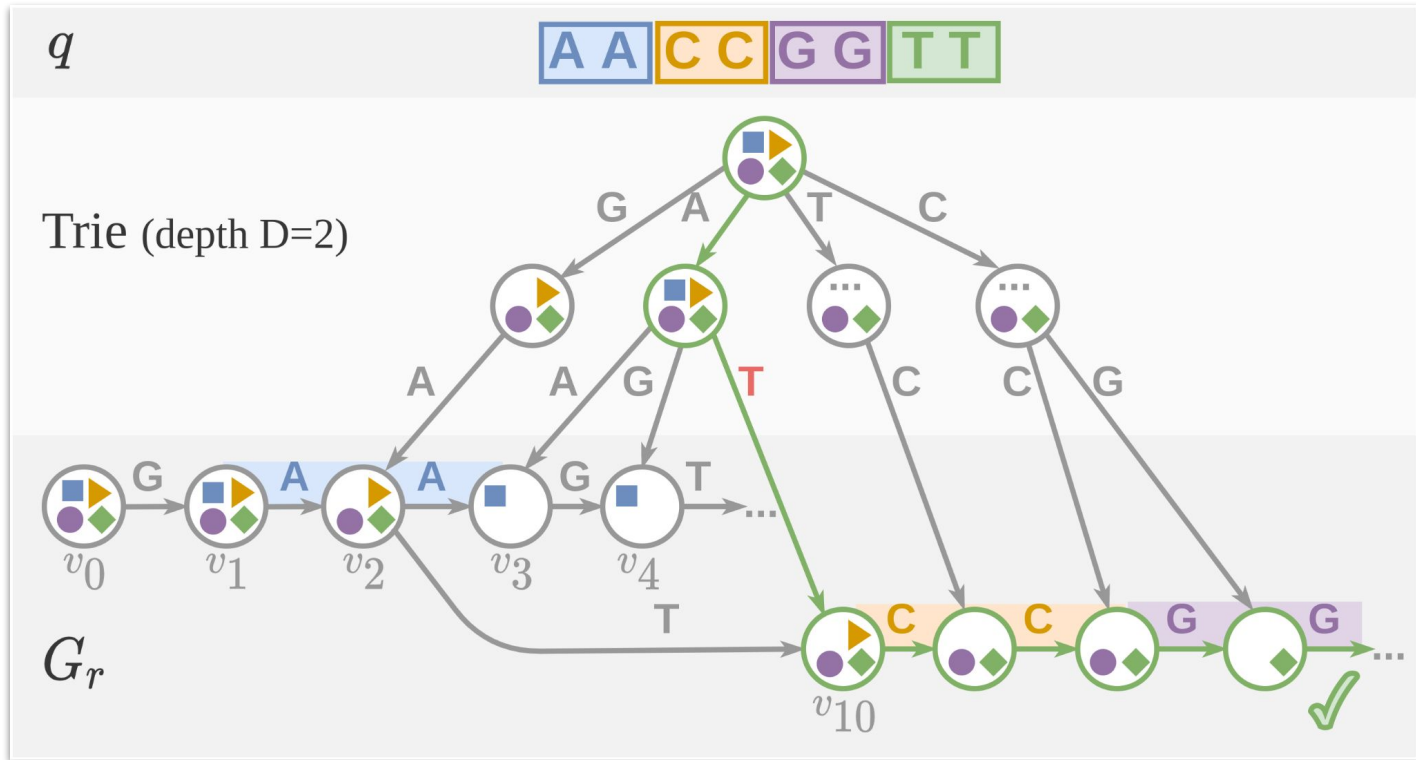


Query → **Seeds** → **Matches** → **Crumbs**

Ivanov, Bichsel and Vechev (RECOMB 2022) –

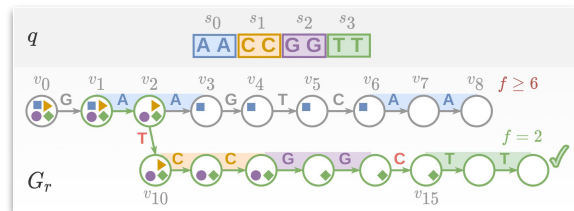
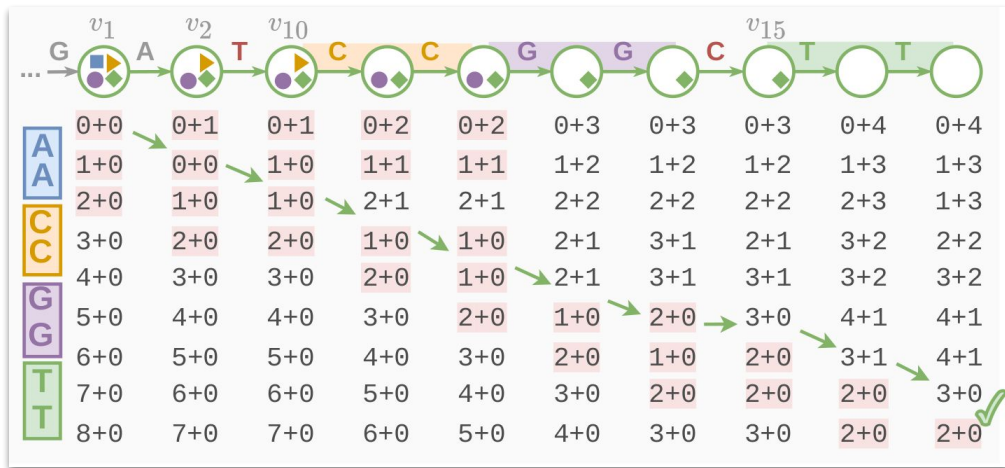
Fast and Optimal Sequence-to-Graph Alignment Guided by Seeds

Seed heuristic on trie



Query \rightarrow Seeds \rightarrow Matches \rightarrow Crumbs

Seed heuristic: query



$$h\langle v, i \rangle = \underbrace{\text{misses}\langle v, i \rangle}_{\left| \{s \in \text{Seeds}_{\geq i} \mid s \text{ has no crumb in } v\} \right|} \cdot \underbrace{\delta_{\min}}_{\min(\Delta_{\text{subst}}, \Delta_{\text{del}}, \Delta_{\text{ins}})}$$

$\left| \{s \in \text{Seeds}_{\geq i} \mid s \text{ has no crumb in } v\} \right|$

$\min(\Delta_{\text{subst}}, \Delta_{\text{del}}, \Delta_{\text{ins}})$

Speedup results

Tool	Illumina		HiFi		
	<i>E. coli</i>	MHC	<i>E. coli</i>	MHC	
Seed heuristic	0.019	0.041	0.001	0.002	s/kbp
Prefix heuristic	269x	180x	n/a	n/a	x slowdown
GRAPHALIGNER	424x	212x	118x	64x	
VARGAS	133x	67x	1 413x	705x	
PASGAL	263x	130x	1 367x	736x	

Seed heuristic
skips >99.99%
of the table cells

References: *E. coli* – linear 4.6M, Major Histocompatibility Complex (MHC) – 5.3M

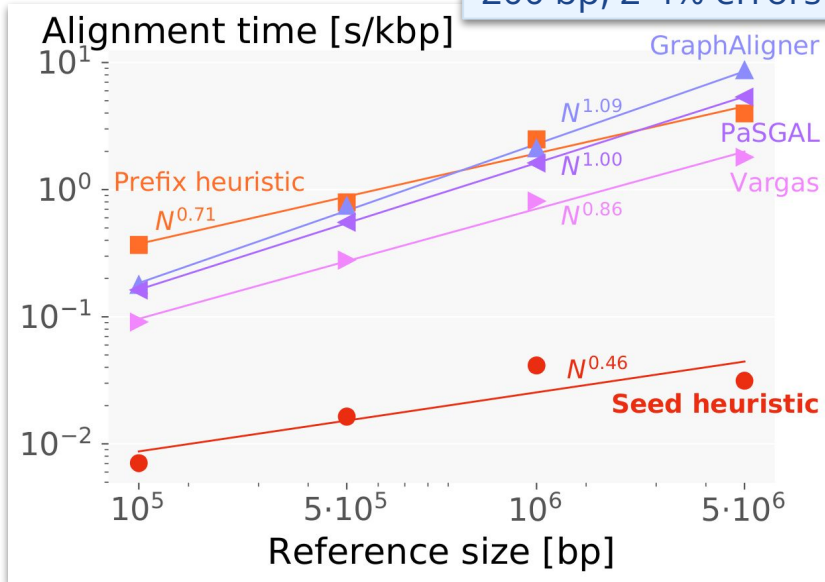
Simulated queries: 200bp Illumina with $\Delta=(0,1,5,5)$; HiFi: 5–25kbp, $e=0.3\%$, $\Delta=(0,1,1,1)$

Prefix heuristic parameters: length cap $d=5$, cost cap $c=5$, trie depth $D=\log(N)$

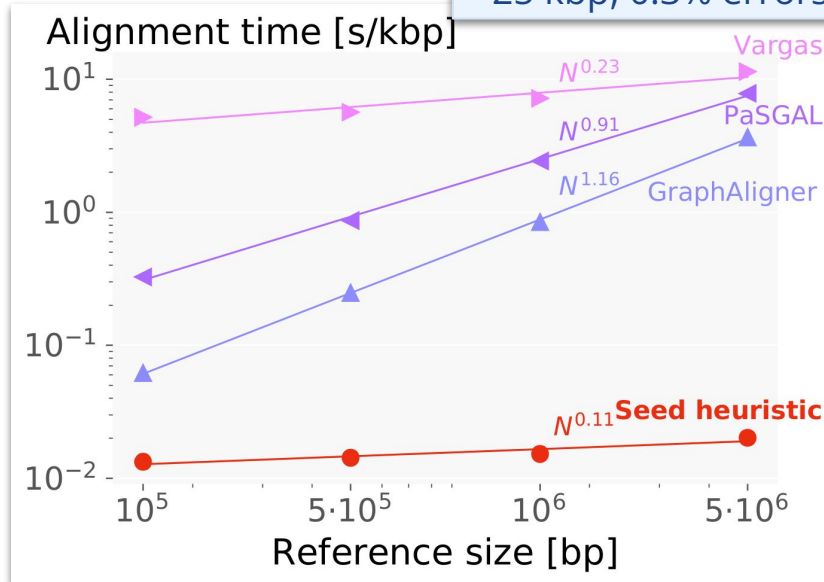
Seed heuristic parameters: $D=14 \approx \log_4 N$; $k=25$ for Illumina, $k=150$ for HiFi reads

