

# Биоинформатика

Коротков Евгений Вадимович  
Институт Биоинженерии, ФИЦ Биотехнологии  
РАН

[bioinf@yandex.ru](mailto:bioinf@yandex.ru)

Множественные выравнивания

Профили

## There are essentially four major approaches to multiple sequence alignment:

1. Optimal global sequence alignment . It is NP-complete problem.

2. Progressive global alignment

3. Block-based global alignment

4. Motif-based local alignment

5. Markov models

Methods 2-5 don't work for sequences which have more 1.5 mutation per amino acid.

# Обобщение парного выравнивания

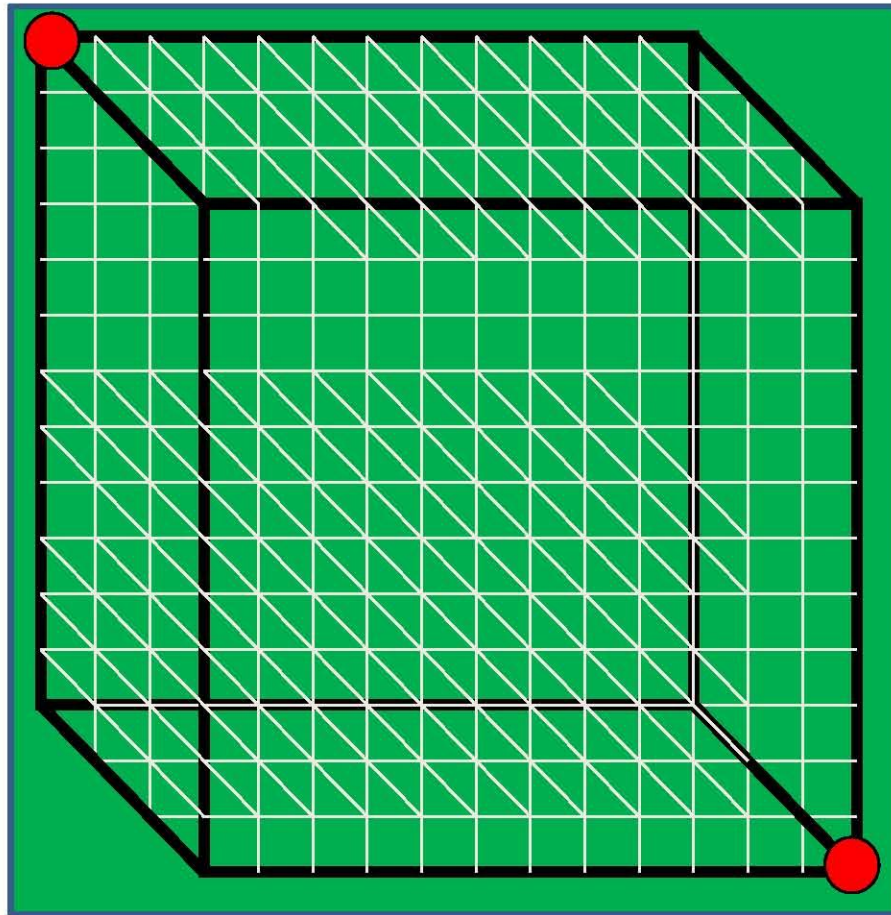
- Выравнивание 2-х последовательностей – двумерная матрица
- 3-х последовательностей – 3-х мерная.

A	T	_	G	C	G	_
A	_	C	G	T	_	A
A	T	C	A	C	_	A

- Задача: больше консервативных столбцов, лучше выравнивание

# Глобальное выравнивание 3-х последовательностей

начало



конец

# Алгоритм

- $$s_{i,j,k} = \max \left\{ \begin{array}{l} s_{i-1,j-1,k-1} + \delta(v_i, w_j, u_k) \\ s_{i-1,j-1,k} + \delta(v_i, w_j, \_) \\ s_{i-1,j,k-1} + \delta(v_i, \_, u_k) \\ s_{i,j-1,k-1} + \delta(\_, w_j, u_k) \\ s_{i-1,j,k} + \delta(v_i, \_, \_) \\ s_{i,j-1,k} + \delta(\_, w_j, \_) \\ s_{i,j,k-1} + \delta(\_, \_, u_k) \end{array} \right\}$$

Нет гэпов

Один гэп

Два гэпа

- $\delta(x, y, z)$  – запись в трехмерной матрице весов

# Время работы алгоритма

- Для 3-х последовательностей длины  $n$ , время работы –  $7n^3$ ;  $O(n^3)$
- Для  $k$  последовательностей -  $(2k-1)(n^k)$ ;  $O(2kn^k)$

# Множественное выравнивание порождает парные выравнивания

**x:** AC-GCGG-C  
**y:** AC-GC-GAG  
**z:** GCCGC-GAG

Порождает:

**x:** ACGCGG-C ;    **x:** AC-GCGG-C ;    **y:** AC-GCGAG  
**y:** ACGC-GAC ;    **z:** GCCGC-GAG ;    **z:** GCCGCGAG

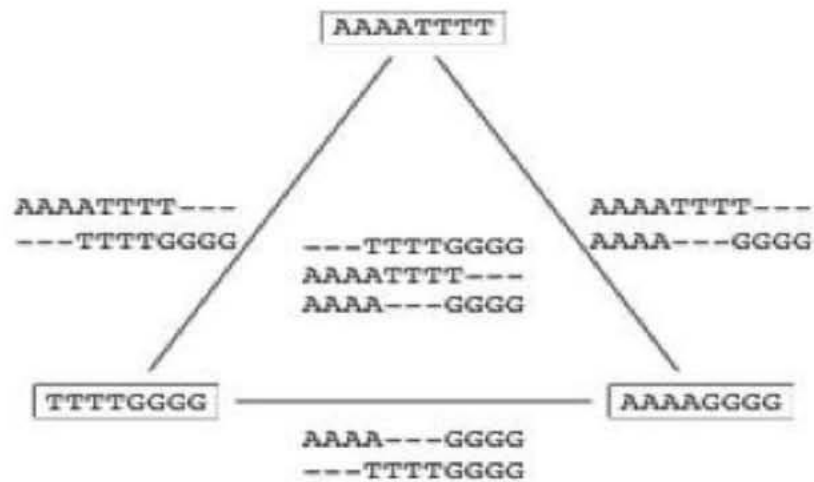


# Обратная проблема

Имея 3 субъективных парных варнивания:

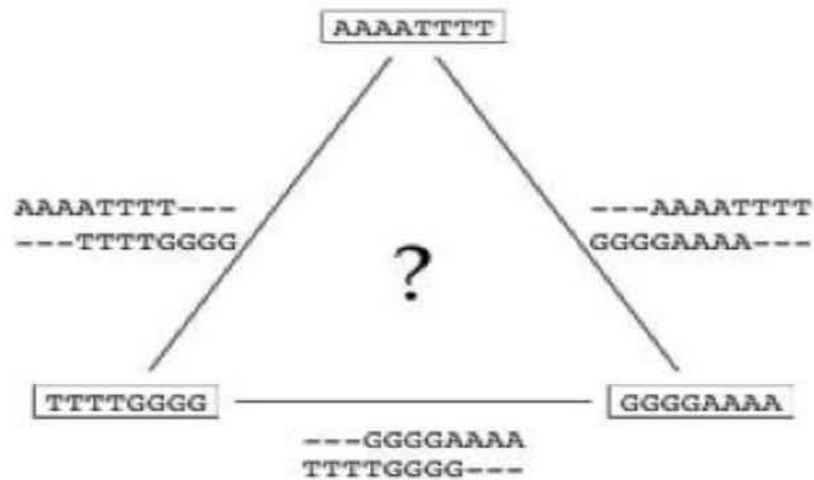
**x**: ACGCGG-C ;    **x**: AC-GCGG-C ;    **y**: AC-GCGAG  
**y**: ACGC-GAC ;    **z**: GCCGC-GAG ;    **z**: GCCGCGAG

Хороший вариант



(a) Compatible pairwise alignments

Плохой вариант



(b) Incompatible pairwise alignments

# Выравнивание выравниваний

x GGGCACTGCAT  
y GGTTACGTC--  
z GGGAACTGCAG

Alignment 1

w GGACGTACC--  
v GGACCT-----

Alignment 2

# Описание выравнивания

**GTC**TA**GA**  
**GTC**AG**C** } **GTC**[TA]**G**[AC] - профиль  
                  [X] [5X]

x GGGCACTGCAT

y GGTTACGTC--

z GGGAAC**TGCAG**

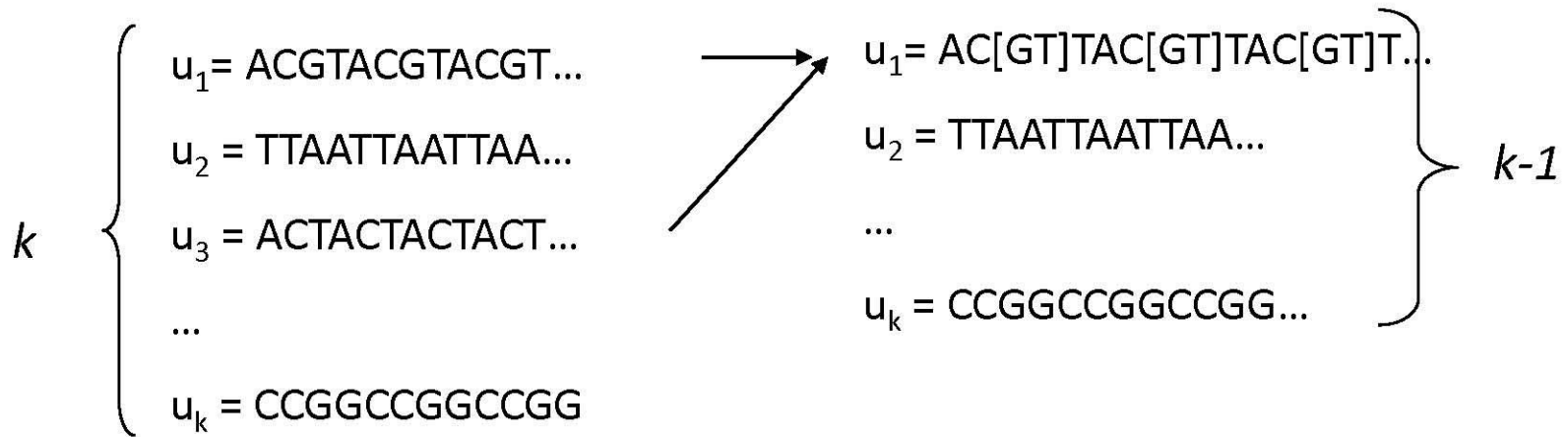
w GGACGTACC--

v GGACCT-----

**GGACACAGCAT** - консенсус

Матрица частот – используется редко

# Множественное выравнивание – жадный алгоритм



Время работы алгоритма на  $k$  последовательностях длины  $n$  –  $O(n^2k^2)$

# Прогрессивное выравнивание ClustalW

- Прогрессивное выравнивание – жадный алгоритм с более «умным» способом выбора пар.
- Три шага
  - 1.) Построить парные выравнивания
  - 2.) Построить дерево-подсказку
  - 3.) Прогрессивное выравнивание по дереву-подсказке

