# A NEW ALGORITHM FOR THE ALIGNMENT OF LONG GENOMIC SEQUENCES

INTERNSHIP SUPERVISED BY DOMINIQUE LAVENIER, GENSCALE TEAM

Igor MARTAYAN
September 20, 2021

DNA sequencing is more and more accessible

ATCGGGCAAT  TAAAAGGATC
TGAAGCGAAG  ACACCGTACC
AGACGTAGCG  AGCCCTATTT
ATCGGAGCAA  TAAAAGGATC
CGAAGCGAGA  GACCACCGTA
CAGGACGTAG  GGAGCCCTAT

DNA sequencing is more and more accessible
but still presents many errors

```
ATCGGGCAAT   TAAAAGGATC
TGAAGCGAAG   ACACCGTACC
AGACGTAGCG   AGCCCTATTT
ATCGGAGCAA   TAAAAGGATC
CGAAGCGAGA   GACCACCGTA
CAGGACGTAG   GGAGCCCTAT
```

How do we correct these errors?

3 types of errors
- insertion: `GCA -> GCTA`
- deletion: `GCA -> GA`
- substitution: `GCA -> GTA`

3 types of errors

- insertion: `GCA -> GCTA`
- deletion: `GCA -> GA`
- substitution: `GCA -> GTA`

several reads from the same sequence

```
TTGAC_TCAAGGGCCA_TCATG
T _ACAT_A _GGCTAATTATG
T GACATTACGGGCCAGTAATG
TCGACA_CA GGGTAAA_AATT
---------------------
T GACATCA GGGCCAATAATG
```

Computing an alignement score

$+\alpha$ by match

$-\beta$ by substitution

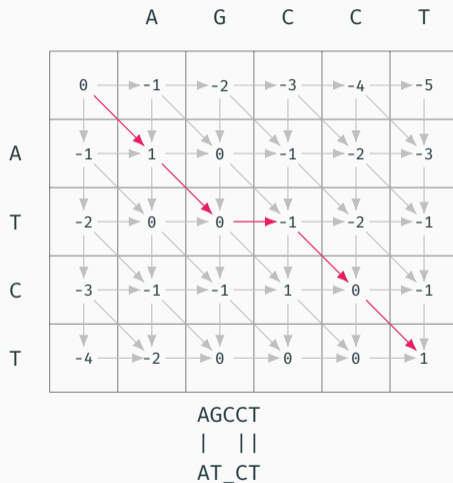$-\delta$ by insertion/deletion

Computing an alignement score

$+\alpha$ by match

$-\beta$ by substitution

$-\delta$ by insertion/deletion

Needleman—Wunsch algorithm
Compute score matrix with dynamic
programming, $\mathcal{O}(nm)$ complexity

Computing an alignement score

$+\alpha$ by match

$-\beta$ by substitution

$-\delta$ by insertion/deletion

Needleman—Wunsch algorithm

Compute score matrix with dynamic programming, $\mathcal{O}(nm)$ complexity



```
AGCCT
|  ||
AT_CT
```
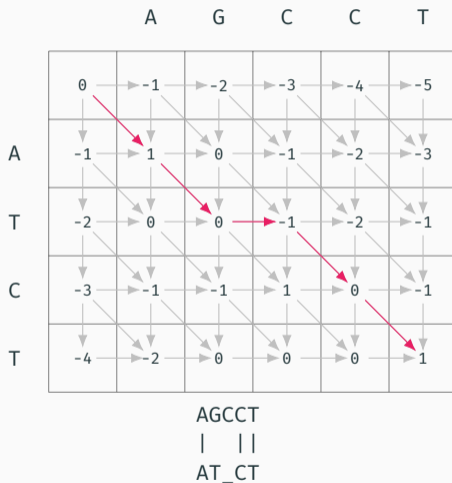
Computing an alignement score

$+\alpha$ by match

$-\beta$ by substitution

$-\delta$ by insertion/deletion

**Needleman—Wunsch algorithm**

Compute score matrix with dynamic
programming, $\mathcal{O}(nm)$ **complexity**



```
AGCCT
|  ||
AT_CT
```

HOW CAN WE GO FASTER?