Scaling: reference size

Illumina reads:
200 bp, 2-4% errors



HiFi reads:
<25 kbp, 0.3% errors



References: *E. coli* – linear 4.6M, Major Histocompatibility Complex (MHC) – 5.3M

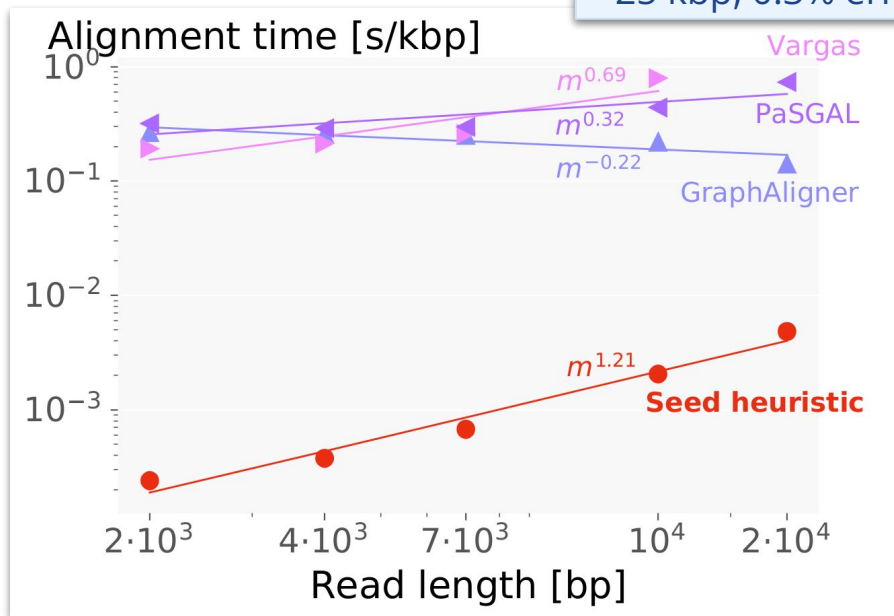
Simulated queries: 200bp Illumina with $\Delta=(0,1,5,5)$; HiFi: 5–25kbp, $e=0.3\%$, $\Delta=(0,1,1,1)$

Prefix heuristic parameters: length cap $d=5$, cost cap $c=5$, trie depth $D=\log(N)$

Seed heuristic parameters: $D=14 \approx \log_4 N$; $k=25$ for Illumina, $k=150$ for HiFi reads

Scaling: query length

HiFi reads:
<25 kbp, 0.3% errors



References: Major Histocompatibility Complex (MHC) – 5.3M

Simulated queries: HiFi: 2–25kbp, $e=0.3\%$, $\Delta=(0,1,1,1)$

Prefix heuristic parameters: length cap $d=5$, cost cap $c=5$, trie depth $D=\log(N)$

Seed heuristic parameters: $D=14 \approx \log_4 N$; $k=150$ for HiFi reads

Future work

Applicability:

- longer indels
- complex references

Performance:

- memory-efficiency
- parallelization
- optimize for linear references

Extensions:

- local alignment
- affine costs

Theoretical analysis:

- scaling proofs

Fast and Optimal Sequence-to-Graph Alignment Guided by Seeds

Pesho Ivanov, Benjamin Bichsel and Martin Vechev

ETH Zurich

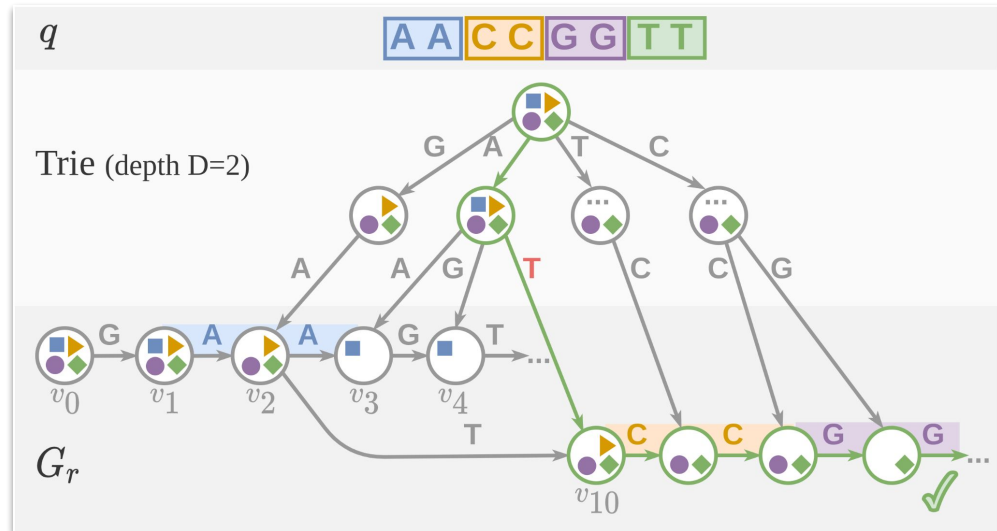
 **SRILAB**

<https://sri.inf.ethz.ch/>



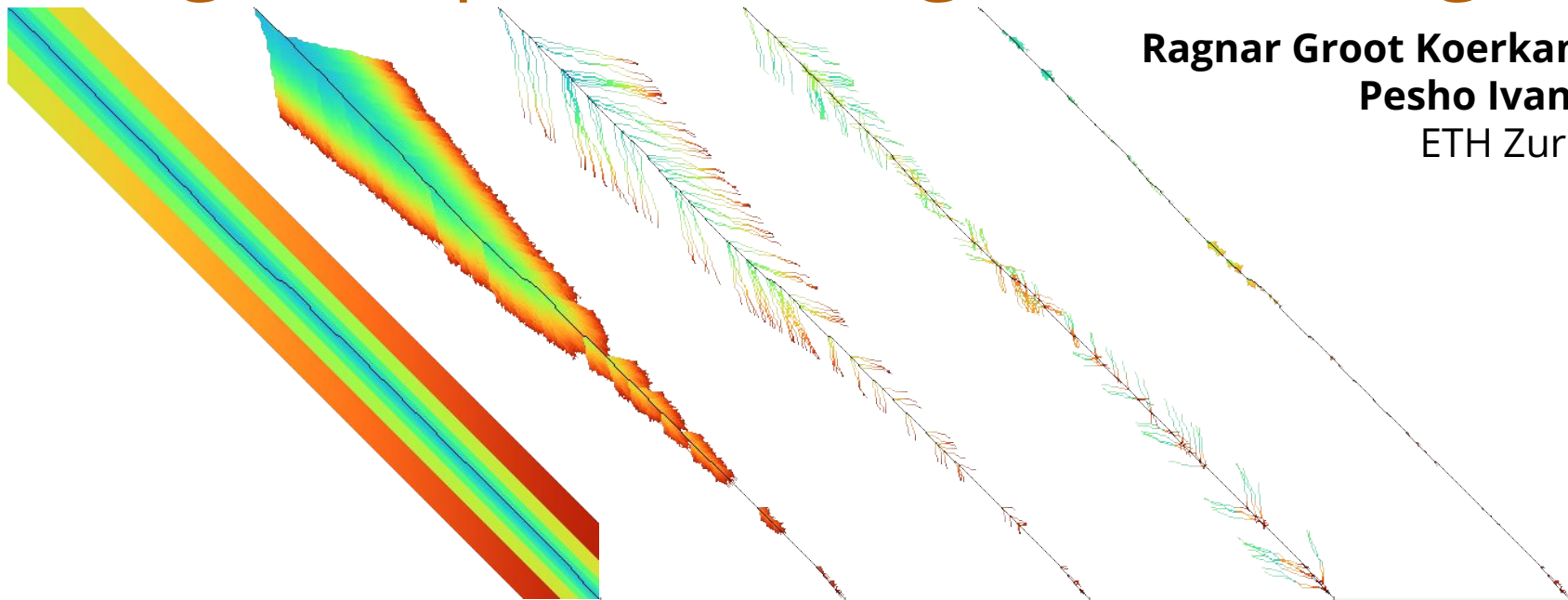
[eth-sri/astarix](https://github.com/eth-sri/astarix)