# Шаг 1: Парные Выравнивания

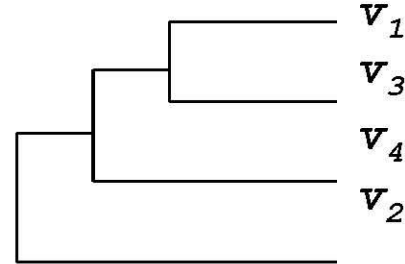
- Выравнивания пар порождают матрицу identity

	$v_1$	$v_2$	$v_3$	$v_4$
$v_1$				
$v_2$	.17	-		
$v_3$	.87	.28	-	
$v_4$	.59	.33	.62	-

(.17 значит идентичны на 17 %)

# Шаг 2: Дерево-подсказка

	$v_1$	$v_2$	$v_3$	$v_4$
$v_1$				
$v_2$	.17	-		
$v_3$	.87	.28	-	
$v_4$	.59	.33	.62	-



Далее вычислить:


- $v_{1,3}$  = выравнивание  $(v_1, v_3)$
- $v_{1,3,4}$  = выравнивание  $((v_{1,3}), v_4)$
- $v_{1,2,3,4}$  = выравнивание  $((v_{1,3,4}), v_2)$



# Шаг 3: Прогрессивное выравнивание

- Выравниванием 2 наиболее близких последовательности.
- Следуя дереву - подсказке, довыравниваем следующую последовательность к имеющемуся выравниванию

```
FOS_RAT      PEEMSVTS-LDLTGGLPEATTPESSEEAFTLPLLNDPEPK-PSLEPVKNISNMELKAEFPD
FOS_MOUSE   PEEMSVAS-LDLTGGLPEASTPESEEAFTLPLLNDPEPK-PSLEPVKISNVELKAEFPD
FOS_CHICK    SEELAAATALDLG----APSPAAAEAFALPLMTEAPPAVPPKEPSG--SGLELKAEFPD
FOSB_MOUSE  PGPGLAEVRDLPG-----STSAKEDGFGWLLPPPPPPP-----LPFQ
FOSB_HUMAN  PGPGLAEVRDLPG-----SAPAKEDGFSWLLPPPPPPP-----LPFQ
.           . : ** . :.. *:. * * . * **:
```



Точки и звезды отображают насколько консервативны столбцы.

# Множественные Выравнивания: Взвешивание

- Количество полных совпадений
- Сумма по парам (SP-Score)
- Энтропия

# Количество полных совпадений

ААА  
ААА  
ААТ  
АТС

- Хорошо только для очень близких последовательностей

# Сумма по парам (SP-Score)

- Построим парное выравнивание по множественному
- Посчитаем веса всех этих парных выравниваний –  $s(a_i, a_j)$
- Просуммируем:

$$s(a_1, \dots, a_k) = \sum_{i,j} s(a_i, a_j)$$

# Энтропия

- Определим вероятности букв в столбцах
  - $p_A = 1, p_T=p_G=p_C=0$  (1-ый столбец)
  - $p_A = 0.75, p_T = 0.25, p_G=p_C=0$  (2-ый столбец)
  - $p_A = 0.50, p_T = 0.25, p_C=0.25, p_G=0$  (3-ий столбец)
- Энтропия столбца будет равна

$$- \sum_{X=A,T,G,C} p_X \log p_X$$

AAA  
AAA  
AAT  
ATC

# Энтропия: Пример

Лучший вариант  $entropy \begin{pmatrix} A \\ A \\ A \\ A \end{pmatrix} = 0$

Худший вариант  $entropy \begin{pmatrix} A \\ T \\ G \\ C \end{pmatrix} = -\sum \frac{1}{4} \log \frac{1}{4} = -4\left(\frac{1}{4} * -2\right) = 2$

# Энтропия: Пример

Энтропия столбца:

$$-(p_A \log p_A + p_C \log p_C + p_G \log p_G + p_T \log p_T)$$

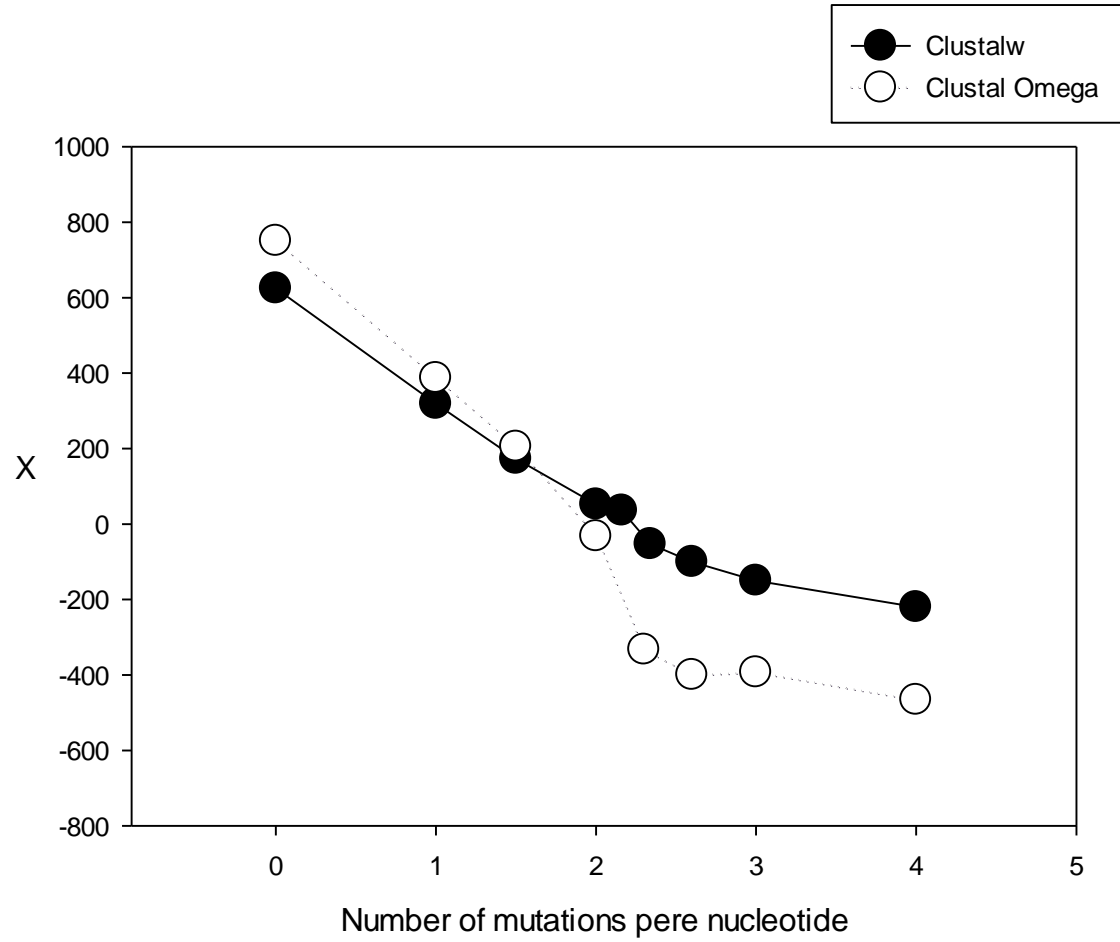
A	A	A
A	C	C
A	C	G
A	C	T

- Столбец 1 =  $-[1 \cdot \log(1) + 0 \cdot \log 0 + 0 \cdot \log 0 + 0 \cdot \log 0]$   
= 0

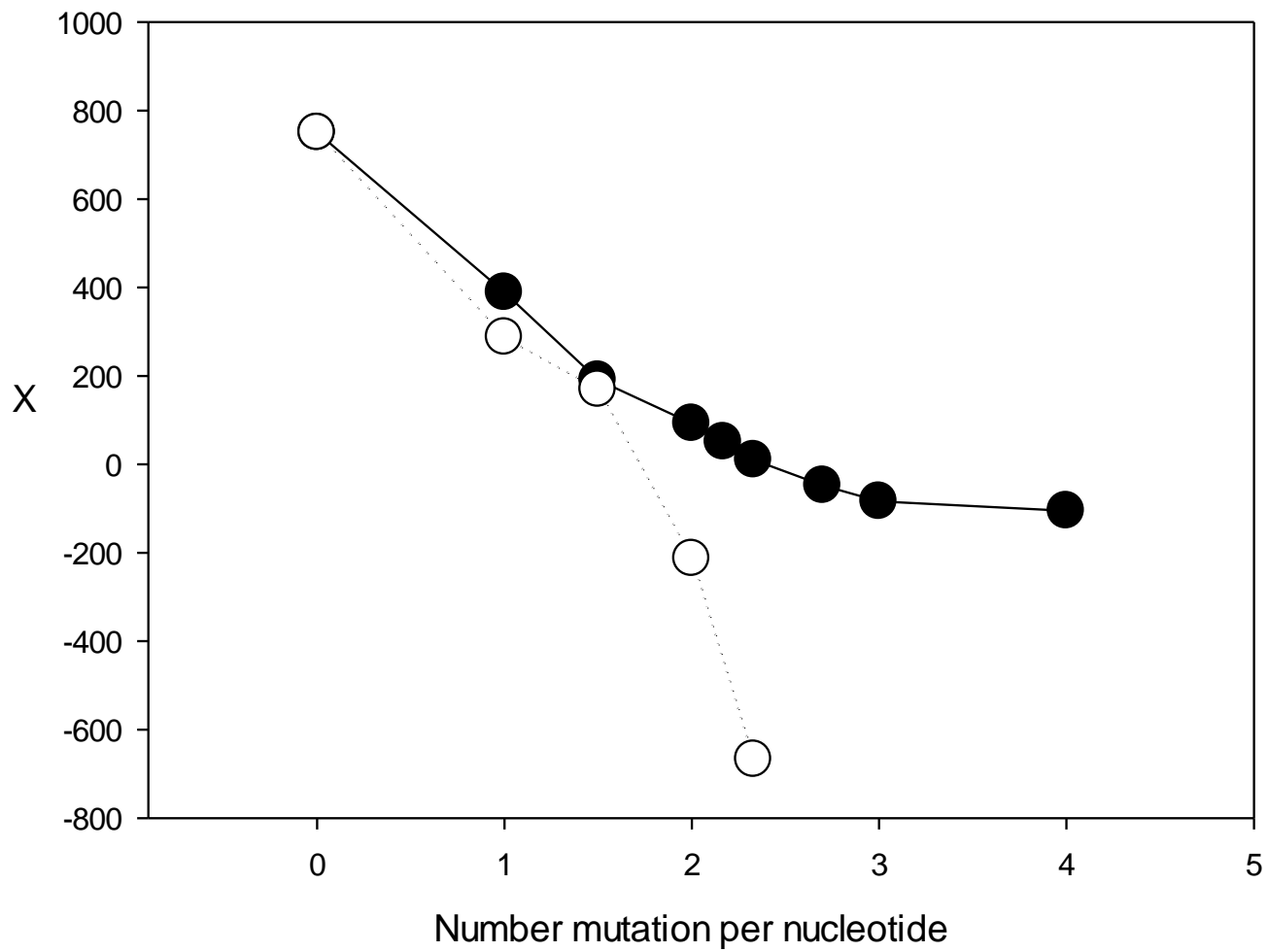
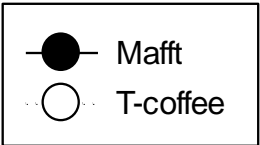
- Столбец 2 =  $-[(1/4) \cdot \log(1/4) + (3/4) \cdot \log(3/4) + 0 \cdot \log 0 + 0 \cdot \log 0]$   
=  $-[(1/4) \cdot (-2) + (3/4) \cdot (-.415)] = +0.811$