Google Scholar

accelerating sequence alignment
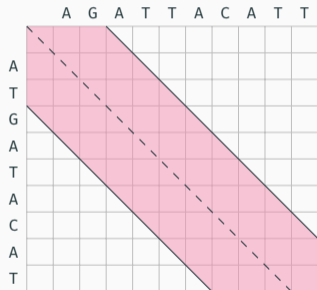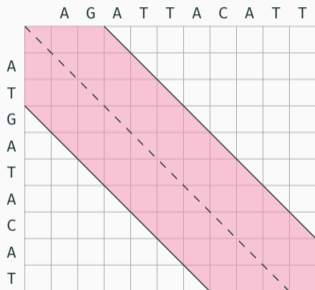
Articles

About 180,000 results (0.03 sec)

Restricting the score matrix:

Restricting the score matrix:



Vectorizing operations:
SSE / AVX instruction set

Restricting the score matrix:



Vectorizing operations:
SSE / AVX instruction set

Hardware acceleration:
using GPU / FPGA

# My contribution

Constraints we have to deal with:

- very long sequences (20000 bases)
- errors can be agglomerated

Constraints we have to deal with:

- very long sequences (20000 bases)
- errors can be agglomerated

*but* sequences are very similar

Constraints we have to deal with:

- very long sequences (20000 bases)
- errors can be agglomerated

*but* sequences are very similar

Desired properties:

- subquadratic time complexity
- low memory footprint

Constraints we have to deal with:

- very long sequences (20000 bases)
- errors can be agglomerated

*but* sequences are very similar

Desired properties:

- subquadratic time complexity
- low memory footprint
- easy parallelization

Constraints we have to deal with:

- very long sequences (20000 bases)
- errors can be agglomerated

*but* sequences are very similar

Desired properties:

- subquadratic time complexity
- low memory footprint
- easy parallelization
- reduced instruction set

ATCGGGCAATTAAAAGGATCTGAAGCGAAGACACCGTACCAGACGTAGCGAGCCCTATTT

ATCGGAGCAATAAAAGGATCCGAGCGAGAGACCACCGTACAGGACTTAGGGAGCCCTATTT

take advantage of the similarity between sequences

ATCGGGCAATTAAAAGGATCTGAAGCGAAGACACCGTACCAGACGTAGCGAGCCCTATTT

ATCGGAGCAATAAAAGGATCCGAGCGAGAGACCACCGTACAGGACTTAGGGAGCCCTATTT

1. mark the words of size $k$ common to both sequences

ATCGG  GCAAT  **TAAAAGG**  ATCTGAAGCGA  AGA    **CACCGTA**  CCA  GACGTAGC  **GAGCCCT**  ATTT

ATCGGAGCAA    **TAAAAGG**  ATCCGA  GCGAGAGAC  **CACCGTA**    CAGGACTTAGG  **GAGCCCT**  ATTT

1. mark the words of size $k$ common to both sequences

ATCGG GCAAT **TAAAAGG** ATCTGAAGCGA AGA    **CACCGTA** CCA GACGTAGC **GAGCCCT** ATTT

ATCGGAGCAA    **TAAAAGG** ATCCGA GCGAGAGAC **CACCGTA**    CAGGACTTAGG **GAGCCCT** ATTT

1. mark the words of size $k$ common to both sequences
2. align sequences between the anchor points

```
ATCGG_GCAAT TAAAAGG ATCTGAAGCGA_AGA_ CACCGTA CCA_GACGTAGC GAGCCCT ATTT
||||| ||||            ||| || |||| |||           || ||| |||           ||||
ATCGGAGCAA_ TAAAAGG ATCCGA_GCGAGAGAC CACCGTA _CAGGACTTAGG GAGCCCT ATTT
```