ETH zürich



Exact global pairwise alignment using A^*

Ragnar Groot Koerkamp
Pesho Ivanov
ETH Zurich



Exponential band
Ukkonen'85

Edlib
Šošić, Šikić'17

$O(ns)$, $O(n)$

Dijkstra
Ukkonen'85

-

$O(ns)$

Diagonal Transition + Divide & Conquer
Ukkonen'85,
Myers'86

WFA
Marco-Sola et al'20

$O(s^2)$

BiWFA
Marco-Sola et al'22

$O(s^2)$, $O(s)$

A^* + SH + pruning
Groot Koerkamp,
Ivanov'22

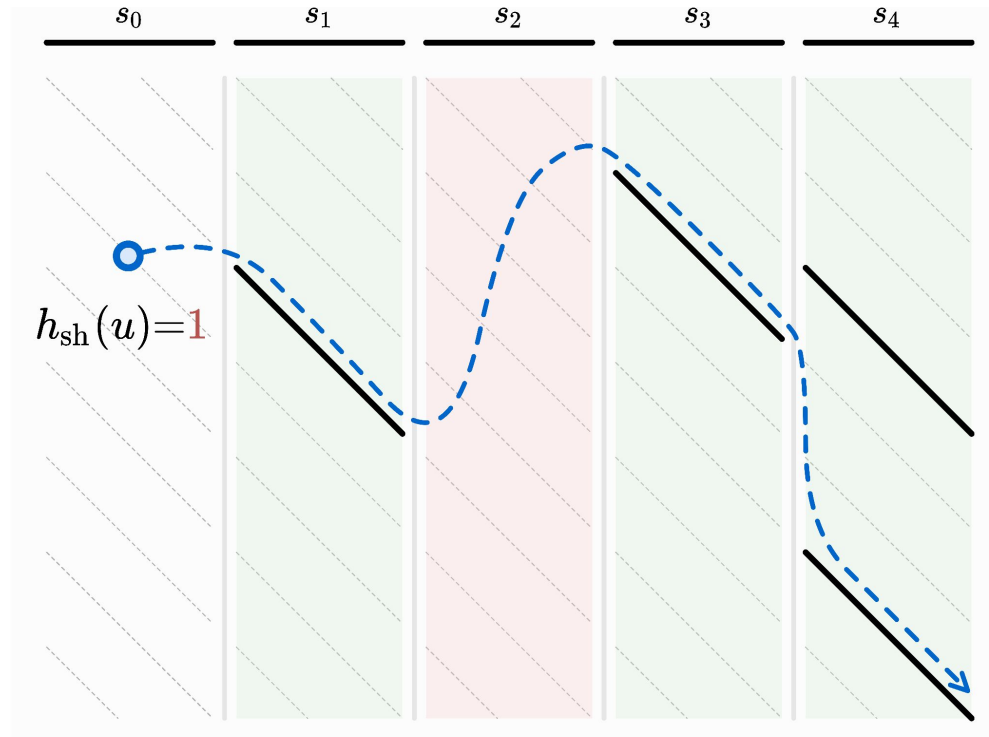
A^*PA

$O(n)?$

$n=500$, $e=20\%$

Visualization by Mykola Akulov

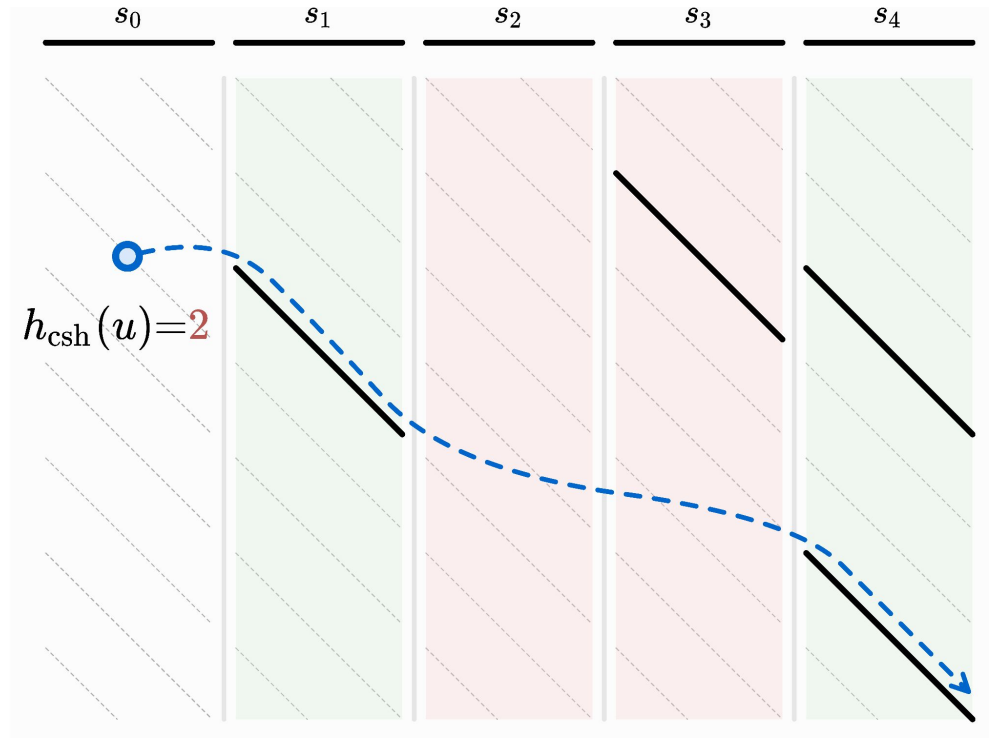
Seed Heuristic (SH)



Number of seeds without a match:

$$h_{sh} = \#\{\text{upcoming seeds}\} - \#\{\text{matching upcoming seeds}\}$$

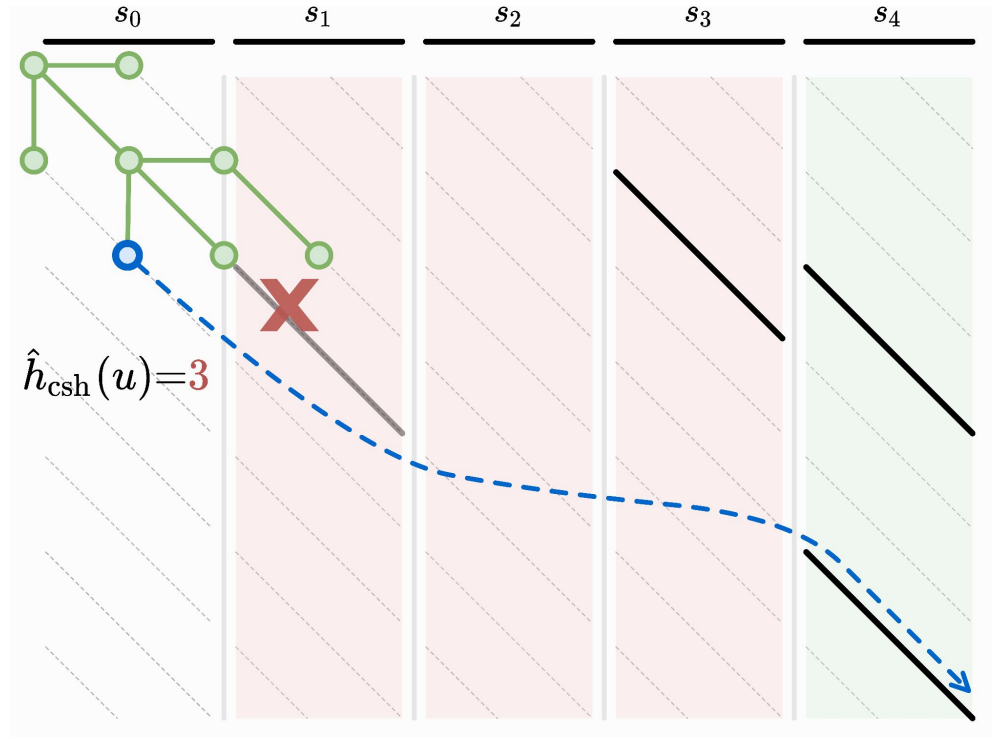
Chained Seed Heuristic (CSH)



Require matches to form a chain:

$$h_{\text{csh}} = \#\{\text{upcoming seeds}\} - \#\{\text{longest **chain** of matches}\}$$

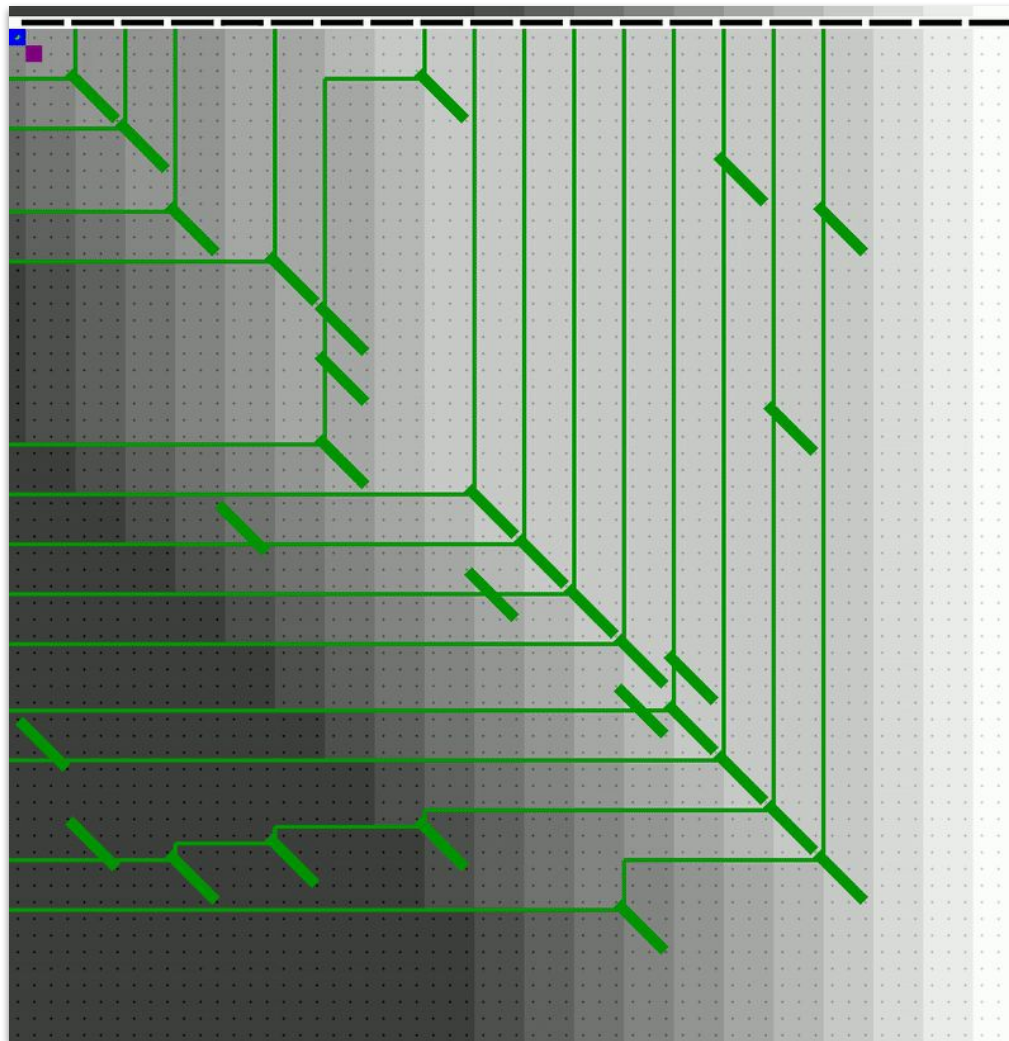
Match pruning



Exclude paths containing visited states => Prune visited matches.

