Project :

# Design of an embedding alignment program by dynamic programming

Mia Legras
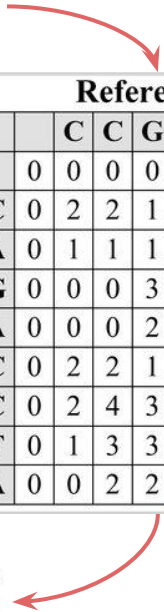
# Aim : Implement an alignment using embeddings

## Alignment

Reference (R):   CCGTACTA
Query (Q):       CAGACCTA



|        |   | C | C | G | T | A | C | T | A |
|--------|---|---|---|---|---|---|---|---|---|
|        | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C      | 0 | 2 | 2 | 1 | 0 | 0 | 2 | 1 | 0 |
| A      | 0 | 1 | 1 | 1 | 0 | 2 | 1 | 1 | 3 |
| G      | 0 | 0 | 0 | 3 | 2 | 1 | 1 | 0 | 2 |
| A      | 0 | 0 | 0 | 2 | 2 | 4 | 3 | 2 | 2 |
| C      | 0 | 2 | 2 | 1 | 1 | 3 | 6 | 5 | 4 |
| C      | 0 | 2 | 4 | 3 | 2 | 2 | 5 | 5 | 4 |
| T      | 0 | 1 | 3 | 3 | 5 | 4 | 4 | 7 | 6 |
| A      | 0 | 0 | 2 | 2 | 4 | 7 | 6 | 6 | 9 |

Reference (R) (top), Query (Q) (side)

C−GTAC−TA
| | || ||
CAG−ACCTA

Wang, Z., Combs, S.A., Brand, R. *et al.* LM-GVP: an extensible sequence and structure informed deep learning framework for protein property prediction
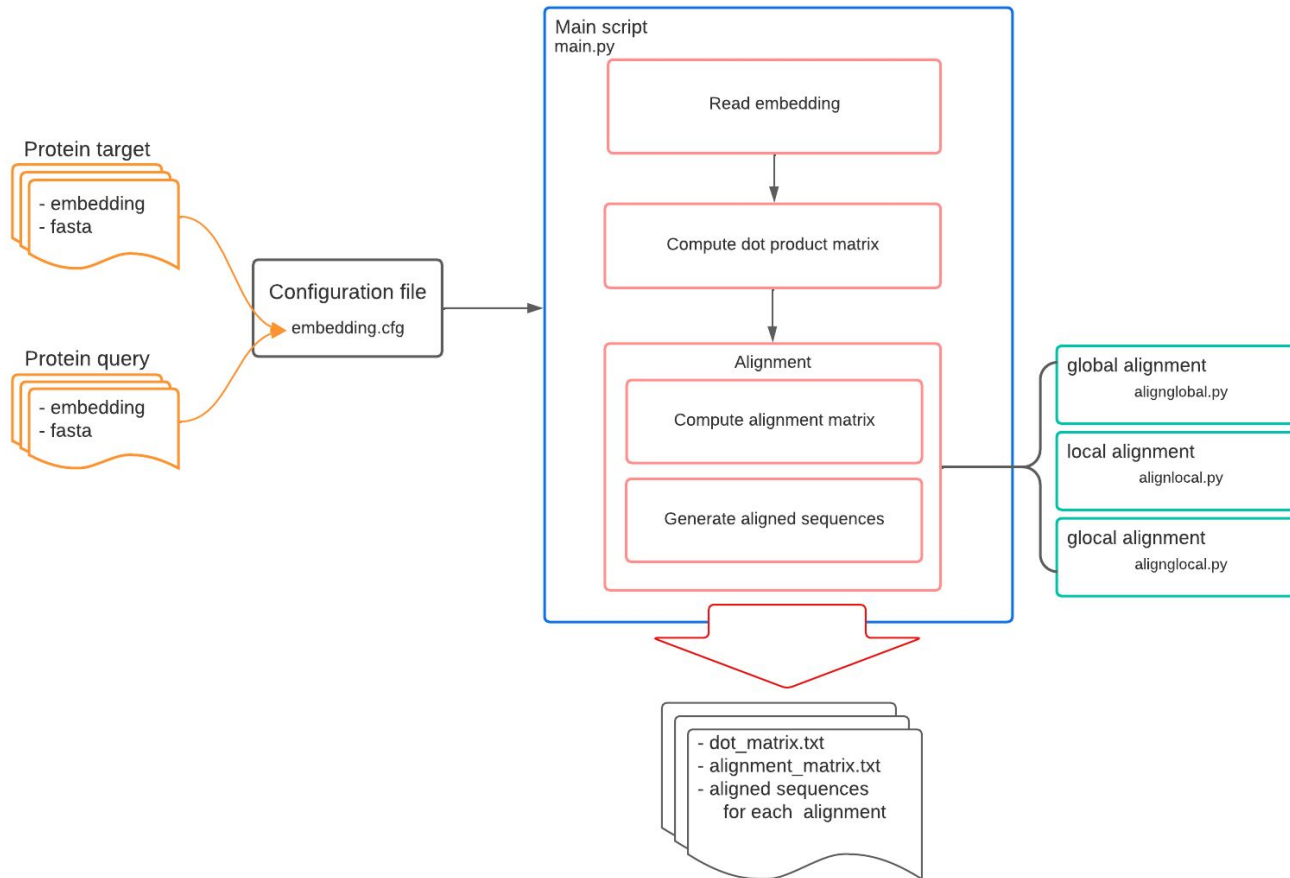
## Embedding

