

Zadanie 4

Wojciech Miśta, 236453

Działanie programu:

G8Drosophilahybri	TGG _ T TGA _ A _ A _ A _ CG _ AA _ _ _ C _ G _ GCA _ G _ GCA	1 - 59
T15Drosophilahybr	TGGACTT _ TCA _ TCA _ A _ A _ TGTTTTTC _ _ _ C _ C _ GCATA _ G _ A	1 - 59
G11Drosophilahybr	T _ G _ TATAATT _ AGCAGCAG _ CG _ AAGCTGCCGCCGCC _ GCAGCA _ GTT _ GGA	1 - 59
C6Drosophilahybri	T _ G _ TATGA _ A _ AT _ ATTCT _ AT _ _ _ CTG _ GAA _ GATATCA	1 - 59
	* * _ * * * _ * _ * _ _ _ _ * _ * _ *	
G8Drosophilahybri	_ T _ AA _ C _ G _ T _ TTG _ C _ CT _ T _ T _ T _ GTTA _ C _	60 - 118
T15Drosophilahybr	_ A _ AACTTTTTTC _ G _ T _ A _ TTGTC _ ATCT _ AT _ TCG _ GTT _	60 - 118
G11Drosophilahybr	GGTGAACAA _ C _ GAAACATACAGAATG _ CGA _ CT _ TAT _ TCATTT _ C _ T	60 - 118
C6Drosophilahybri	_ TG _ AAT _ CATTGAG _ T _ TCA _ C _ CTC _ G _ T _ T _ CCTACCA _	60 - 118
	_ _ _ ** _ * _ * _ * _ _ _ * _ ** _ * _ *	
G8Drosophilahybri	C _ G _ G _ GC _ GC _ A _ G _ A _ _ _ A _ C _ A _ GA _ AT _ C _ C _ G _ T _ G	119 - 177
T15Drosophilahybr	_ G _ TG _ GCATC _ A _ TCGAA _ TTGTCA _ _ _ CC _ A _ GA _ AT _ T _ GTT _ T	119 - 177
G11Drosophilahybr	C _ _ A _ AC _ G _ A _ _ AC _ _ A _ TTTT _ C _ A _ TTTACTTCC _ ACAGAT _ G	119 - 177
C6Drosophilahybri	CTGTT _ GCGC _ CTAC _ A _ A _ _ ATG _ _ _ CTTAT _ AA _ AT _ CA _ T _ G _ TAGG	119 - 177
	_ _ _ * _ * _ * _ * _ _ _ * _ * _ * _ * _ *	
G8Drosophilahybri	_ _ _ _ _ AGA _ GGA _ A _ A _ AC _ TCAG _ G _ GA _ A _ AAG _ AT _ C _ CT _ CT	178 - 236
T15Drosophilahybr	_ _ _ _ _ ATA _ ATC _ A _ AT _ TC _ _ _ GATACT _ TGGTAT _ AC _ TTGGTACT _	178 - 236
G11Drosophilahybr	_ _ _ _ _ CTACACGCATATC _ GCGCCAG _ GTTTC _ A _ C _ AA _ A _ _ CT _ GTA	178 - 236
C6Drosophilahybri	TCTTACAA _ AGA _ AGA _ A _ AGA _ _ AGCA _ GA _ A _ AGAA _ ATA _ CA _ CA _ C _	178 - 236
	_ _ _ * _ * _ _ _ _ _ * _ _ _ *	