$$\hat{h}_{\text{csh}} = \#\{\text{upcoming seeds}\} - \#\{\text{longest chain of **unpruned** matches}\}$$

Demo



Blue: expanded

Purple: explored

Lime path: current state

Background: heuristic

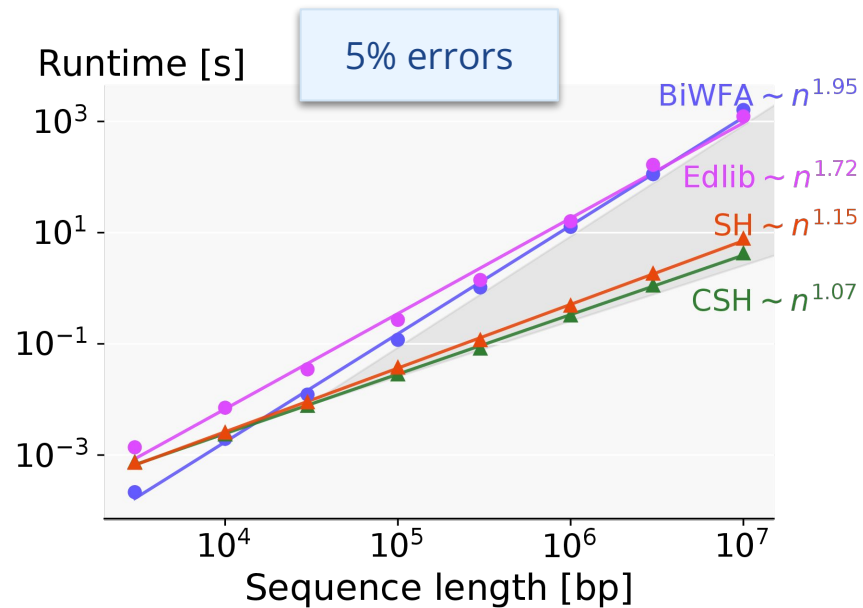
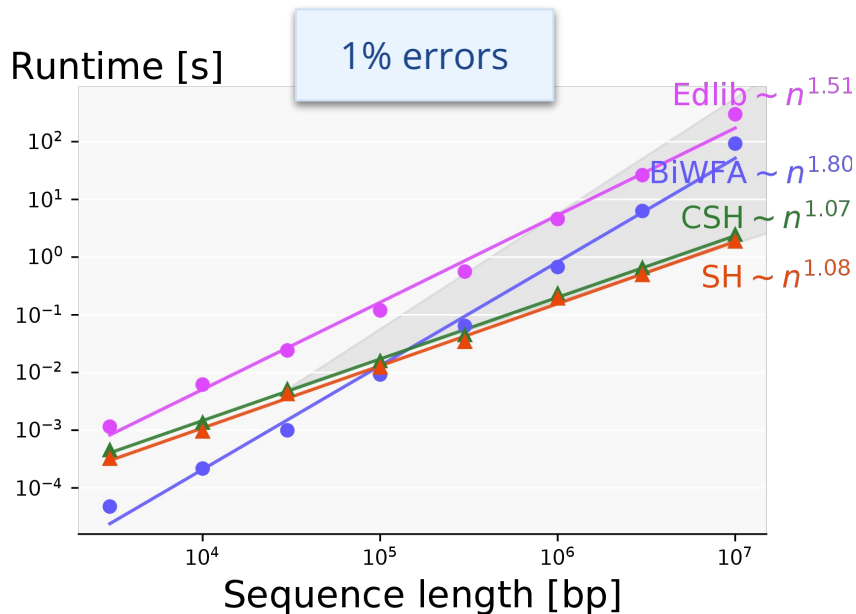
Black: seeds

Green: seed matches

Green lines: contours

Lime matches: pruned

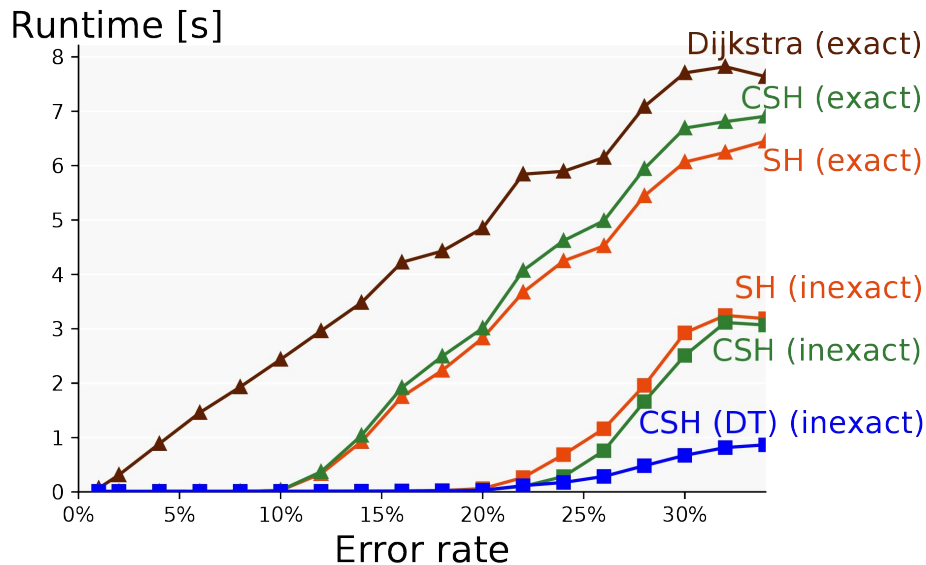
Comparison on random data



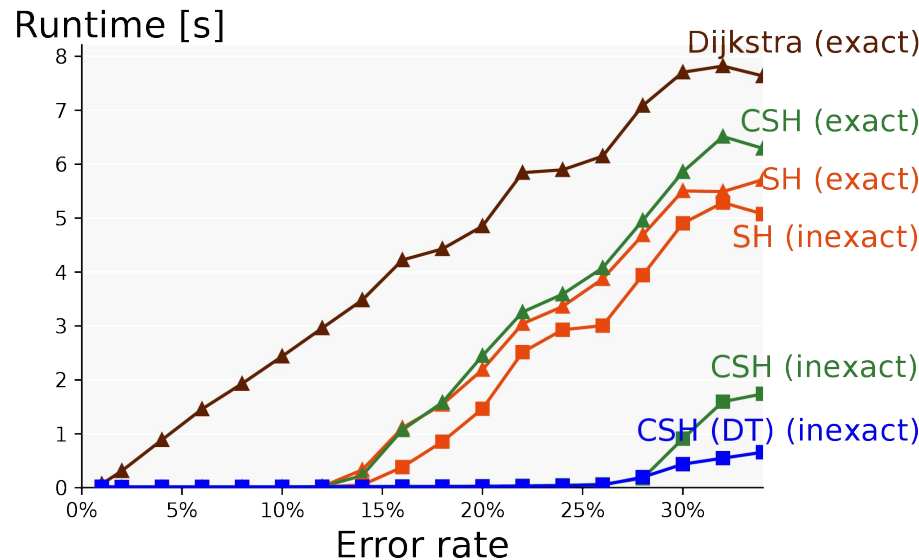
Seed length $k=15$
Exact matches

Randomness is important!
On human data:
~2x faster to ~100x slower

Scaling with error rate



$n=10^4$, seed length $k=11$



$n=10^4$, seed length $k=9$

Limitations & future work

Limitations

- Sequences must not be repetitive
- Mutations must be uniform random: No long indels!

Performance

- Use diagonal transition (work in progress)
- Variable seed length
- Gap cost for joining chains

Scope

- More cost models (affine costs)
- Semi-global alignment
- Extend to sequence-to-graph alignment

Prove expected linear time on random input

Q&A: A* for Optimal Sequence Alignment

Pesho Ivanov and Ragnar Groot Koerkamp

ETH Zurich



[eth-sri/astarix](https://github.com/eth-sri/astarix)

[@peshotrie](https://twitter.com/peshotrie)



[RagnarGrootKoerkamp/astar-pairwise-aligner/](https://github.com/RagnarGrootKoerkamp/astar-pairwise-aligner/)

[@curious_coding](https://twitter.com/curious_coding)