Vectorial matrix

Encoding

Input sequence

M
K
G
E
E
L
...
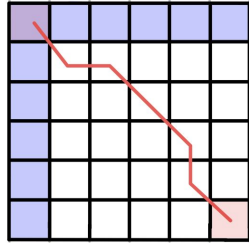
value associated to a characteristic of the protein

Liao, Yi-Lun & Li, Yu-Cheng & Chen, Nae-Chyun & Lu, Yi-Chang. (2018). Adaptively Banded Smith-Waterman Algorithm for Long Reads and Its Hardware Accelerator
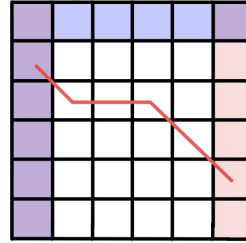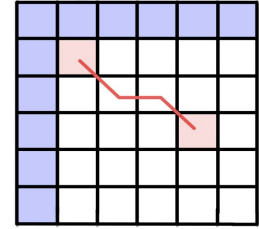
# Workflow of the alignment program

# Differences between alignments



|  | Global Alignment | Glocal Alignment | Local Alignment |
|---|---|---|---|
| alignment matrix | Needleman & Wunsch algorithm | Needleman & Wunsch algorithm | Smith & Waterman algorithm |
| **Traceback** |  |  |  |
| Initialization | Bottom right | Right column | Anywhere |
| Terminaison | Top left | Left column | First position with 0 encountered |

# Comparison of the 1RDS/1DE3A and 1BIF/1GC7A global alignments

Target : 1RDS     Query : 1DE3A

```
1 - ESC-E-YTCGS------T-C--Y-------WSS-DVS-AA-K-AKGY-SL-Y-E-S-G-DTIDDYPHEY-H-D-Y-E--GFD-F--P----------------V----------S-G-TYYEYPIMSDY-D-V-
1 - AV-T-W-TCLNDQKNPKTN-KY-ETKRLLYNQNK-AES--N-S---HH-AP-L-S-D-GKTGSSYPHW-FTN-G-Y-DG-D-G-KLPKGRTPIKFGKSDCDRPP-KHSKDGNGKT-D-HYLLEFPTFPDGH-D-Y
```

```
2 - YTG------G-SPGADRVIFNG-D-D-ELAGVITHTGASG-DDFVACSS-S
2 - -KFDSKKPK-ENPGPARVIYTYP-N-KVFCGIIAHTKENQG-ELKLCS-H-
```