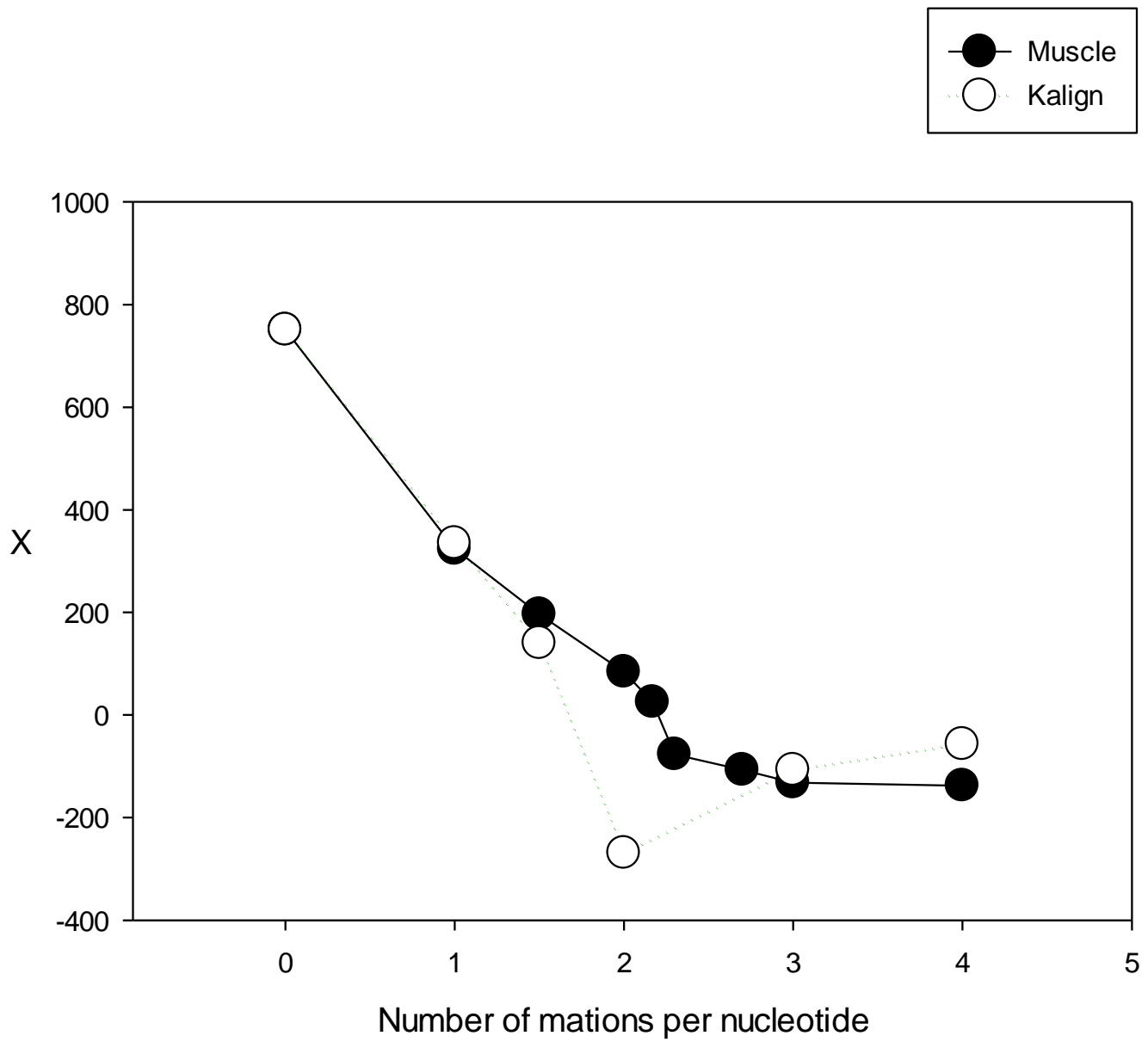
- Столбец 3 =  $-[(1/4) \cdot \log(1/4) + (1/4) \cdot \log(1/4) + (1/4) \cdot \log(1/4) + (1/4) \cdot \log(1/4)]$   
=  $4 \cdot -[(1/4) \cdot (-2)] = +2.0$

- Энтропия выравнивания =  $0 + 0.811 + 2.0 = +2.811$

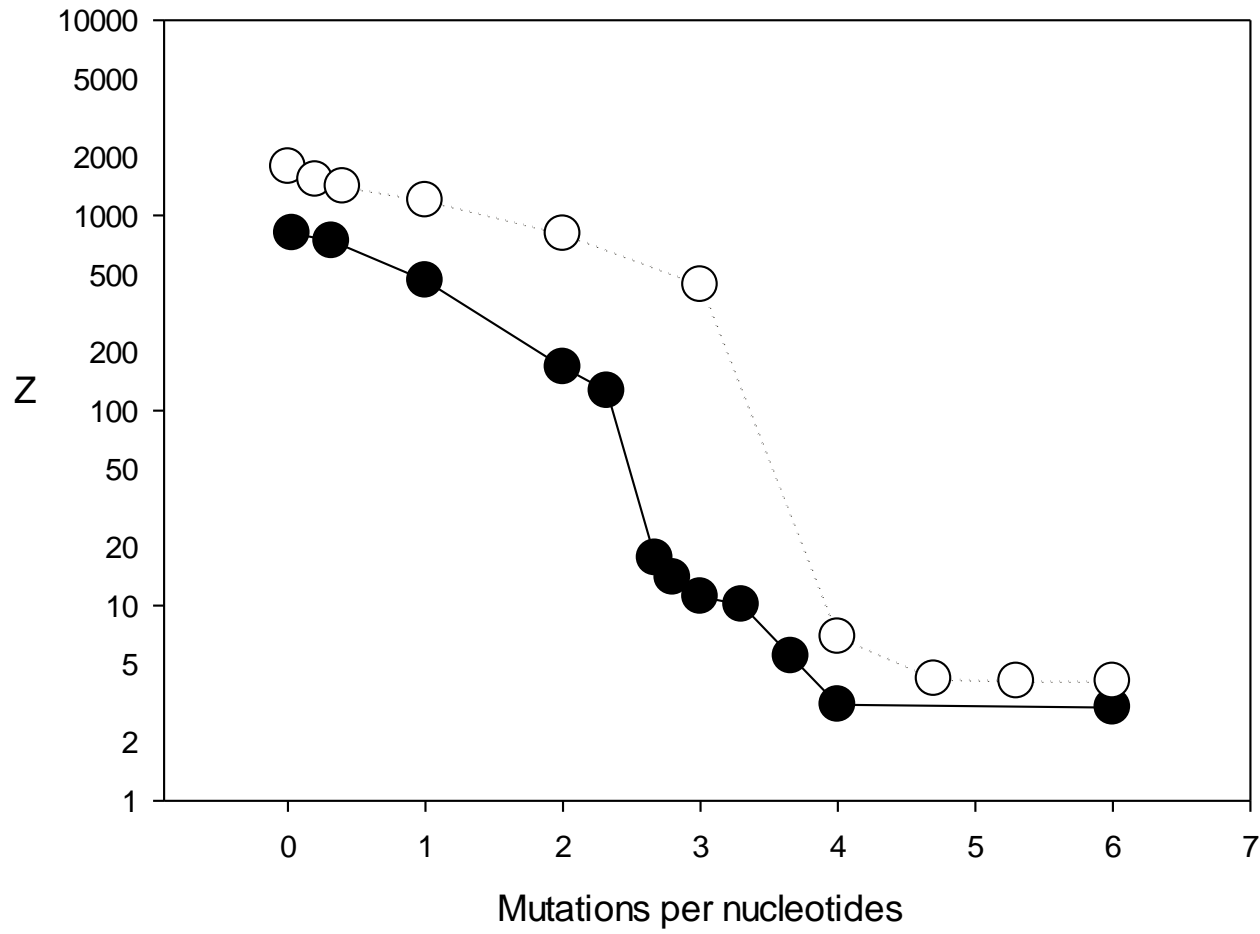








100\_50, 500\_250



# Algorithm

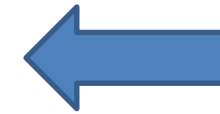
- How we can find the multiple alignment if the sequences have big number of mutations (more than 1.5 per amino acid or nucleotide)?
- Statistical important pairwise alignments, common “seeds” and “words” are absent.