research.curiouscoding.nl

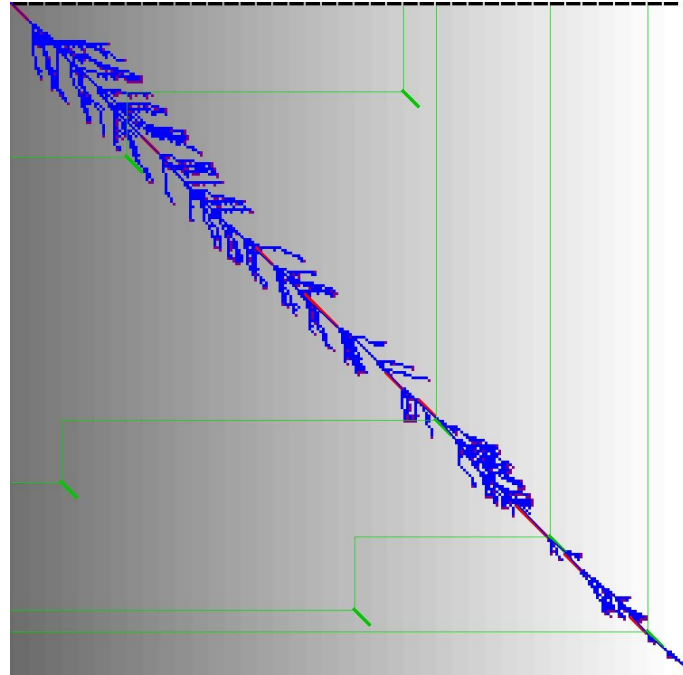
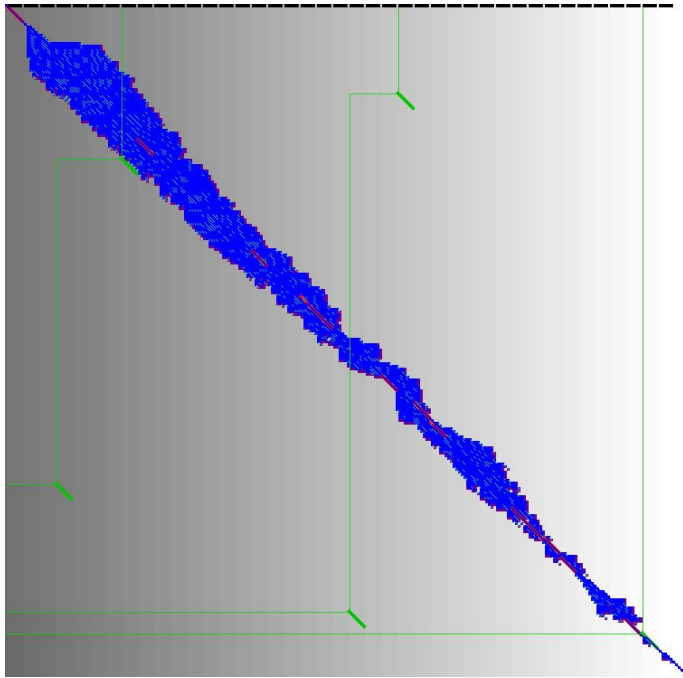
Q&A

	DP Needleman, Wunsch'69	Exponential band Ukkonen'85	Dijkstra Ukkonen'85	Diagonal Transition + Divide & Conquer Ukkonen'85, Myers'86		A* + SH + pruning Groot Koerkamp, Ivanov'22
		Edlib Šošić, Šikić'17		WFA Marco-Sola et al'20	BiWFA Marco-Sola et al'22	A*PA
Expected [#] runtime	$O(n^2)$	$O(ns)$	$O(ns)$	$O(n+s^2)$	$O(n+s^2)$	$O(n)^{*?}$
Expected [#] memory	$O(n^2)$	$O(ns) / O(n)$	$O(ns)$	$O(s^2)$	$O(s)$	$O(n)^{*?}$

[#]: random string with random mutations

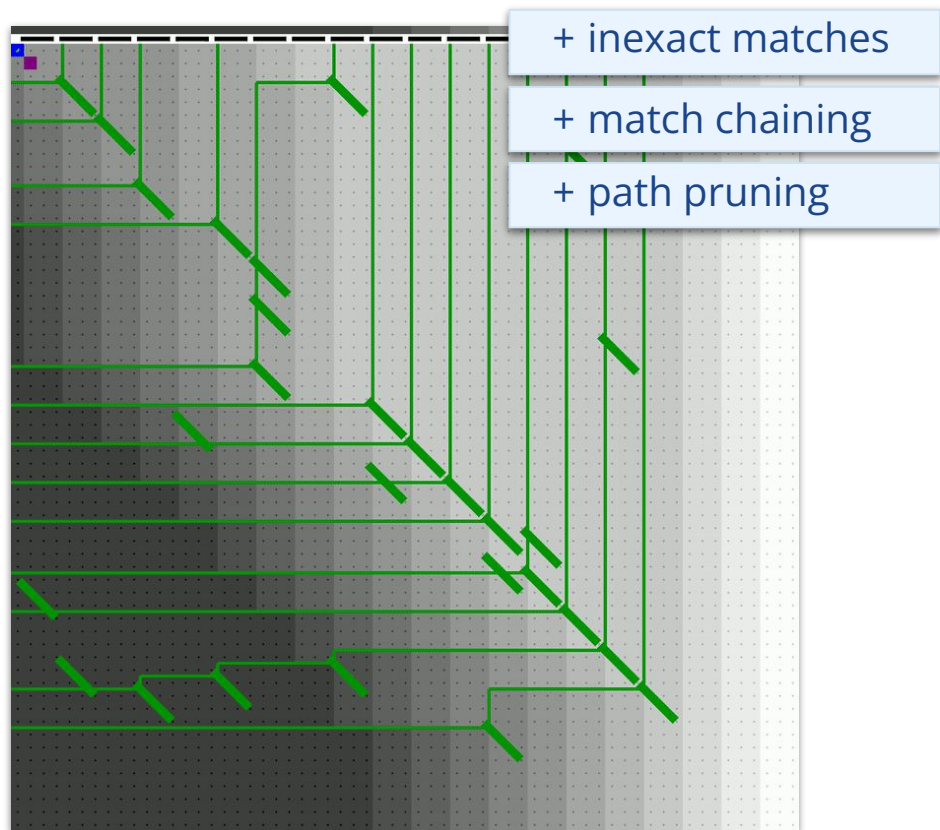
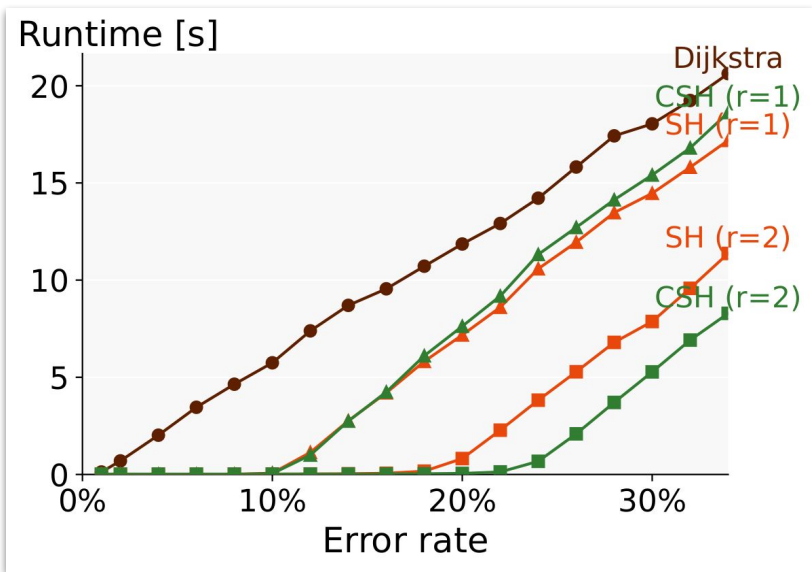
^{*}: For n and e small enough

A* + Diagonal transition



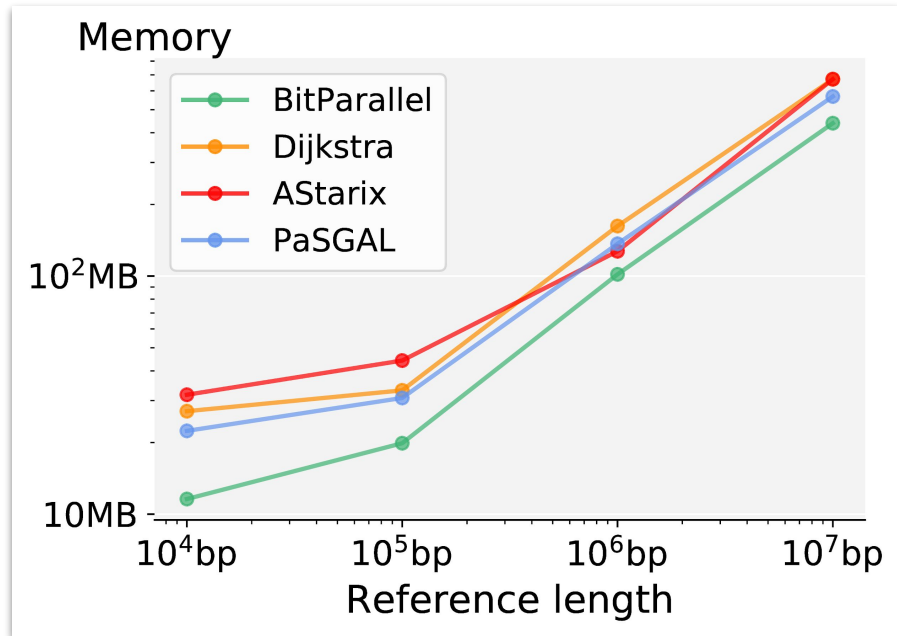
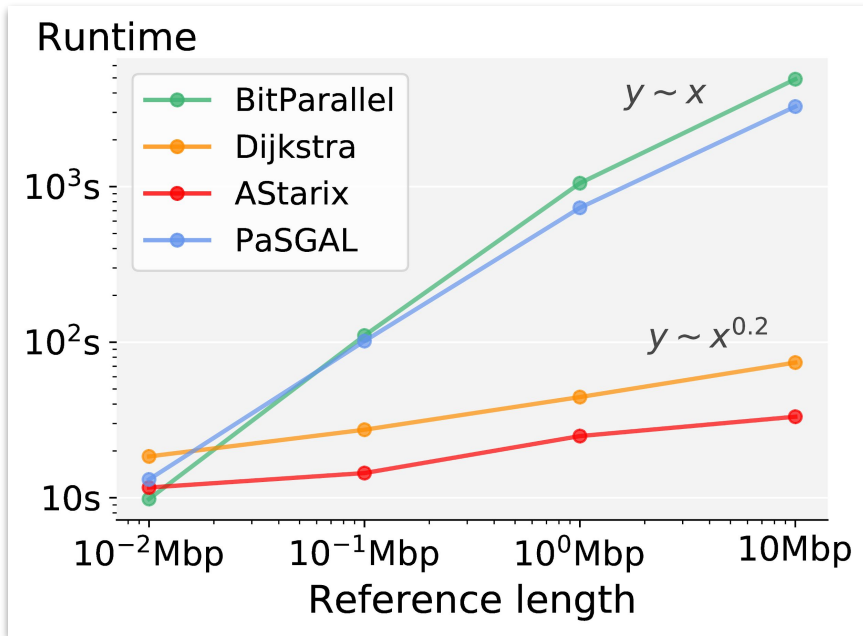
without/with DT
($n=250$, $e=30\%$, $k=6$)