Bottom right alignment matrix score : 2173.32          TMscore = 0.78741

Target : 1BIF     Query : 1GC7A

```
1 - CPT--LI--V--M---VGLPARGK------T-YI-SKKLTR-Y-L-NFIGVPTREF--N-V-GQ-Y-RRDMVK-TYKS-FEFFL-PDN-EEGLKIRK--Q-CALAALND-VRK-FLSEE-G-GHV-A-VF--DAT
1 - MPKPI--NV-RV-TTM---DAELEFAIQPN-TT-G-KQLFDQ-V-VKTVGLR--EVWF-F-G-LQYVD--S--KGYS-T--WLKL-N-K--------KVTQQ------DVK-KE----N-P-L-Q-F--KFRA-
```

```
2 - NTTRERRAMIFNFGEQNGYKTF-FVESICVDPEVIAANIVQVKLGSPDYVNRDSDEATEDFMRRIECYENSYESLDEEQDRDLSYIKIMDVGQSYVVNRVADHIQSRIVYYLMNIHVTPR
2 - --------K-------F-F------P------------E-------------------------------------------------------------------------------
```

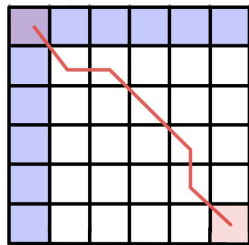Bottom right alignment matrix score : 580.04          TMscore = 0.16355

# Conclusion & prospect

- Successful implementation of alignments using embeddings
- Consistent results with TMscores

- Need to compute the right value for the gap penalty
- Implement affine gap penalty to compute alignment matrix