1. mark the words of size $k$ common to both sequences
2. align sequences between the anchor points

```
ATCGG_GCAAT TAAAAGG ATCTGAAGCGA_AGA_ CACCGTA CCA_GACGTAGC GAGCCCT ATTT
||||| ||||           ||| || |||| |||            || ||| |||            ||||
ATCGGAGCAA_ TAAAAGG ATCCGA_GCGAGAGAC CACCGTA _CAGGACTTAGG GAGCCCT ATTT
```

1. mark the words of size $k$ common to both sequences
2. align sequences between the anchor points
3. merge the results

```
ATCGG_GCAATTAAAAGGATCTGAAGCGA_AGA_CACCGTACCA_GACGTAGCGAGCCCTATTT
||||| |||| ||||||||| || |||| ||| ||||||| || ||| ||| ||||||||||||
ATCGGAGCAA_TAAAAGGATCCGA_GCGAGAGACCACCGTA_CAGGACTTAGGGAGCCCTATTT
```

1. mark the words of size $k$ common to both sequences
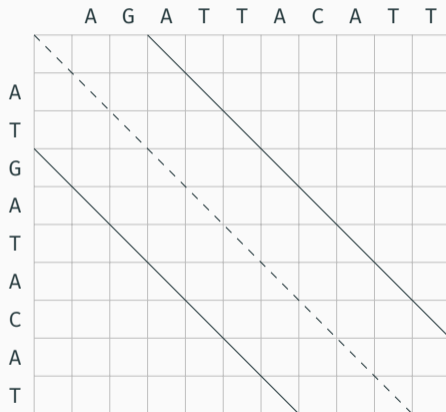2. align sequences between the anchor points
3. merge the results

Apply existing optimizations:

|   |   | A | G | A | T | T | A | C | A | T | T |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |

Apply existing optimizations:

- restricting the matrix

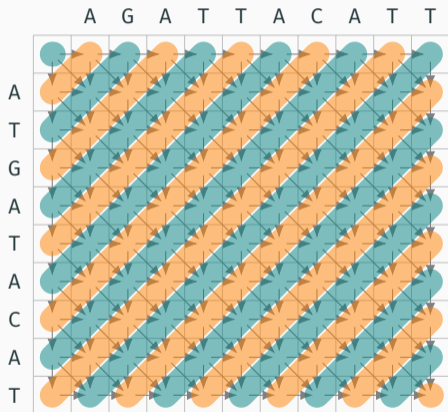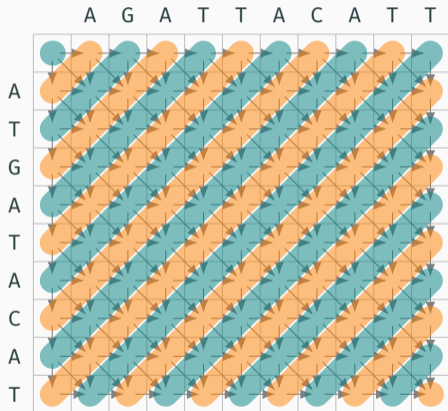Apply existing optimizations:

- restricting the matrix
- vectorizing operations

$$S_{i,j} = \max \left\{ S_{i-1,j} - \delta, S_{i,j-1} - \delta, S_{i-1,j-1} + m_{i,j} \right\}$$

with $m_{i,j} = \begin{cases} \alpha \text{ if } u_i = v_j \\ -\beta \text{ otherwise} \end{cases}$

$$S_{i,j} = \max \left\{ S_{i-1,j} - \delta, S_{i,j-1} - \delta, S_{i-1,j-1} + m_{i,j} \right\}$$