Scaling: error rate (in progress)



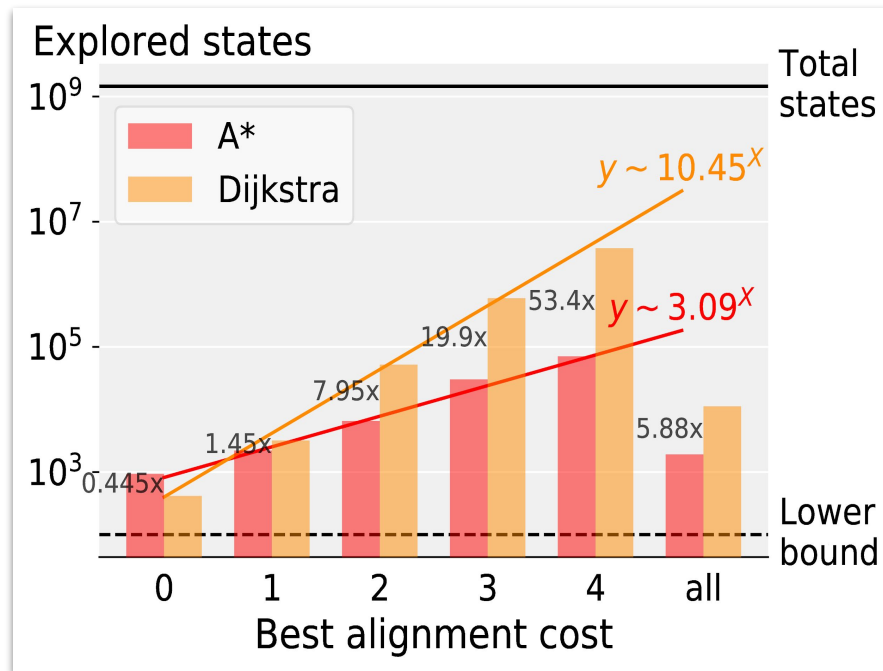
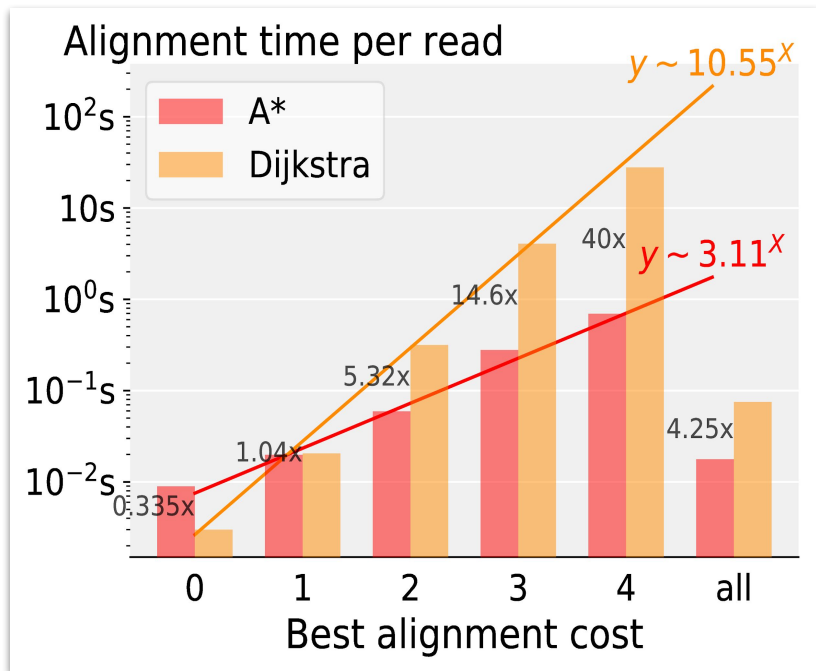
Joint work with **Ragnar Groot Koerkamp**
Visualization by Mykola Akulov
ETH Zurich

Prefix heuristic: Scaling with reference size




The reference length is the prefix of the linear E.coli used as a reference reference

Prefix heuristic: Alignment cost ↗



Dijkstra reuses all codebase, optimizations and parameters for A*, except for the A*-specific.



Abstract. We present a novel A* *seed heuristic* enabling fast and optimal sequence-to-graph alignment, guaranteed to minimize the edit distance of the alignment assuming non-negative edit costs.

We phrase optimal alignment as a shortest path problem and solve it by instantiating the A* algorithm with our novel *seed heuristic*. The key idea of the *seed heuristic* is to extract *seeds* from the read, locate them in the reference, mark preceding reference positions by *crumbs*, and use the crumbs to direct the A* search. We prove admissibility of the *seed heuristic*, thus guaranteeing alignment optimality.

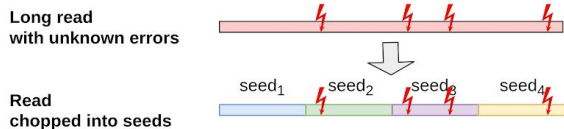
Our implementation extends the free and open source ASTARIX aligner and demonstrates that the *seed heuristic* outperforms all state-of-the-art optimal aligners including GRAPHALIGNER, VARGAS, PASGAL, and the *prefix heuristic* previously employed by ASTARIX. Specifically, we achieve a consistent speedup of $>60\times$ on both short Illumina reads and long HiFi reads (up to 25kbp), on both the *E. coli* linear reference genome (1Mbp) and the MHC variant graph (5Mbp). Our speedup is enabled by the *seed heuristic* consistently skipping $>99.99\%$ of the table cells that optimal aligners based on dynamic programming compute.

AStarix 2.0: Scaling optimal sequence-to-graph alignment to long reads

presented by Pesho Ivanov
ETH Zurich

ETH zürich

1 Chop the read into seeds



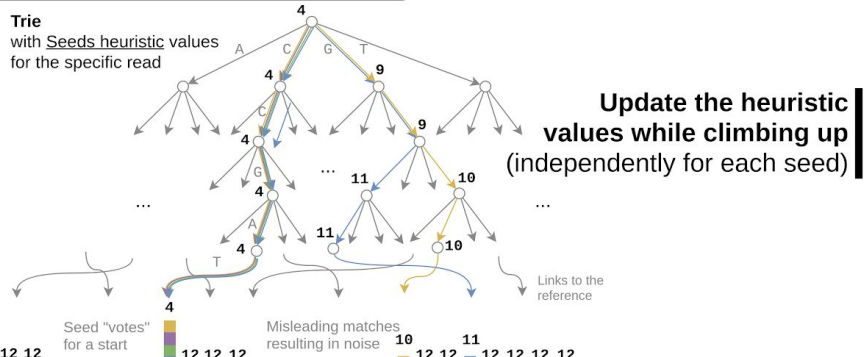
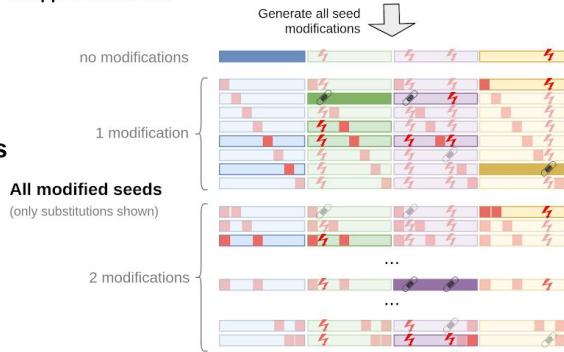
Parameters

Match cost: 0
 Mismatch cost: 1
 Indel cost: $+\infty$ (no indels)

Seeds: 4
 Max. #corrections: 2
 \Rightarrow Max. penalty: 3

Potential = Seeds * Max.Penalty * Max.edit cost
 This heuristic can "pay" for ≤ 12 errors

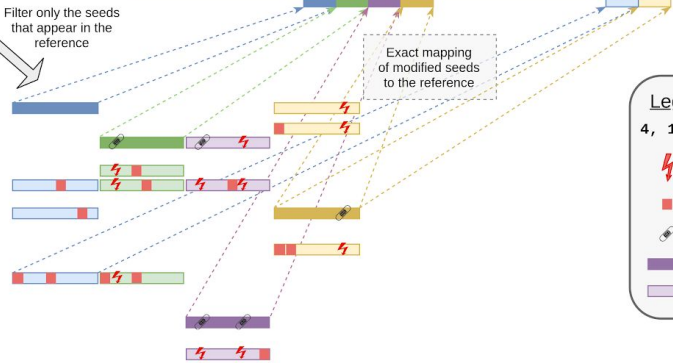
2 Generate all seeds at distance ≤ 2



4

Throw away most seeds as they do not appear in the reference

`/dev/null`
 Default penalty ≥ 3



3 Map the modified seeds exactly

exact match \Rightarrow no penalty

1 modification \Rightarrow penalty ≥ 1

2 modifications \Rightarrow penalty ≥ 2

Modified seeds that appear in the reference

Legend

4, 12 Heuristic value from 0 to the potential

⚡ Read errors not known in advance (only substitutions shown for simplicity)