**Thank you**

# Differences between alignments



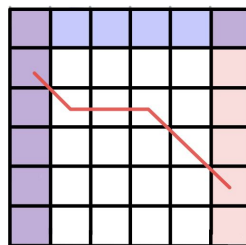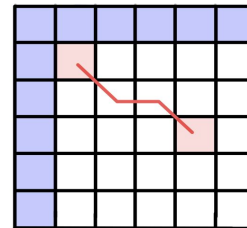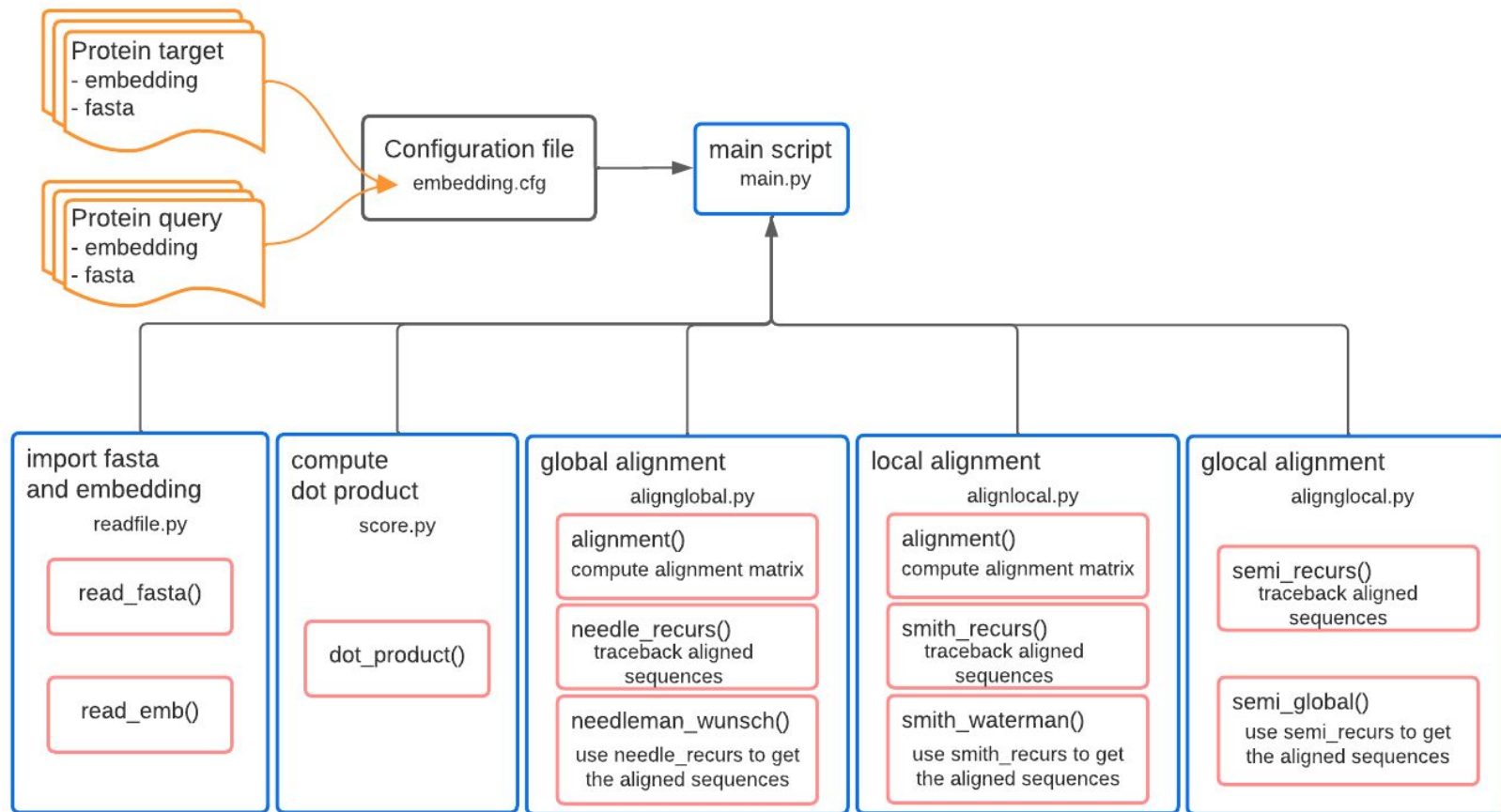| | Global Alignment | Glocal Alignment | Local Alignment |
|---|---|---|---|
| alignment matrix | $\max \begin{cases} \text{mat}[i-1, j-1] + \text{dot}[i, j] \\ \text{mat}[i, j-1] + \text{gap} \\ \text{mat}[i-1, j] + \text{gap} \end{cases}$ | $\max \begin{cases} \text{mat}[i-1, j-1] + \text{dot}[i, j] \\ \text{mat}[i, j-1] + \text{gap} \\ \text{mat}[i-1, j] + \text{gap} \end{cases}$ | $\max \begin{cases} \text{mat}[i-1, j-1] + \text{dot}[i, j] \\ \text{mat}[i, j-1] + \text{gap} \\ \text{mat}[i-1, j] + \text{gap} \\ 0 \end{cases}$ |
| **Traceback** | | | |
| Initialization | Bottom right | Right column | Anywhere |
| Terminaison | Top left | Left column | First position with 0 encountered |

# Structure of the alignment program

# Configuration file

```ini
1   [paths]
2       to_data = ../data/
3       to_res = ../results/
4
5   [files]
6       # protein target
7       prot_int_emb = rnase_1rds.t5emb
8       prot_int_fasta = RNASE_1RDS.fasta
9       # protein query
10      prot_comp_emb = RNase_U2_1de3a.t5emb
11      prot_comp_fasta = RNASE_U2_1DE3A.fasta
12
13      dot = dotprod_matrice.txt
14      align = align_matrice.txt
15
16  [alignment]
17      # True or False
18      global = True
19      local = True
20      glocal = True
21      # int
22      gap = 0
23
```

if True : the alignment is generated
all alignments can be generated in one time

the value of the gap penalty