## Position Weight Matrix (PWM)

1	2	3	4	5
A	1.2	-2.3	3.1	-1.0
F	-1.2	2.2	2.5	-0.8
K	1.8	-0.9	0.1	3.3
M	2.3	3.9	-0.8	1.8
...	...	...	...	...

Multi alignment

MAHNV  
RASSG  
LSYPE  
RFYAA  
ASYVL



Sequence A 123451234512345123451234512345...

Sequence B 

Sequence 1	Sequence 2	Sequence 3	Sequence 4	Sequence 5	...
------------	------------	------------	------------	------------	-----

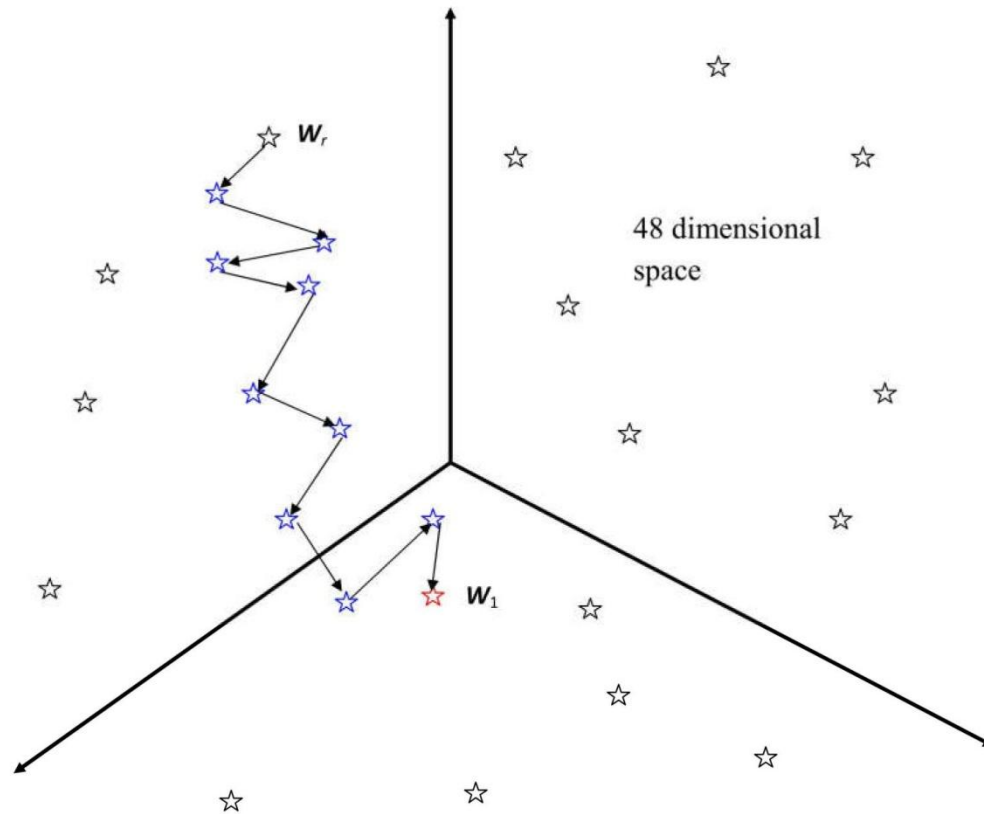
Task: to find PWM which describes multi alignment of analyzed sequences by west way.

# Algorithm for search latent periodicity with indels

- **Main idea:**
- Generate the set of the random position weights matrixes (PWM)
- Calculate the alignment of the sequence with PWM's
- Improve the matrixes by genetic algorithm

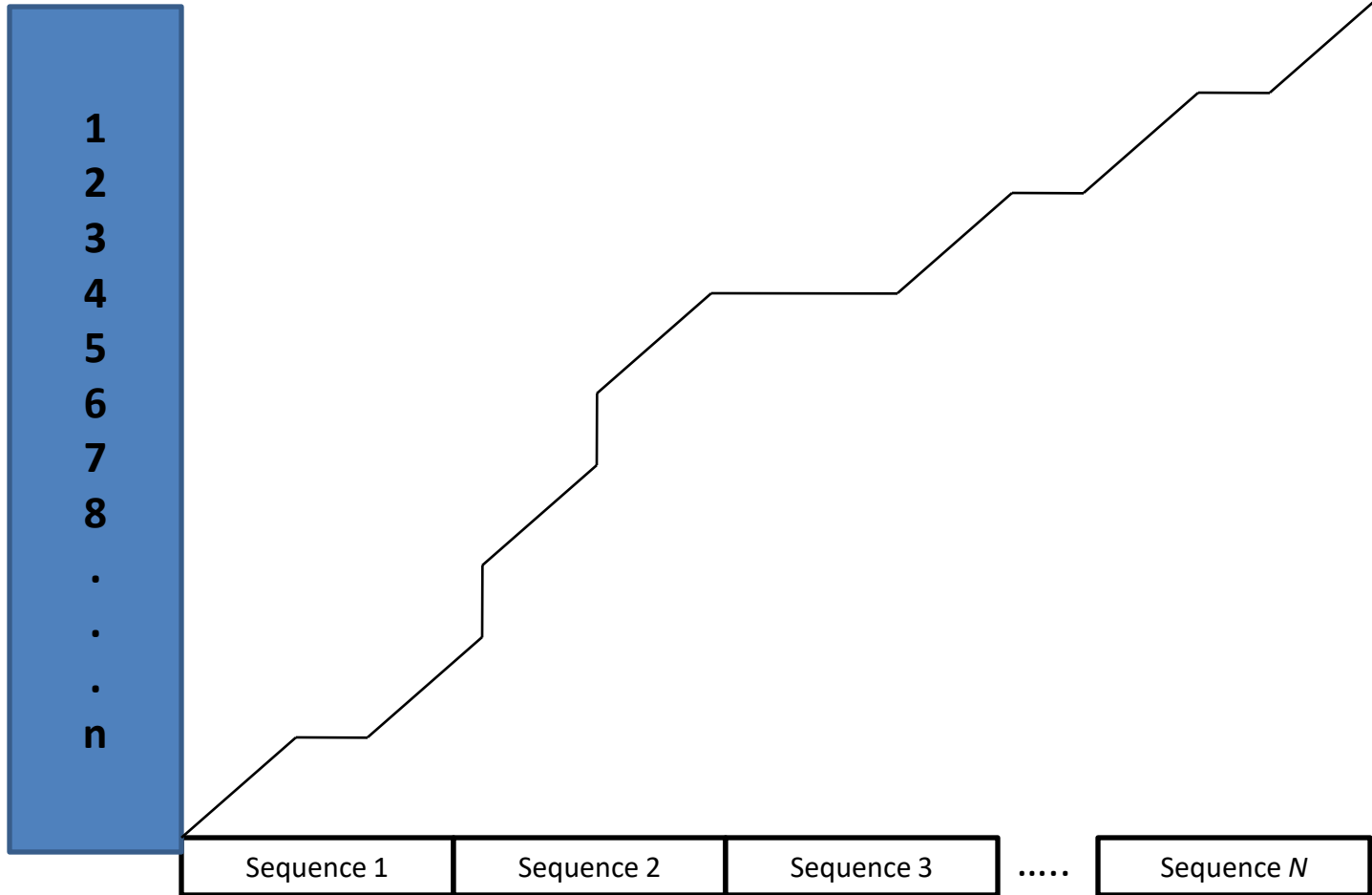
# Algorithm

- 1 Creation of the set of random PWM's. Each PWM corresponds to some set of multi alignments.
- 2 Union of all sequences in a single sequence one after another
- 3 Global alignment of united sequence with each random PWM.
- 4 Optimization of PWM by the genetic algorithm.
- 5 Selection of PWM with maximum of similarity function  $F$



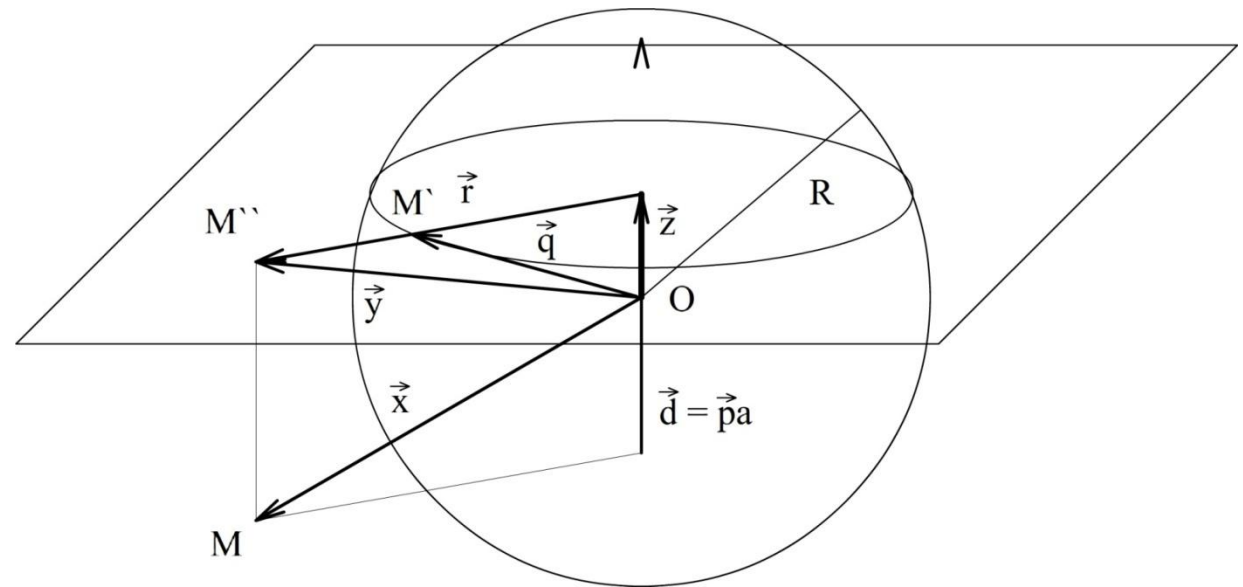


**PWM**



## Creation of the set of random PWM's

Space  $20n$



$$A = \sum_{i=1}^{20} \sum_{j=1}^n m(i, j)^2$$

$$B = \sum_{i=1}^{20} \sum_{j=1}^n m(i, j) p(i) f(j)$$

where  $f(j)=1/n$ ;  $p(i)=n(i)/N$ ,  $n(i)$ - number of  $a_i$  in sequence  $A$ ,  $N$  – length of sequence  $A$ ,  $n$  – length of period.