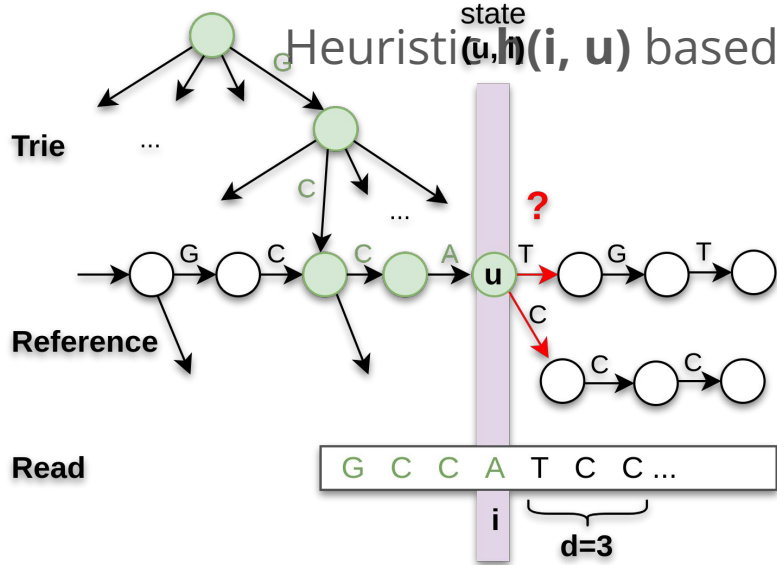
❌ Incorrect correction in a seed

🔧 Correct correction in a seed

🟡 Corrected seed seen in the reference

🟠 Corrected seed not seen in the reference

AStarix: A* prefix heuristic

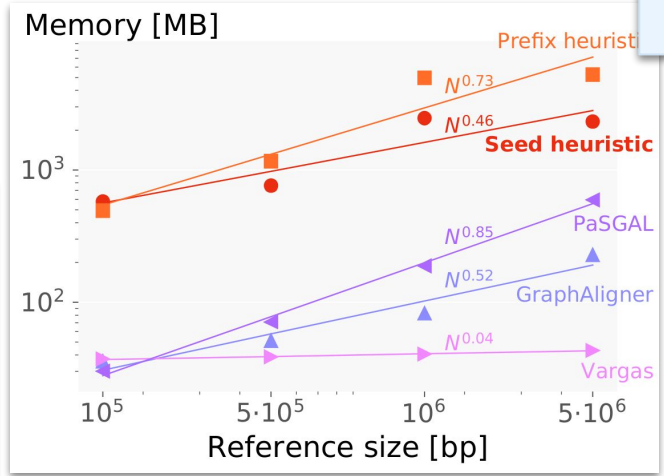
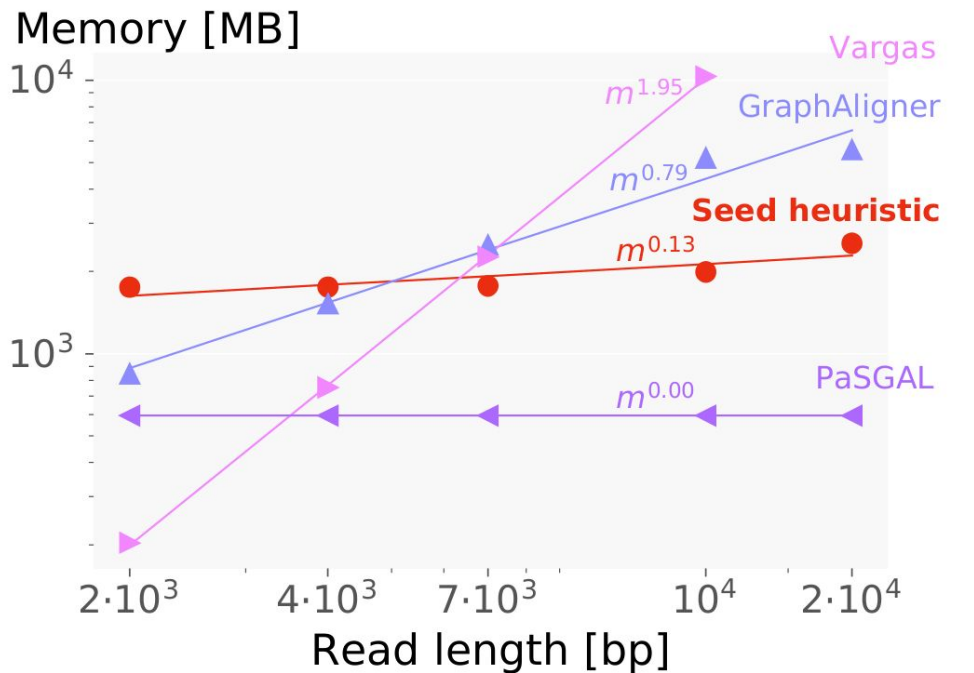


Algorithm 3 Recursive alignment used by Heuristic in Algorithm 1.

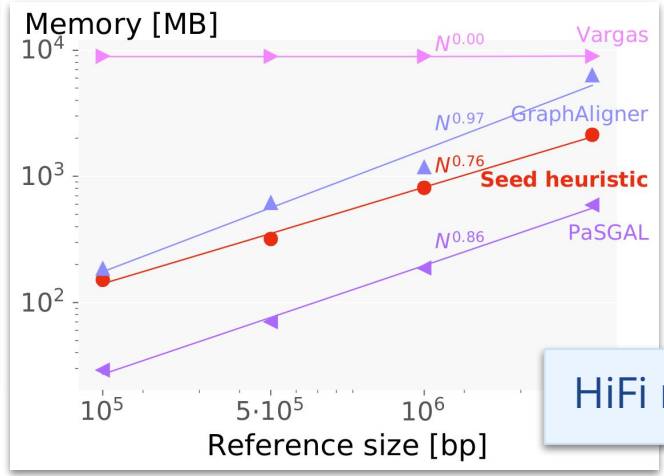
```

1: function RECURSIVEALIGN( $u, s, curr, best$ ) ▷ Return value is  $\leq best$ 
2:   if  $curr \geq best$  then
3:     return  $best$  ▷ Branch and bound: bounding
4:   if  $s = \epsilon$  then ▷ Reached a target
5:     return  $curr$ 
6:   for all  $(u, v, l, w) \in E_e$  where  $l \in \{s[0], \epsilon\}$  do
7:      $suff = s[1 :]$  if  $l \neq \epsilon$  else  $s$ 
8:      $best = RECURSIVEALIGN(u, suff, curr + w, best)$ 
9:   return  $best$ 
  
```

Memory



Illumina reads



HiFi reads

Results

Tool	Illumina		HiFi		
	<i>E. coli</i>	MHC	<i>E. coli</i>	MHC	
<i>Seeds heuristic</i> (this work)	0.019	0.041	0.001	0.002	s/kbp
	2.4	2.6	2.4	1.7	GB
	99.9996	99.9981	99.9989	99.9984	% skipped states
<i>Prefix heuristic</i>	269x	180x	n/a	n/a	x slowdown
	7.7	9.6	>20	>20	
	99.9501	99.9501	n/a	n/a	
GRAPHALIGNER	424x	212x	118x	64x	
	0.2	0.2	3.6	3.4	
VARGAS	133x	67x	1 413x	705x	
	<0.1	<0.1	7.3	7.3	
PASGAL	263x	130x	1 367x	736x	
	0.6	0.6	0.6	0.6	

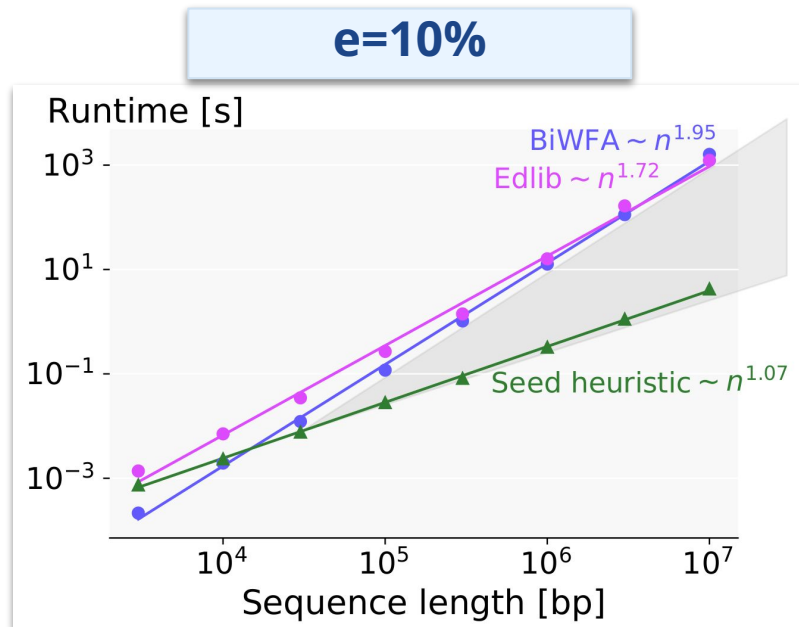
References: *E. coli* – linear 4.6M, Major Histocompatibility Complex (MHC) – 5.3M

Simulated queries: 200bp Illumina with $\Delta=(0,1,5,5)$; HiFi: 5–25kbp, $e=0.3\%$, $\Delta=(0,1,1,1)$

Prefix heuristic parameters: length cap $d=5$, cost cap $c=5$, trie depth $D=\log(N)$

Seed heuristic parameters: $D=14 \approx \log_4 N$; $k=25$ for Illumina, $k=150$ for HiFi reads

Near-linear scaling



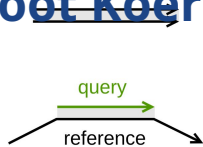
Seed heuristic for global alignment (unpublished)