with $m_{i,j} = \begin{cases} \alpha \text{ if } u_i = v_j \\ -\beta \text{ otherwise} \end{cases}$
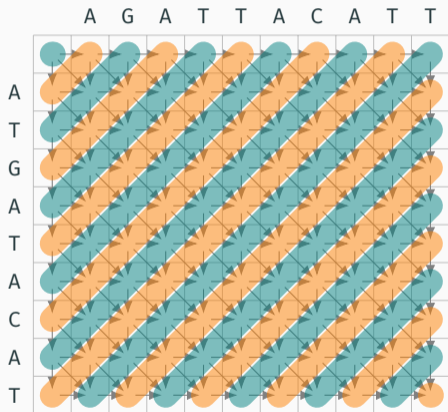
1. sum (match, substitution, indel)

$$S_{i,j} = \max \left\{ S_{i-1,j} - \delta, S_{i,j-1} - \delta, S_{i-1,j-1} + m_{i,j} \right\}$$

with $m_{i,j} = \begin{cases} \alpha \text{ if } u_i = v_j \\ -\beta \text{ otherwise} \end{cases}$

1. `sum` (match, substitution, indel)
2. `compare`

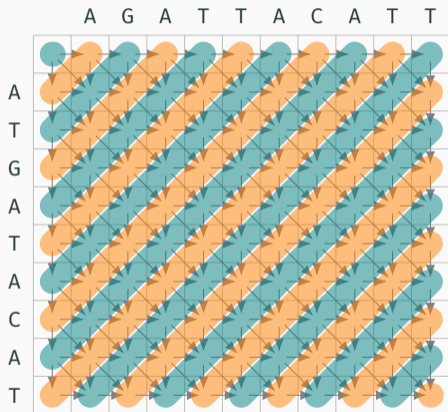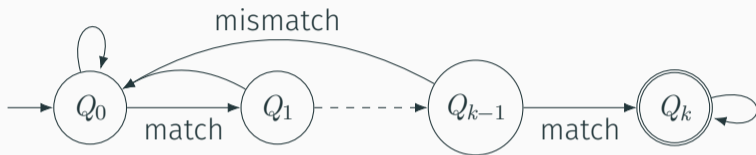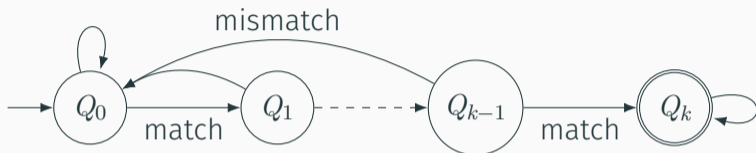$$S_{i,j} = \max\left\{S_{i-1,j} - \delta, S_{i,j-1} - \delta, S_{i-1,j-1} + m_{i,j}\right\}$$

$$\text{with } m_{i,j} = \begin{cases} \alpha \text{ if } u_i = v_j \\ -\beta \text{ otherwise} \end{cases}$$
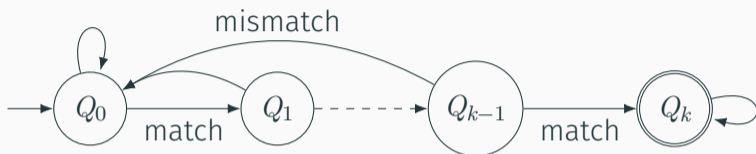
1. `sum` (match, substitution, indel)
2. `compare`
3. `blend` (keep the maximum)

