

Scalable String Similarity Search/Join with Approximate Seeds and Multiple Backtracking

Enrico Siragusa, David Weese, Knut Reinert
FU Berlin, Algorithmic Bioinformatics

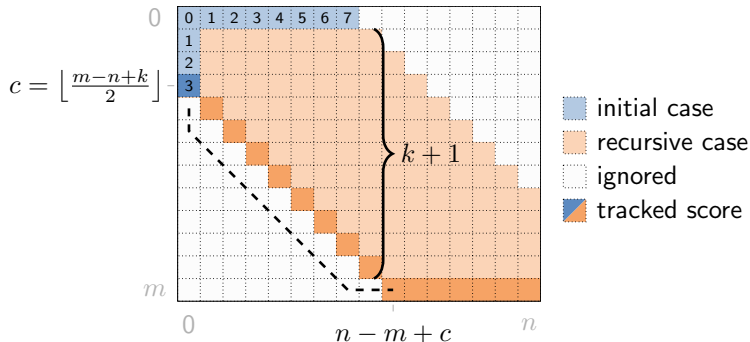
String Similarity Search/Join Competition, 22 March 2013

- ▶ This work is a variation of our new DNA read mapper¹
- ▶ We follow the seminal work of Navarro and Baeza-Yates²
- ▶ We combine multiple methods
 - online Banded Myers bit-vector algorithm
 - indexed Multiple backtracking
 - filtering Approximate seeds
- ▶ We preprocess both database and query strings
- ▶ Our approach solves search and join

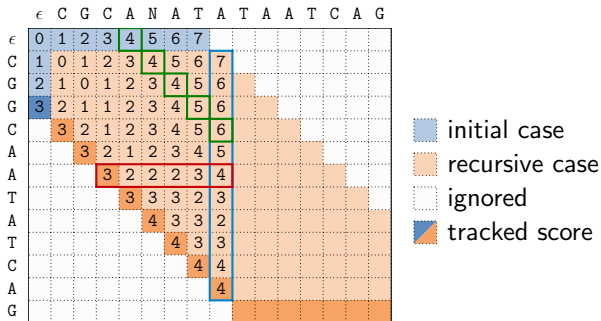
¹Enrico Siragusa, David Weese, and Knut Reinert. “Fast and accurate read mapping with approximate seeds and multiple backtracking”. In: **Nucleic Acids Res.** (2013).

²Gonzalo Navarro and Rircardo A. Baeza-Yates. “A hybrid indexing method for approximate string matching”. In: **Journal of Discrete Algorithms** 1.1 (2000), pp. 205–239.

- ▶ Check the edit distance between each query and any database string: **global k-differences problem**
- ▶ A band of $k + 1$ cells is sufficient