Dopasowanie dla gatunków z rodzaju 'Drosophila'.

```

>G8Drosophilahybri
TGG _ T TGA _ A _ A _ A _ CG _ AA _ _ _ C _ G _ GCA _ G _ GCA _ T _ AA _
_ C _ G _ T _ TTG _ C _ CT _ T _ T _ T _ GTTA _ C _ C _ G _ G _ GC _ GC _ A _ G _ A _
_ A _ _ C _ A _ GA _ AT _ C _ C _ G _ T _ G _ _ _ AGA _ GGA _ A _ A _ AC _ TCAG _ G _
GA _ A _ AAG _ AT _ C _ CT _ CT _ C _ CC

>T15Drosophilahybr
TGGACTT _ TCA _ TCA _ A _ A _ TGTTTTTC _ _ _ C _ C _ GCATA _ G _ A _ A _ AACT
TTTTT _ G _ T _ A _ TTGTC _ ATCT _ AT _ TCG _ GTT _ _ _ G _ TG _ GCATC _ A _ TCGAA _
TTGTCA _ _ CC _ A _ GA _ AT _ T _ GTT _ T _ _ _ ATA _ ATC _ A _ AT _ TC _
GATACT _ TGGTAT _ AC _ TTGGTACT _ C _ AATA

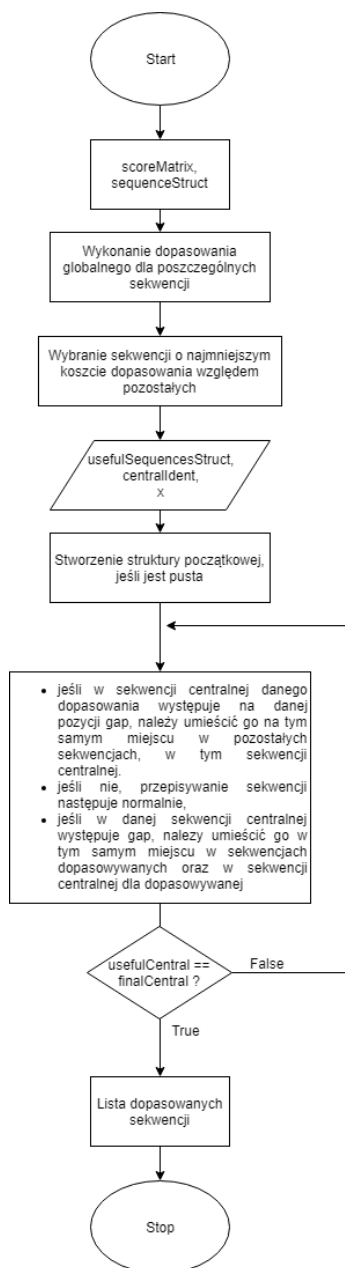
>G11Drosophilahybr
T _ G _ TATAATT _ AGCAGCAG _ CG _ AAGCTGCCGCCGCC _ GCAGCA _ GTT _ GGAGGTGAACAA _
_ C _ GAAACATACAGAATG _ CGA _ CT _ TAT _ TCATTT _ C _ TC _ _ A _ AC _ G _ A _ A _ AC
_ A _ TTTT _ C _ A _ TTTACTTCC _ ACAGAT _ G _ _ CTACACGCATATC _ GCGCCAG _ GTT
TC _ A _ C _ AA _ A _ CT _ GTAGC _ _ CC

>C6Drosophilahybri
T _ G _ TATGA _ A _ AT _ ATTCT _ AT _ _ _ CTG _ GAA _ GATATCA _ TG _ AAT _
_ CATTGAG _ T _ TCA _ C _ CTC _ G _ T _ T _ CCTACCA _ CTGTT _ GCGC _ CTAC _ A _ A _
_ ATG _ CTTAT _ AA _ AT _ CA _ T _ G _ TAGGCTTACAA _ AGA _ AGA _ A _ AGA _ _ AGCA _
GA _ A _ AGAA _ ATA _ CA _ CA _ C _ CAAG _ GG
  
```

Plik .fasta

Działanie programu dla krótkich sekwencji z zajęć:

i dB	__CGT__	1 - 8
i dA	AACGT__	1 - 8
i dC	__CGTAAA	1 - 8
i dD	__CAT__	1 - 8
		* *



Komentarze znajdujące się wewnątrz funkcji *multipleSequenceAlignment.m* tłumaczą dokładnie działanie algorytmu.