$$X_n = \begin{bmatrix} p & \cdots & p & 0 \\ 1-p & & (0) & \vdots \\ & \ddots & & 0 \\ (0) & & 1-p & 1 \end{bmatrix}^n \times \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$X_n = \begin{bmatrix} p & \cdots & p & 0 \\ 1-p & & (0) & \vdots \\ & \ddots & & 0 \\ (0) & & 1-p & 1 \end{bmatrix}^n \times \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$
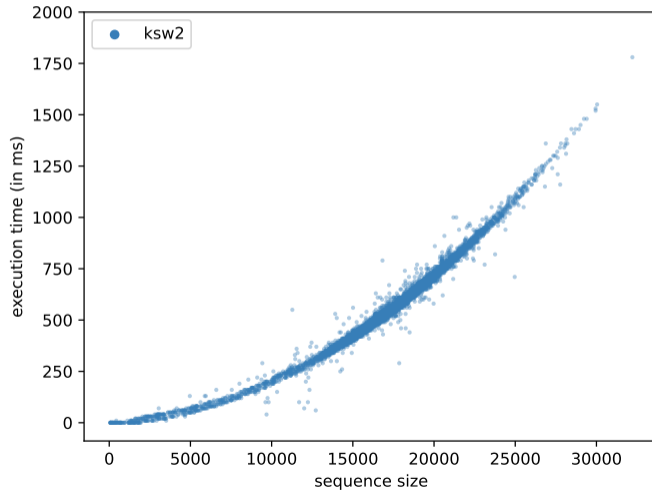
with $k = 16$ and $p = 0.07$,

$\mathbb{E}$ (# steps to reach $Q_k$) $\approx 31$

# RESULTS

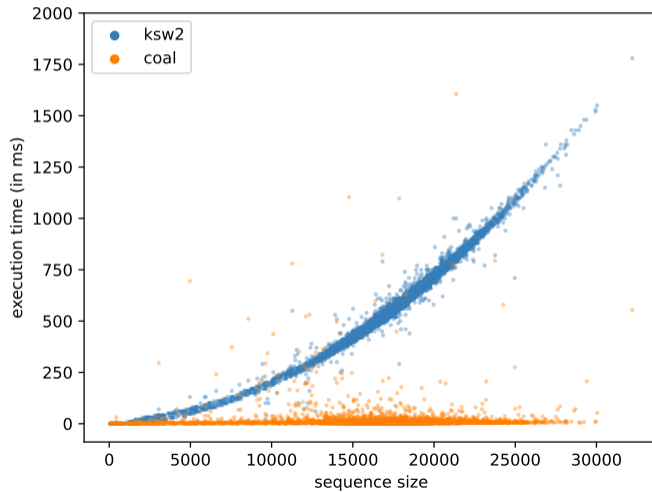|  | sequence size | | | |
|---|---|---|---|---|
|  | 14000–16000 | 16000–18000 | 18000–20000 | 20000–22000 |
|  | average run time (in ms) | | | |
| ksw2 | 430 | 531 | 645 | 776 |
| coal | 10.2 | 8.0 | 7.2 | 10.8 |
| acceleration | × 42 | × 74 | × **90** | × 72 |

|  | sequence size | | | |
|---|---|---|---|---|
|  | 14000–16000 | 16000–18000 | 18000–20000 | 20000–22000 |
|  | average score | | | |
| ksw2 | 18967 | 21955 | 24543 | 26878 |
| coal | 18304 | 21359 | 23931 | 26071 |
| relative gap | 3.5% | 2.7% | 2.5% | 3.0% |

# NEXT STEPS

Improve scores:

- refine the choice of anchor points
- try different alignment methods

Improve scores:

- refine the choice of anchor points
- try different alignment methods

Push performances further:

- parallelize subalignments
- improve hashing methods
- port on other architectures (*processing-in-memory*)

processors directly integrated into the DRAM

processors directly integrated into the DRAM

minimize data access time,
speed up data intensive calculations

many applications in genomics

processors directly integrated into the DRAM



minimize data access time,
speed up data intensive calculations

many applications in genomics

starting in 2022:

- CIFRE PhD
- Inria GenoPIM project

# Conclusion

(+) better performances ($\times$ 90)

(+) low memory footprint

(+) adaptable with different alignment methods

(-) lower scores (-3%)

(+) better performances ($\times$ 90)

(+) low memory footprint

(+) adaptable with different
alignment methods

(-) lower scores (-3%)

https://github.com/imartayan/coal