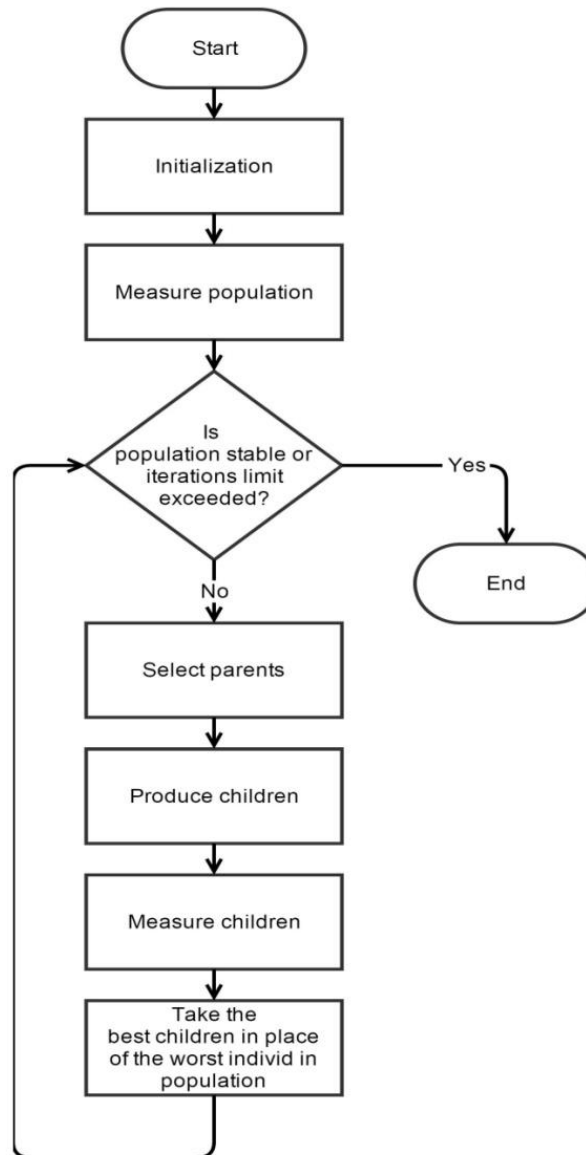
$$F(i, j) = \max \left\{ \begin{array}{l} F(i-1, j-1) + m(a(i), b(j)) \\ F_1(i-1, j-1) - d \\ F_2(i-1, j-1) - d \end{array} \right\}$$

$$F_1(i, j) = \max \left\{ \begin{array}{l} F(i-1, j) - d \\ F_1(i-1, j) - e \end{array} \right\}$$

$$F_2(i, j) = \max \left\{ \begin{array}{l} F(i, j-1) - d \\ F_2(i, j-1) - e \end{array} \right\}$$

$$Z = \frac{F_{\max} - M \left[ \vec{F} \right]}{\sqrt{D \left[ \vec{F}_{\max} \right]}}$$

# Improving the matrixes by genetic algorithms



Volume of PWM population is 500 matrixes for each  $n$ .  
 $n$  – period length

Goal is to increase the  $F_{\max}$ .

## Types of tandem repeats

### 1. Perfect periods

**GCAT GCAT GCAT GCAT GCAT...**

**GCAT  
 GCAT  
 GCAT  
 GCAT  
 GCAT**

There are no changes  
 per symbol

### 2. Periods with indels and changes of the symbols

**GGAT GCAG \*CAT GCGT G\*AT...**

**GGAT  
 GCAG  
 \*CAT  
 GCGT  
 G\*AT**

There are 1.0 or less  
 changes per symbol

## Types of tandem repeats

### 3. Latent periods without indels

**GCAT CGAC TACT AGGT GGTA...**

**GCAT  
CGAC  
TACT  
AGGT  
GGTA...**

There are 1.0 or more changes per symbol

### 4. Latent periods with indels

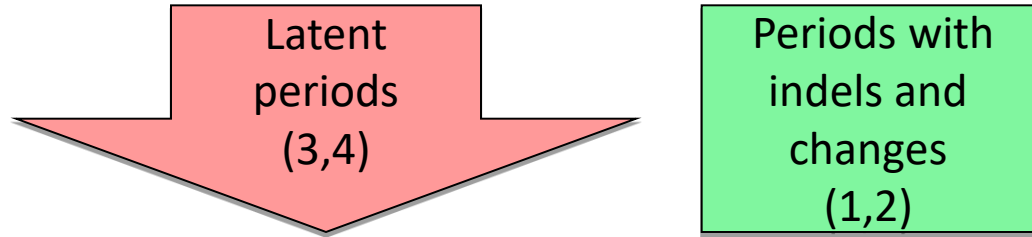
**G\*AT CGAC TA\*T AGGT GG\*A...**

**G\*AT  
CGAC  
TA\*T  
AGGT  
GG\*A...**

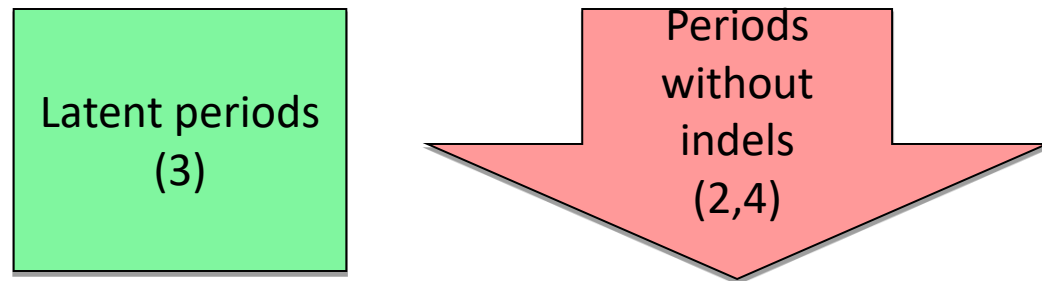
There are 1.0 or more changes per symbol

## Methods applied for periodicity search

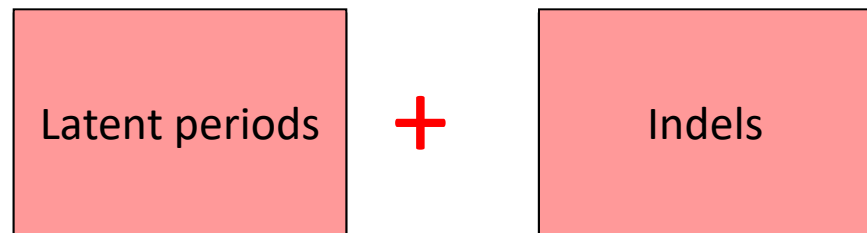
1. The dynamic programming



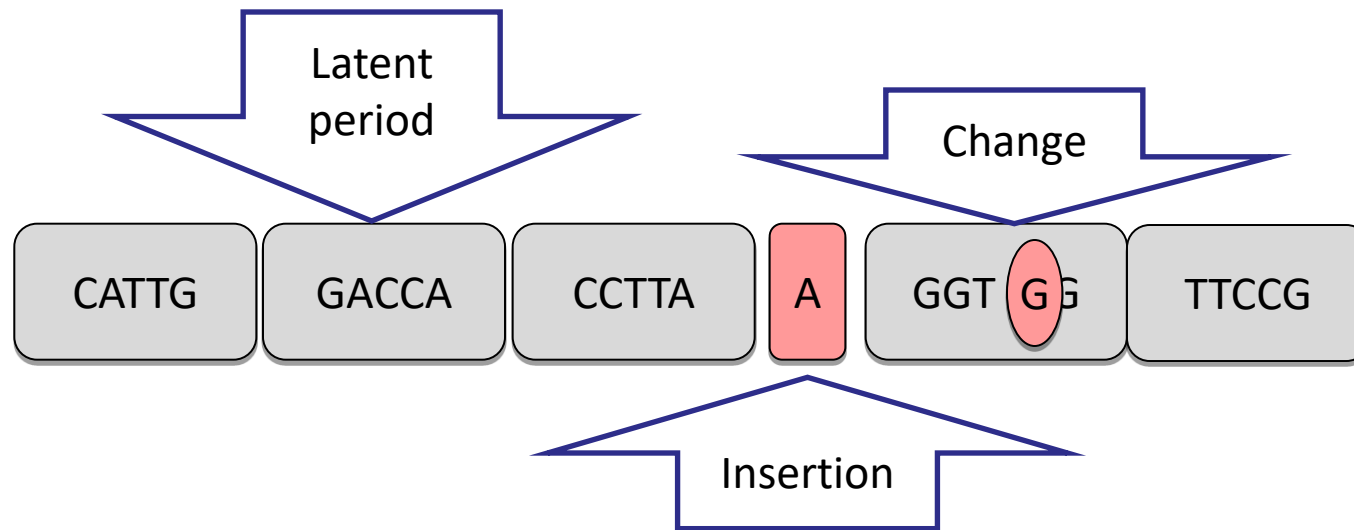
1. Spectral Approach (Fourier transform, Information Decomposition,...)



1. ?