multipleSequenceAlignm.m

```
|function [finalSequenceStruct] = multipleSequenceAlignm(usefulSequencesStruct,x,centralIdent)
|%MULTIPLESEQUENCEALIGNMENT Summary of this function goes here
|% Detailed explanation goes here
|finalSequenceStruct = struct;

|%je?li dopasowywany powoduje gapa w centralnym -> dodaj gap do wszystkich
|%poza dopasowywanym
|%je?li mismatch to nic
|%je?li centralny pierwszy ma gap, dodaj gap do dopasowywanego

|loop przez ilosc dopasowan jakie b?d? robione
|for i = 1:x
|    backup = usefulSequencesStruct;
|    %jesli nie ma nic w finalnej strukturze
|    if (isempty(fieldnames(finalSequenceStruct)))
|        finalSequenceStruct.(centralIdent) = usefulSequencesStruct.("s" + i).(centralIdent);
|        usefulSequencesStruct.("s" + i) = rmfield(usefulSequencesStruct.("s" + i),centralIdent);
|        lastField = char(fieldnames(usefulSequencesStruct.("s" + i)));
|        finalSequenceStruct.(lastField) = usefulSequencesStruct.("s" + i).(lastField);
|    else
|        %jesli struktury do dopasowania istnieja
|        if (isfield(usefulSequencesStruct,char("s"+i)))
|            %Centralna sekwencja z danego dopasowania
|            usefulCentral = usefulSequencesStruct.("s" + i).(centralIdent);

|            %usun z dopasowan zeby uzyskac pozosa??
|            usefulSequencesStruct.("s" + i) = rmfield(usefulSequencesStruct.("s" + i),centralIdent);

|            %Pozosta?a nazwa field sekwencji z danego dopasowania
|            lastField = char(fieldnames(usefulSequencesStruct.("s" + i))); %dopasowywana

|            %pozosta?a sekwencja
|            lastSequence = char(usefulSequencesStruct.("s" + i).(lastField));

|            %Field names w finalnej strukturze
|            existingSequencesFields = fieldnames(finalSequenceStruct);

|
|
|    c = 1;
|    %porównujemy czy finalCentral == usefulCentral(którą będziemy
|    %nadpisywac)
|    while ~strcmp(usefulCentral,char(finalSequenceStruct.(char(existingSequencesFields(1)))))
|        finalCentral = char(finalSequenceStruct.(char(existingSequencesFields(1))));
|        %patrzmy czy jest gap w starym dopasowaniu
|        if (numel(usefulCentral) >= c)
|            if (usefulCentral(c) == '_')
|                %sprawdzamy czy usefulCentral > finalCentral bo
|                %outofbound
|                if (numel(usefulCentral) <= numel(finalCentral))
|                    %patrzmy czy jest gap w NOWYM dopasowaniu
|                    if (finalCentral(c) == '_')
|                        %nie rób nic jest gap jest w obu dopasowaniach
|                    else
|                        %jeśli w nowym dopasowaniu nie ma gapa
|                        %nadpisz nowe dopasowanie + wszystkie
|                        for n = 1:numel(existingSequencesFields)
|                            currentSeq = char(finalSequenceStruct.(char(existingSequencesFields(n))));
|                            currentField = char(existingSequencesFields(n));

|                            if (c ~= 1)
|                                currentSeq = char(strcat(currentSeq(1:c-1) + "_" + currentSeq(c:end)));
|                            elseif (c == numel(currentSeq))
|                                currentSeq = char(strcat(currentSeq + "_"));
|                            elseif (c == 1)
|                                currentSeq = char(strcat("_" + currentSeq));
|                            end
|                            finalSequenceStruct.(currentField) = currentSeq;
|                        end
|                    end
|                else
|                    %finalCentral jest mniejsze of useful <-
|                    %outofbound
|                    for n = 1:numel(existingSequencesFields)
|                        currentSeq = char(finalSequenceStruct.(char(existingSequencesFields(n))));
```

```

currentSeq = char(finalSequenceStruct.(char(existingSequencesFields(n))));
currentField = char(existingSequencesFields(n));
if(c ~= 1)
    currentSeq = char(strcat(currentSeq(1:c-1) + "_" + currentSeq(c:end)));
elseif(c == numel(currentSeq))
    currentSeq = char(strcat(currentSeq + "_"));
elseif(c == 1)
    currentSeq = char(strcat("_" + currentSeq));
end
finalSequenceStruct.(currentField) = currentSeq;
end
end
%jeśli finalna ma gapa, a w starym nie ma
elseif(finalCentral(c) == '_')
    %dodaj gapy do dopasowywanej
    if(c ~= 1)
        lastSequence = char(strcat(lastSequence(1:c-1) + "_" + lastSequence(c:end)));
        usefulCentral = char(strcat(usefulCentral(1:c-1) + "_" + usefulCentral(c:end)));
    elseif(c == numel(lastSequence))
        lastSequence = char(strcat(lastSequence + "_"));
        usefulCentral = char(strcat(usefulCentral + "_"));
    elseif(c == 1)
        lastSequence = char(strcat("_" + lastSequence));
        usefulCentral = char(strcat("_" + usefulCentral));
    end
    %dopisz gap do nowej sekwencji
elseif(numel(usefulCentral) < numel(finalCentral))
    lastSequence = char(strcat(lastSequence + "_"));
    usefulCentral = char(strcat(usefulCentral + "_"));
    %dodac warunek
end
end
c = c + 1;
end
finalSequenceStruct.(lastField) = lastSequence;
end
end
usefulSequencesStruct = backup;
end
end

```

Złożoność czasowa:

n – ilość sekwencji wejściowych

k – górna granica sekwencji wejściowych (np. przy *for* w *while*)

- n – iteracji pierwszej pętli
 - maksymalnie k iteracji pętli *while*
 - maksymalnie $n * k$ iteracji obu pętli *for*, więc $2nk$
 - maksymalnie $3k$ możliwych sprawdzeń warunku

$$n * (k * (2nk + 3k)) = 2k^2n^2 + 3k^2n = O(k^2n^2 + k^2n)$$

Złożoność przestrzenna:

Maksymalny rozmiar sekwencji wejściowych: $2n * k$ <- n par sekwencji

Maksymalny rozmiar finalnej struktury: nk

$$2nk + nk = 3nk \rightarrow O(nk)$$