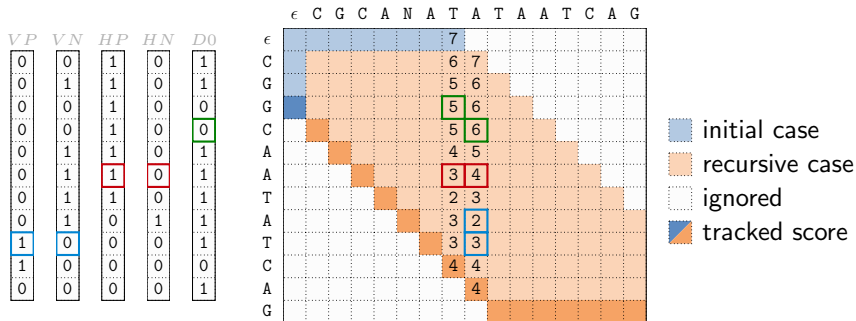


- ▶ The DP matrix has **diagonal** & **adjacency** properties



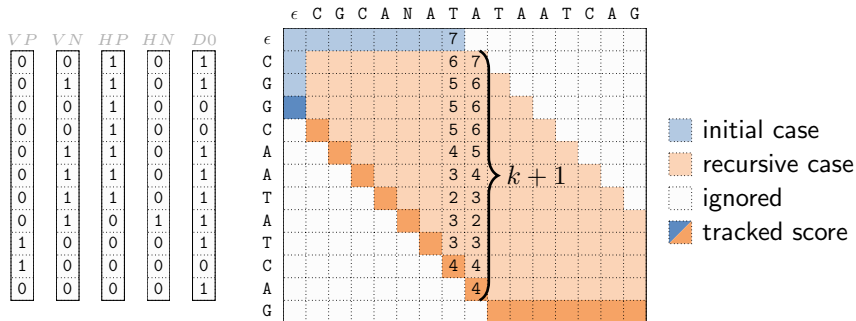
³Gene Myers. "A fast bit-vector algorithm for approximate string matching based on dynamic programming". In: **J. ACM** 46.3 (1999), pp. 395–415.

- ▶ The DP matrix has **diagonal** & **adjacency** properties
- ▶ Myers algorithm³ encodes each column as **bit-vectors**



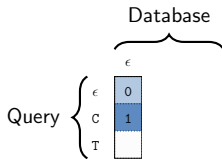
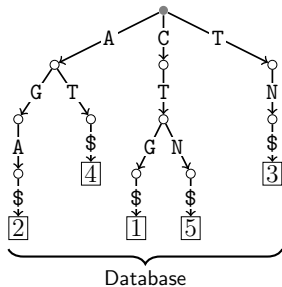
³Gene Myers. "A fast bit-vector algorithm for approximate string matching based on dynamic programming". In: **J. ACM** 46.3 (1999), pp. 395–415.

- ▶ The DP matrix has **diagonal** & **adjacency** properties
- ▶ Myers algorithm³ encodes each column as **bit-vectors**
- ▶ Our banded version computes a column in time $\mathcal{O}(k/w)$

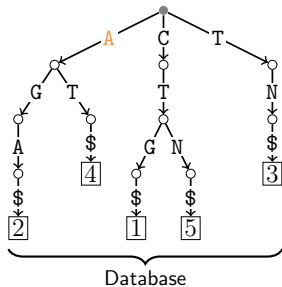


³Gene Myers. “A fast bit-vector algorithm for approximate string matching based on dynamic programming”. In: **J. ACM** 46.3 (1999), pp. 395–415.

- ▶ Index database strings using a **radix tree**
- ▶ Perform a **top-down traversal** on the radix tree



- ▶ Index database strings using a **radix tree**
- ▶ Perform a **top-down traversal** on the radix tree
- ▶ Compute the distance between query and edge labels

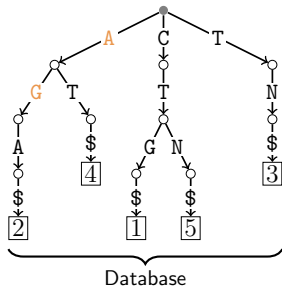


Database

	ϵ	A
ϵ	0	1
C	1	1
T		1

Query

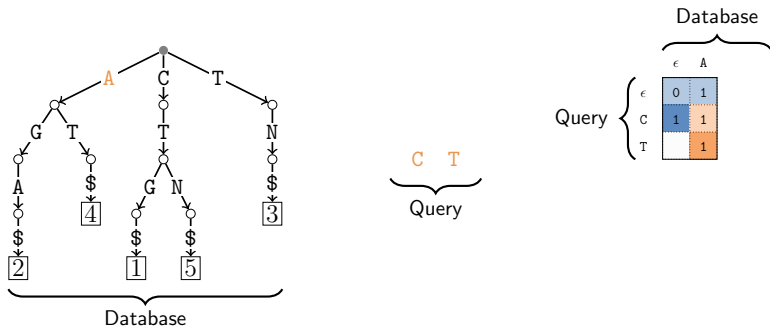
- ▶ Index database strings using a **radix tree**
- ▶ Perform a **top-down traversal** on the radix tree
- ▶ Compute the distance between query and edge labels
- ▶ **Cut** a branch when the minimum distance exceeds k



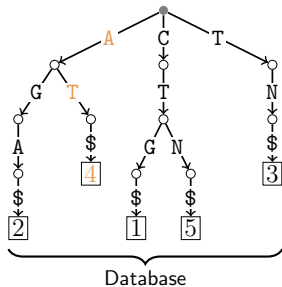
Database

	ϵ	A	G
Query $\left\{ \begin{array}{l} \epsilon \\ C \\ T \end{array} \right.$	0	1	
1	1	2	
	1	2	