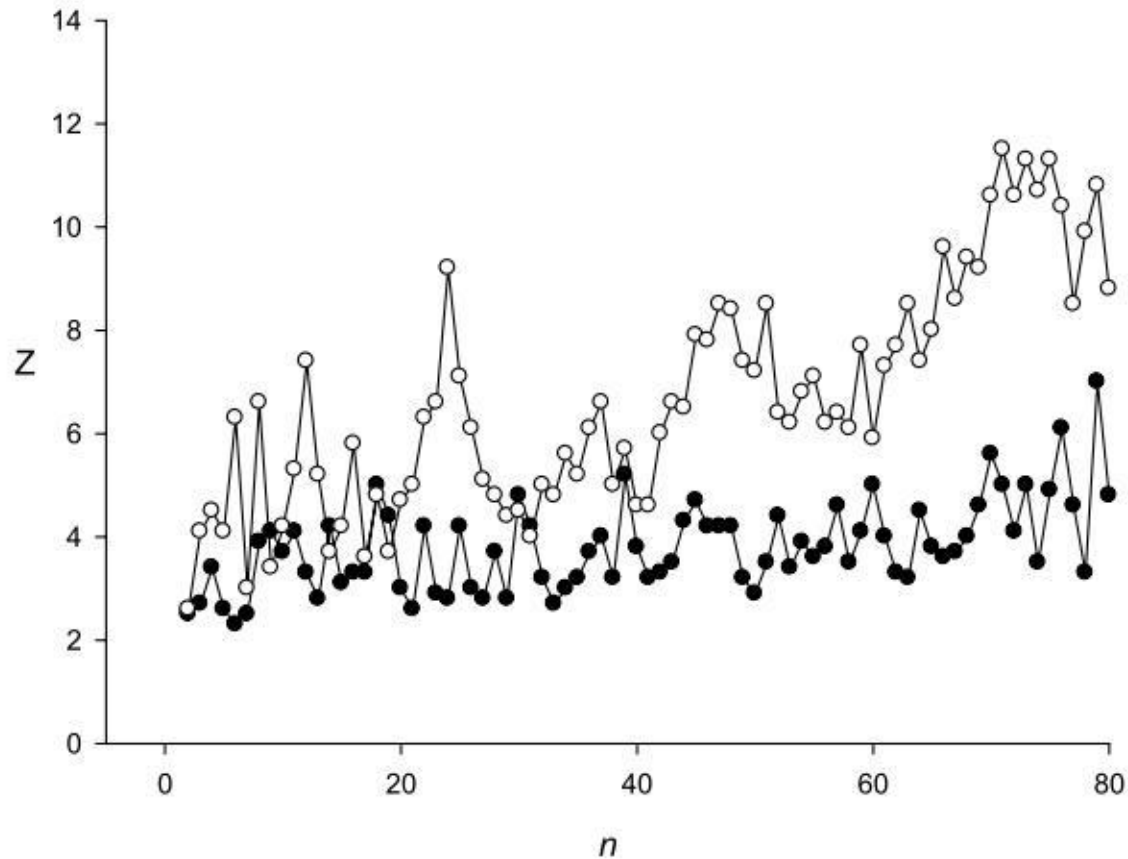
## Type periodicity which we are searching







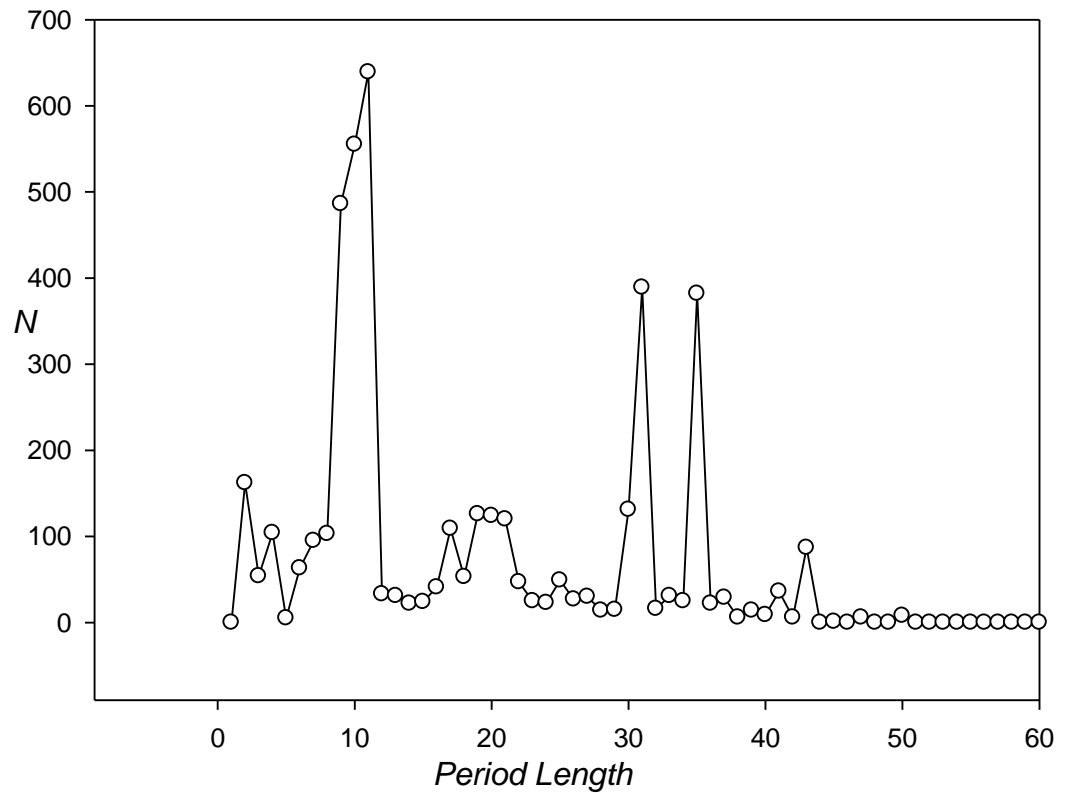
Test sequence  $(actgattcgagctagtaacgggct)_{30}$  with 1.5 random mutation per nucleotide.

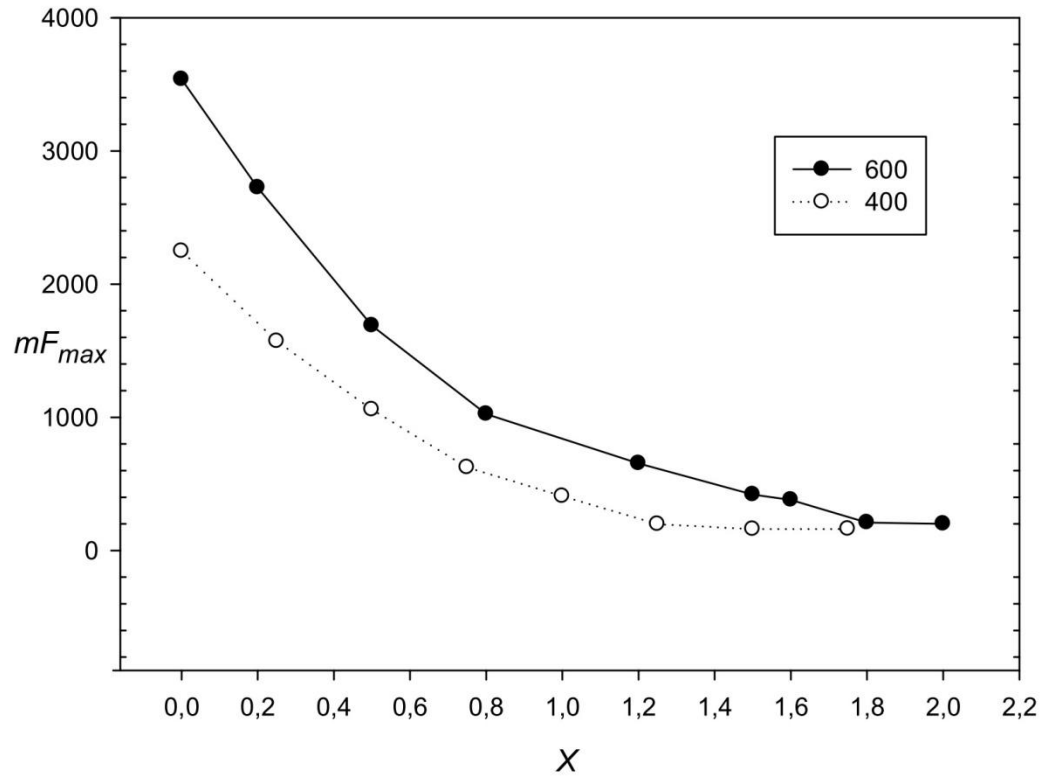


# Regions with periodicity in C.elegans genome

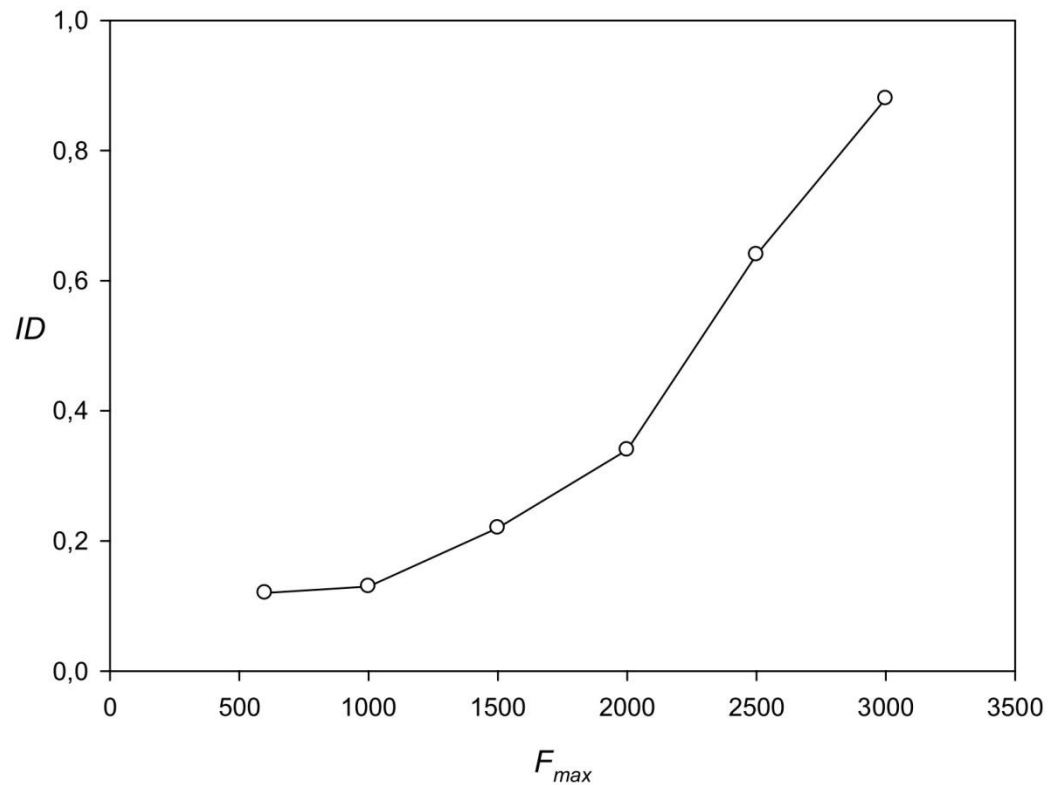
25360 regions having a periodicity with length of 2 to 50 bases