Ragnar Groot Koerkamp and Pesho Ivan

ETH Zurich

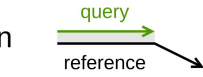
Pairwise alignment types

Global
Semi-global



Global-local (glocal),
mapping, infix, SHW

Global-extension



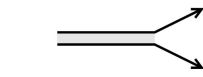
Global-prefix

Overlap



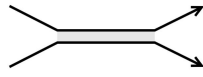
Prefix-suffix,
suffix-prefix, dovetail

Extension

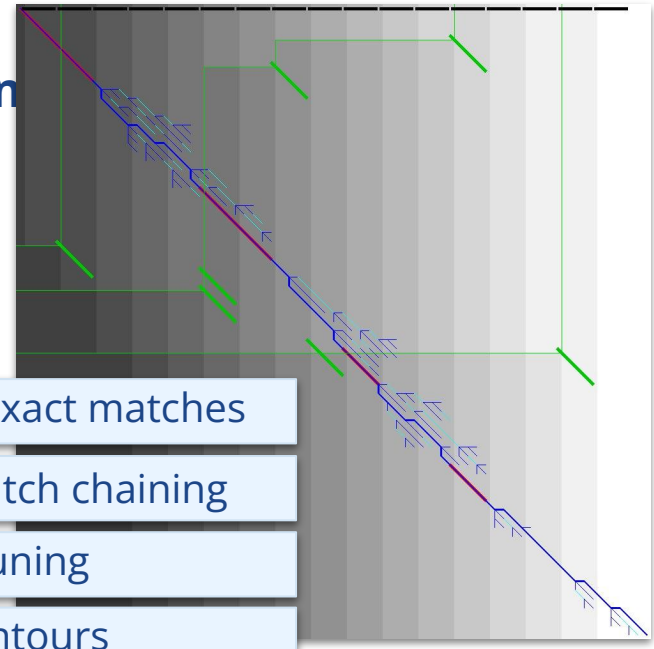


Prefix-prefix

Local

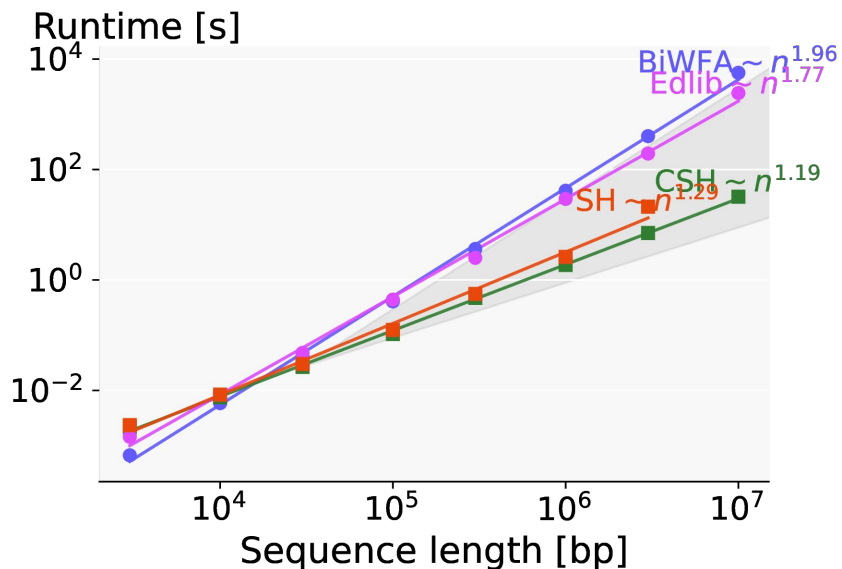


SW

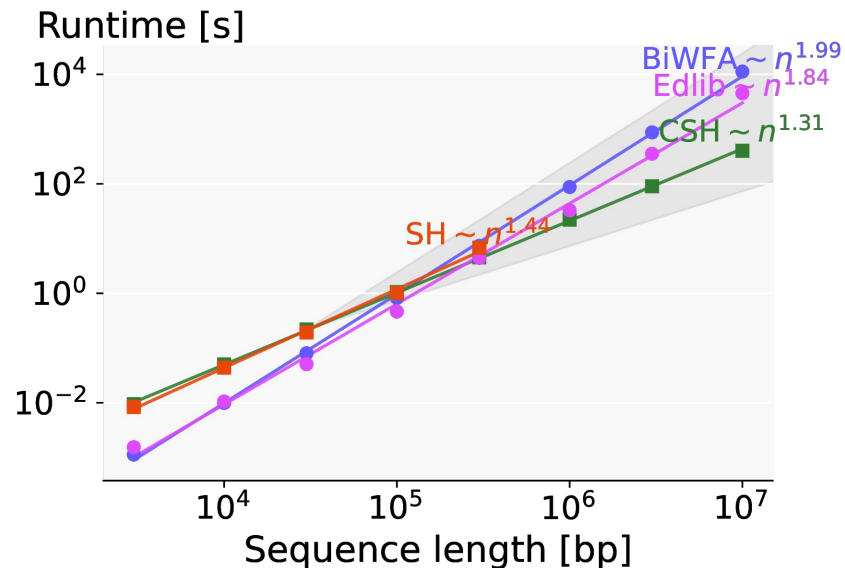


- + inexact matches
- + match chaining
- + pruning
- + contours

Comparison for high error rate

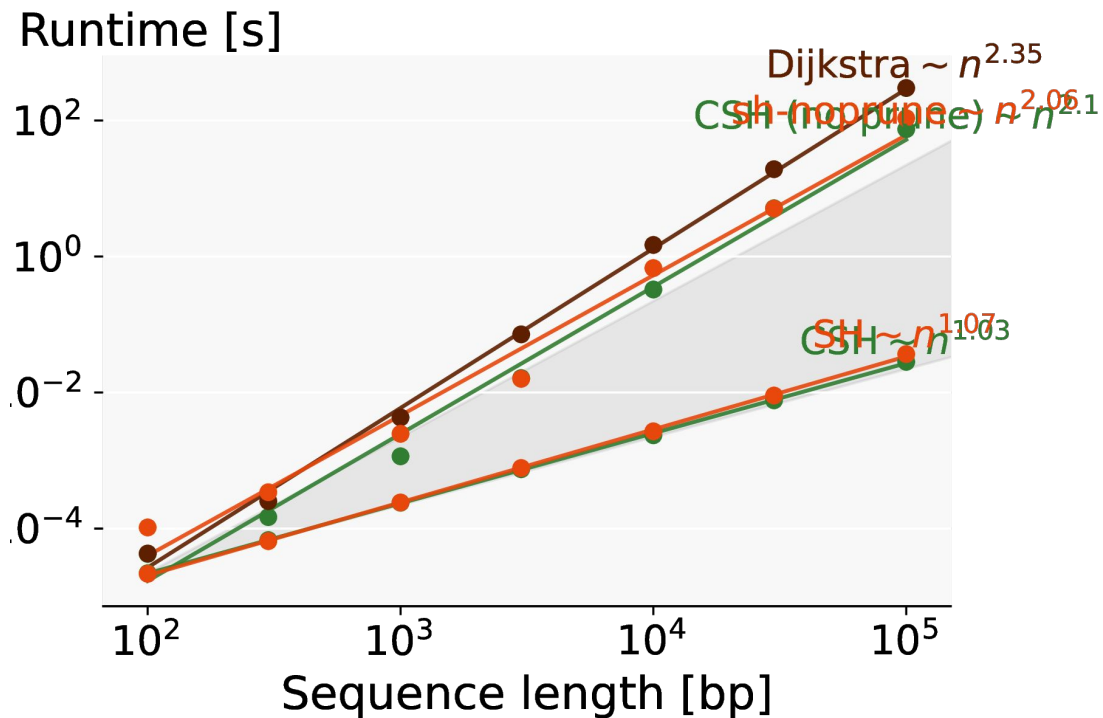


Error rate: 10%
Seed length: 15
Inexact matches



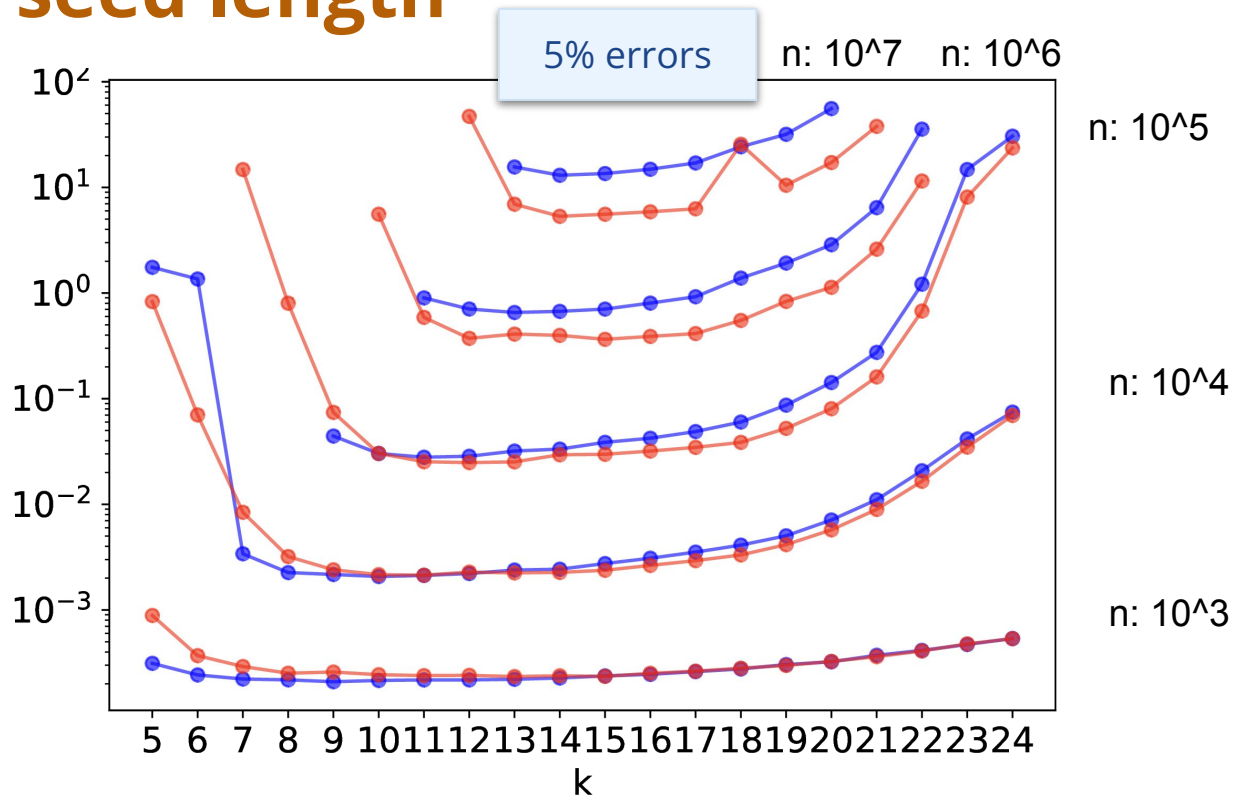
Error rate: 15%
Seed length: 15
Inexact matches

Effect of pruning



Error rate: 5%, seed length: 15

Choosing seed length



$$\log(n) < k < 1/e$$