- ▶ Index database strings using a **radix tree**
- ▶ Perform a **top-down traversal** on the radix tree
- ▶ Compute the distance between query and edge labels
- ▶ **Cut** a branch when the minimum distance exceeds k



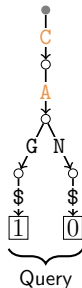
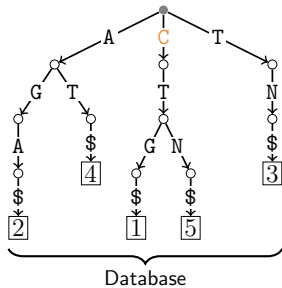
- ▶ Index database strings using a **radix tree**
- ▶ Perform a **top-down traversal** on the radix tree
- ▶ Compute the distance between query and edge labels
- ▶ **Cut** a branch when the minimum distance exceeds k
- ▶ Backtracking takes time **exponential** in k



Database

	ϵ	A	T
Query $\left\{ \begin{array}{l} \epsilon \\ C \\ T \end{array} \right.$	0	1	
1	1	2	
	1	1	

- ▶ Any two queries sharing a common prefix also share part of their backtracking
- ▶ Index the queries using another **radix tree**
- ▶ Backtrack the query tree in the database tree

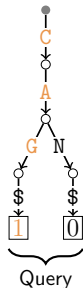
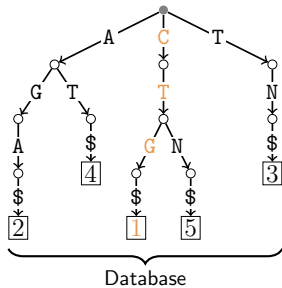


Database

	ϵ	C
ϵ	0	1
C	1	0
A		1

Query

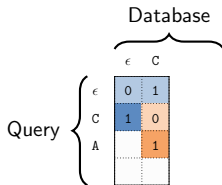
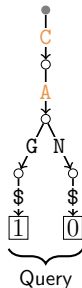
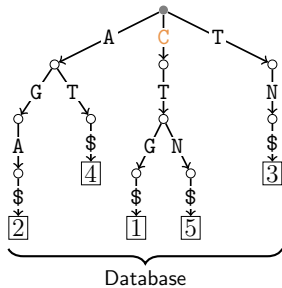
- ▶ Any two queries sharing a common prefix also share part of their backtracking
- ▶ Index the queries using another **radix tree**
- ▶ Backtrack the query tree in the database tree



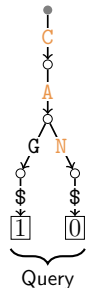
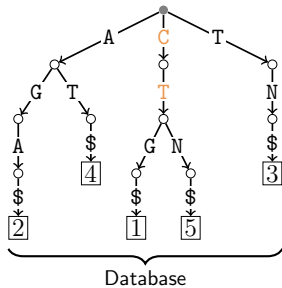
Database

	ϵ	C	T	G
Query ϵ	0	1		
Query C	1	0	1	
Query A		1	1	2
Query G			2	1

- ▶ Any two queries sharing a common prefix also share part of their backtracking
- ▶ Index the queries using another **radix tree**
- ▶ Backtrack the query tree in the database tree



- ▶ Any two queries sharing a common prefix also share part of their backtracking
- ▶ Index the queries using another **radix tree**
- ▶ Backtrack the query tree in the database tree