About 12% genome have periodicity of different length



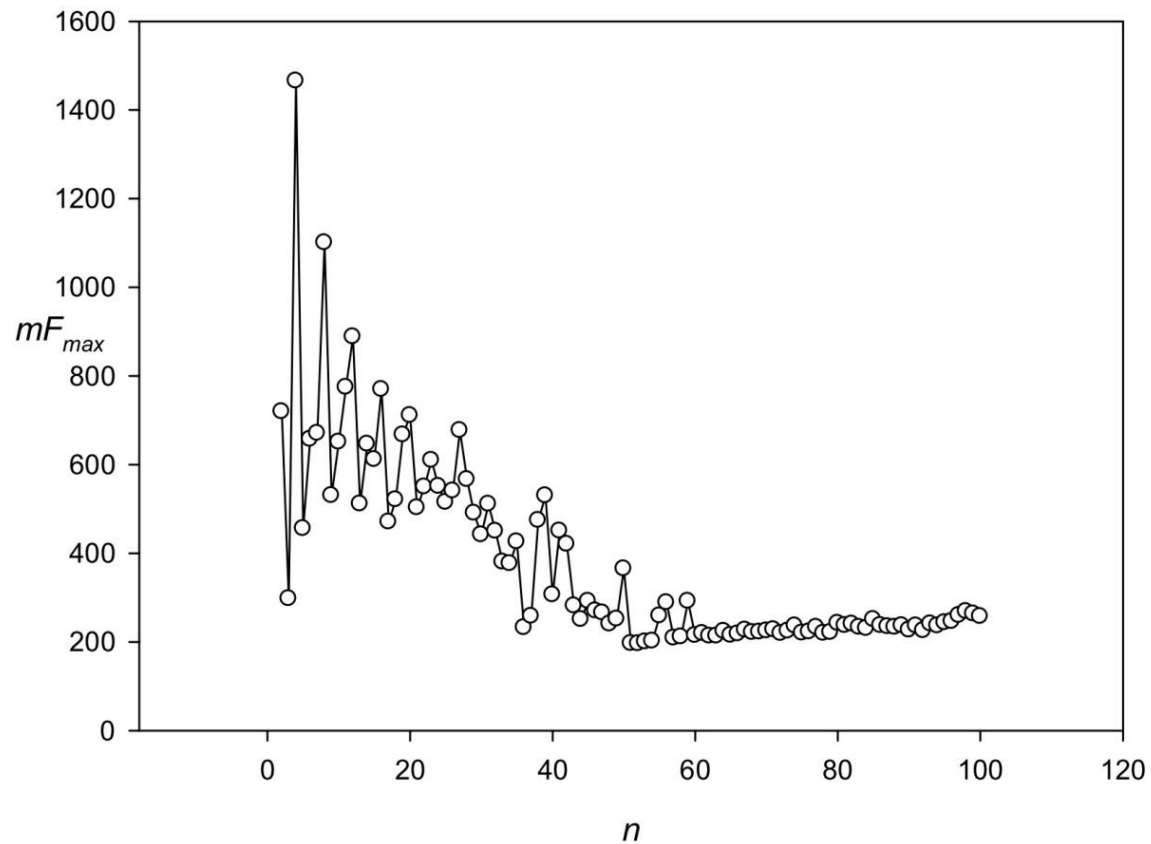


Influence of base changes on  $mF_{max}(20)$  for sequences 400 and 600 base pairs.  $X$  is the number of base changes per 1 nucleotide. The period length equals to 20 b.p.



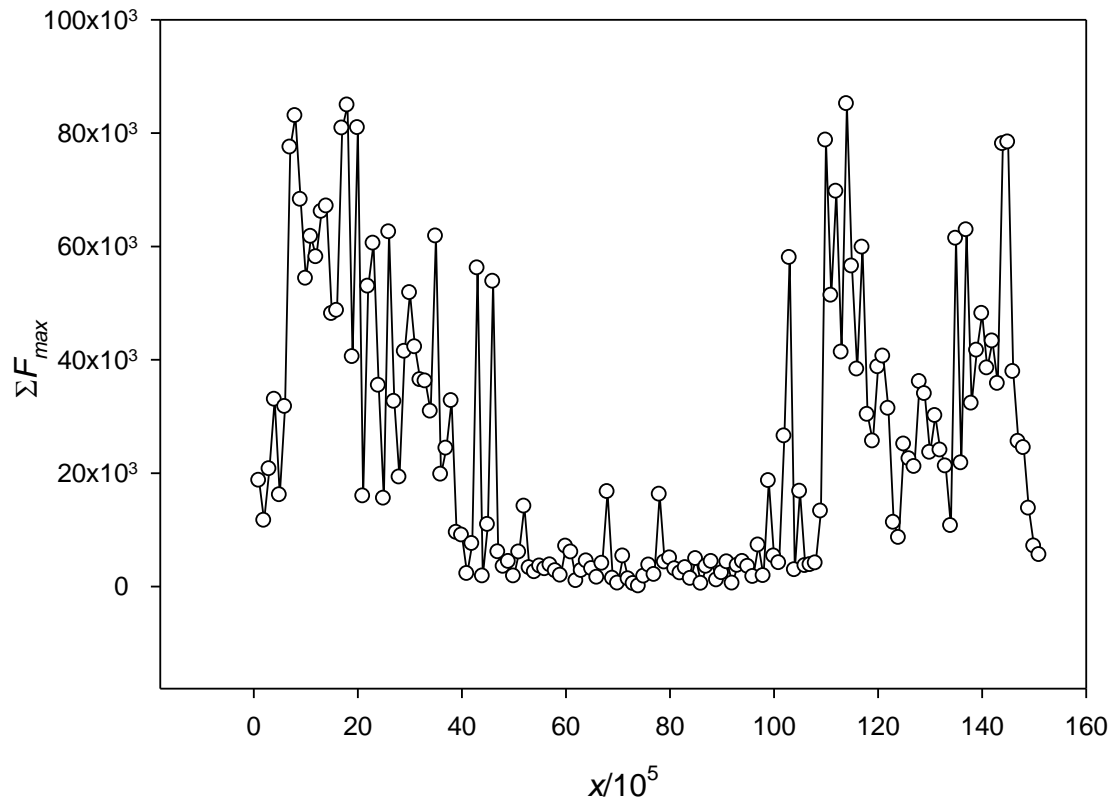
Comparison of developed algorithm with the program T-REKs (Jorda & Kajava 2009).  $ID$  shows the part of periodicities regions which can find the T-REKs. We can assume that the results are the same if the T-REKs detects at least 50% of the number of periods and the period length differs not by more than one base.

$mF_{max}(n)$  spectrum for fragment of the sequence NC\_003279.8 from chromosome 1 of the C.elegans genome. The coordinates of fragment are: 887101-887688.

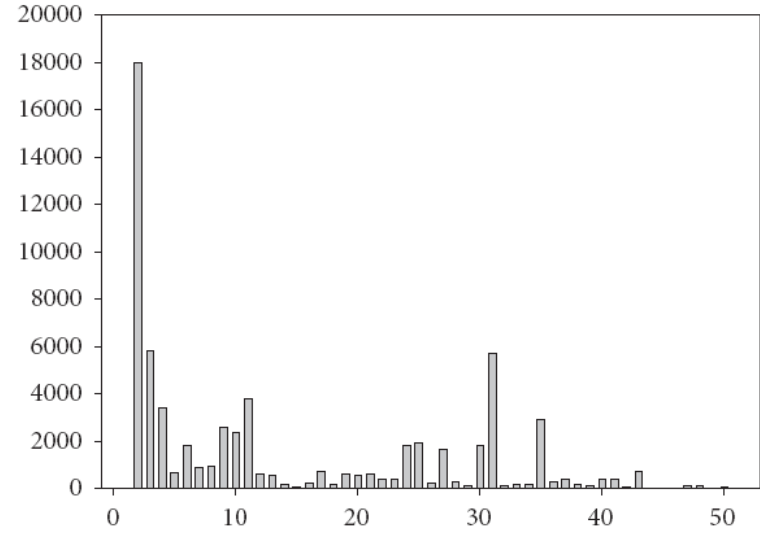
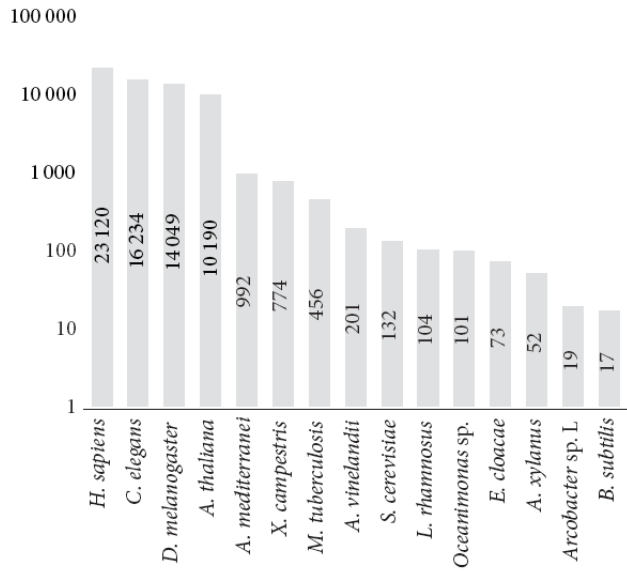




$\Sigma F_{max}$  for the chromosome 1 of the genome *C.elegans*.  $x$  is the position in the chromosome.  $\Sigma F_{max}$  is the sum of  $mF_{max}$  for all sequences with periodicity which are found in the region with coordinates form  $(x-1)/10^5$  to  $x/10^5$ .