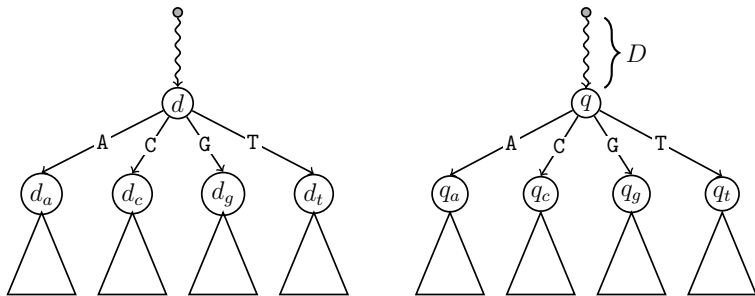


Database

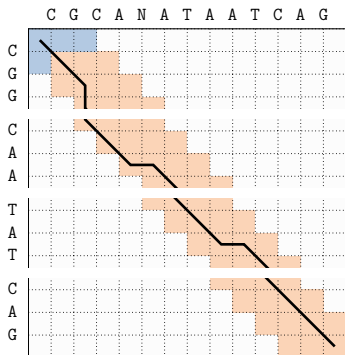
	ϵ	C	T
Query ϵ	0	1	
Query C	1	0	1
Query A		1	1
Query N			2

- ▶ Parallelization of multiple backtracking is non-trivial
- ▶ We collect pairs of subtrees at fixed depth D
- ▶ Each pair can be processed independently

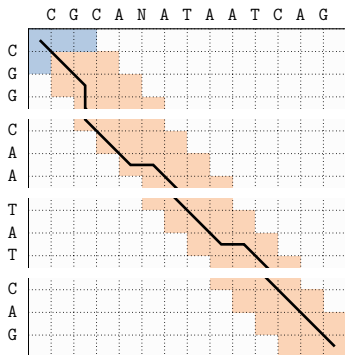
$$\mathcal{Q} \leftarrow \mathcal{Q} \cup \{(d_a, q_a), (d_a, q_c), (d_a, q_g), \dots, (d_t, q_t)\}$$



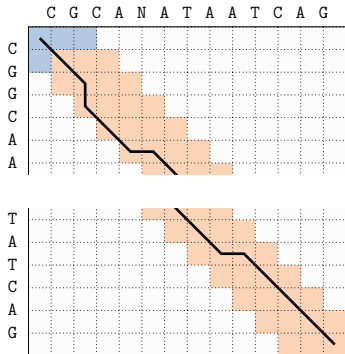
- ▶ Partition query Q in $k + 1$ seeds
 - ▶ One seed occurs without errors
1. Index database strings with a **generalized suffix tree**
 2. Search all seeds
 3. Verify seeds inside a k -band



- ▶ Partition query Q in $k + 1$ seeds
 - ▶ One seed occurs without errors
1. Index database strings with a **generalized suffix tree**
 2. Search all seeds
 3. Verify seeds inside a k -band
- ▶ Seed length $L = \frac{|Q|}{k + 1}$

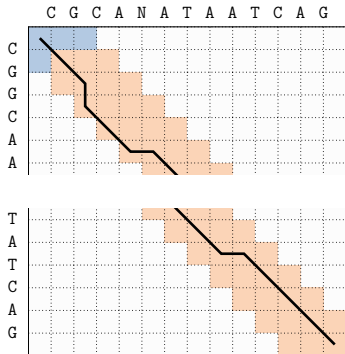


- ▶ Partition query Q in s seeds
- ▶ One seed is within $\lfloor k/s \rfloor$ errors



- ▶ Partition query Q in s seeds
- ▶ One seed is within $\lfloor k/s \rfloor$ errors

- ▶ Seed length $L = |Q|/s$
- ▶ The seed length is a parameter



- ▶ Implemented in C++ using SeqAn and OpenMP
- ▶ Indices implemented as generalized suffix arrays
- ▶ Index construction combines bucket sort and quicksort
- ▶ Backtracking implemented with banded DP






David Weese & Knut Reinert

The SeqAn team



**International
Max Planck Research School**
for Computational Biology
and Scientific Computing

-  Gene Myers. “A fast bit-vector algorithm for approximate string matching based on dynamic programming”. In: **J. ACM** 46.3 (1999), pp. 395–415.
-  Gonzalo Navarro and Rircardo A. Baeza-Yates. “A hybrid indexing method for approximate string matching”. In: **Journal of Discrete Algorithms** 1.1 (2000), pp. 205–239.
-  Enrico Siragusa, David Weese, and Knut Reinert. “Fast and accurate read mapping with approximate seeds and multiple backtracking”. In: **Nucleic Acids Res.** (2013).