# Database of periodic DNA regions in major genomes



<http://victoria.biengi.ac.ru/cgi-bin/indelper/index.cgi>



# Coding of the sequence $s(i)=x_2(i)-x_1(i)$ for Euro/US\$ rates

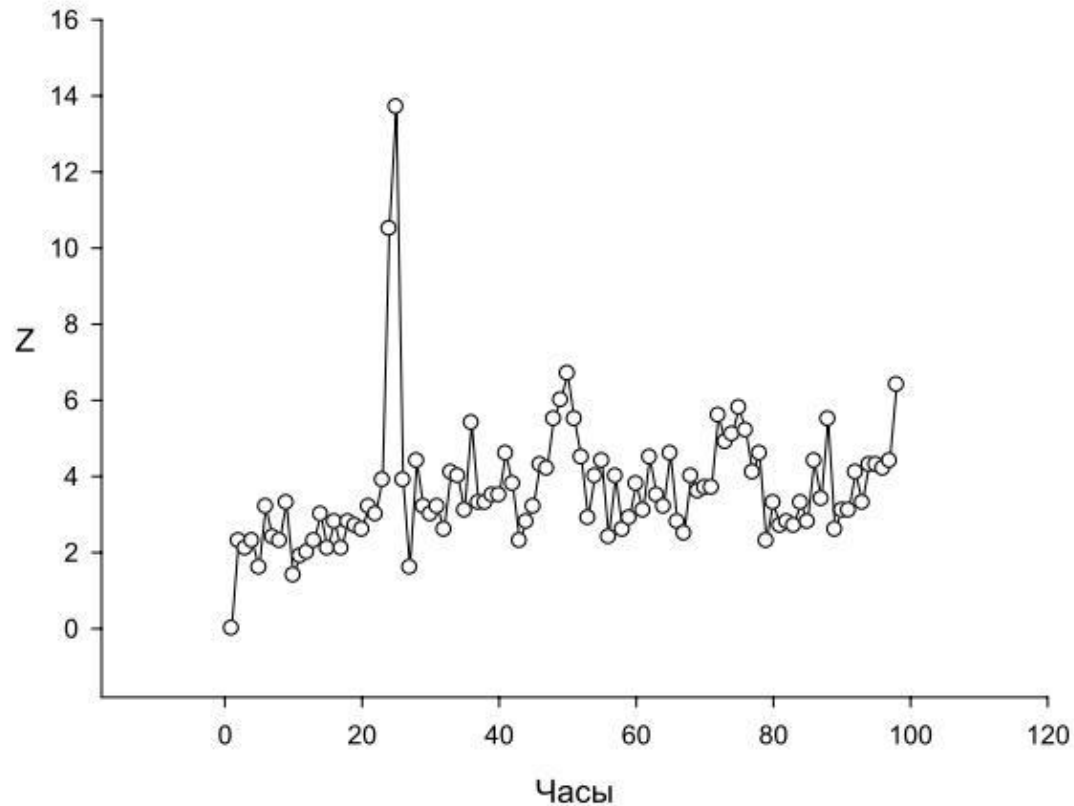


$x_1(i)$  -beginning of the candle  
(occurs at the beginning of each hour)

$x_2(i)$  -end of the candle (occurs at the end of each hour)

K	N	I	M	T	R	S	L	Y	F
-0.03700 -0.01140	-0.01140 -0.00780	-0.00780 -0.00550	-0.00550 -0.00390	-0.00390 0.00260	-0.00260 0.00160	-0.00160 0.00100	0.00100 0.00040	0.00040 0.00010	-0.00010 0.00020
C	W	P	H	Q	V	A	D	E	G
0.00020 0.00060	0.00060 0.00100	0.00100 0.00150	0.00150 0.00230	0.00230 0.00320	0.00320 0.00460	0.00460 0.00610	0.00610 0.00820	0.00820 0.01140	0.01140 5.00000

## $Z(n)$ for Eur/US\$ for candle equal to 1 hour

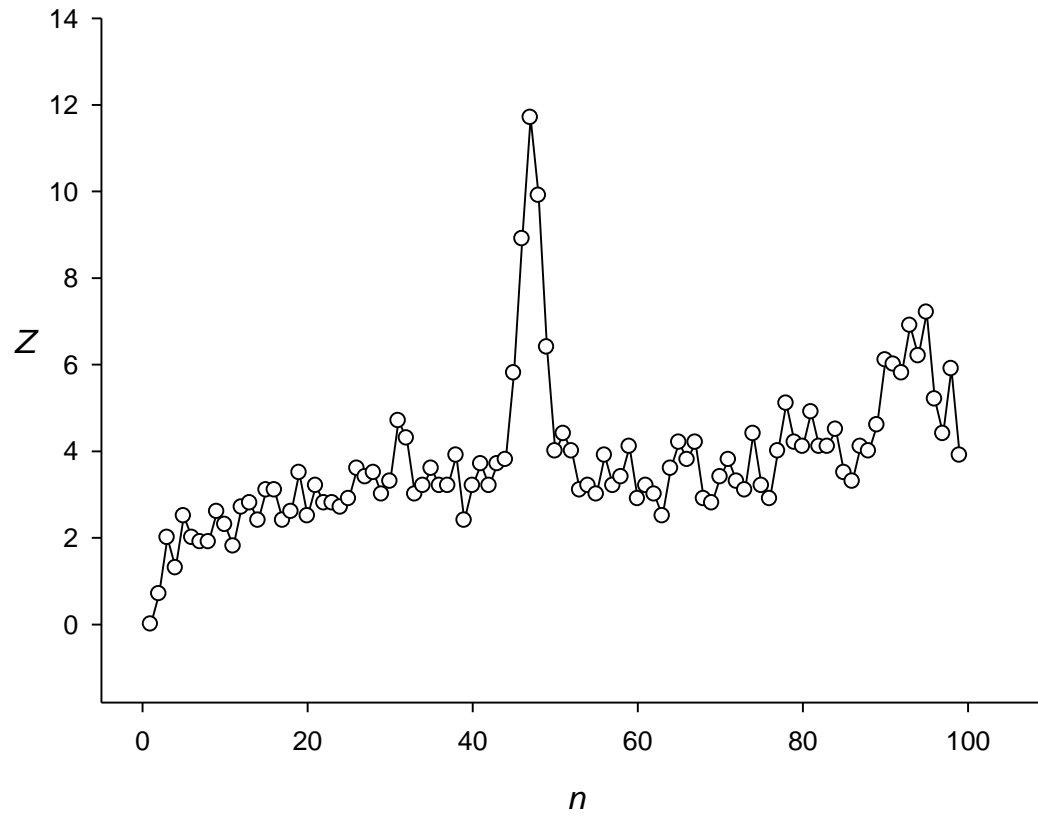


Period equals to 24 hours was found for last 20 days. We have same period for another intervals. Alignment contains 46 indels of different size.

Fragment of the PWM for n=24 hours from 1-12  
position of period

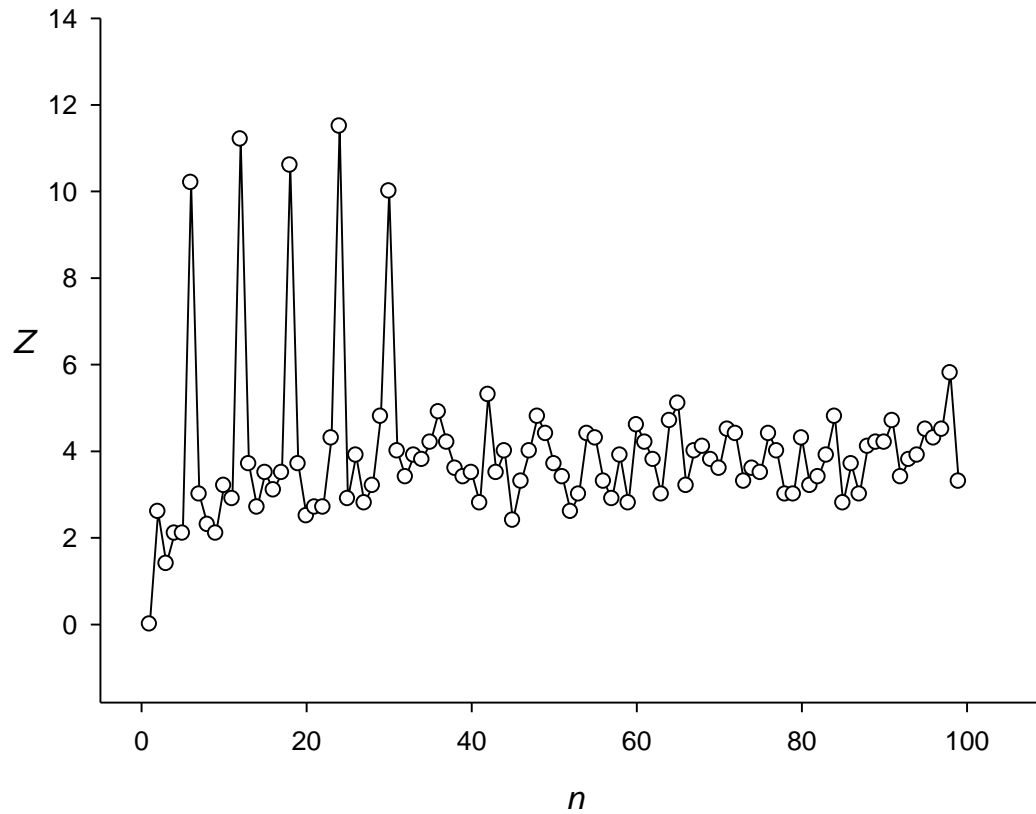
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>
<b>K</b>	2.2	2.5	1.9	0.1	1.9	2.4	2.4	2.5	-3.5	1.3	-0.5	-1.7
<b>N</b>	2.6	-1.5	1.4	-2.0	1.7	0.9	0.9	1.7	2.4	1.8	2.1	-3.8
<b>I</b>	1.6	0.7	-1.1	2.0	2.2	-1.1	0.7	1.3	2.1	-0.5	-1.1	0.1
<b>M</b>	-0.6	-0.6	-0.6	-2.4	0.0	0.6	1.2	1.2	-2.4	0.6	2.1	-0.6
<b>T</b>	-0.6	-1.8	1.3	0.7	-1.8	1.6	-3.6	-3.0	-2.4	0.7	2.0	-1.8
<b>R</b>	-2.9	0.9	-2.2	-1.6	0.9	1.5	0.2	-0.4	1.7	1.8	-1.6	-1.0
<b>S</b>	-1.8	1.8	-0.6	-2.4	-1.2	-3.0	-0.6	-2.4	-1.8	-1.8	0.6	1.2
<b>L</b>	-2.3	-2.3	-1.7	0.7	-1.1	-1.7	-0.5	-1.1	-3.5	-0.5	2.1	0.1
<b>Y</b>	0.4	-3.4	1.9	-1.5	-0.9	-3.4	-2.7	-3.4	-2.1	-2.1	0.4	-0.9
<b>F</b>	-2.2	1.0	-2.8	-1.5	-2.8	-2.2	-0.9	-2.2	-2.2	-3.4	-2.8	0.3
<b>C</b>	0.9	1.5	-0.3	-4.1	-1.6	-3.4	-2.8	-2.2	-1.0	-0.3	1.5	-3.4
<b>W</b>	-3.4	-1.6	-0.3	-1.6	-1.6	-1.6	-2.8	-2.8	-2.8	-1.6	0.3	-1.0
<b>P</b>	-2.9	-1.6	-2.3	0.2	-1.6	-0.4	-4.1	-1.0	1.4	0.2	-2.3	-2.3
<b>H</b>	0.3	-0.3	-3.4	1.9	-1.0	-1.6	-3.4	1.5	-1.6	-0.3	-2.8	-2.2
<b>Q</b>	-1.2	-1.2	1.8	-0.6	-3.6	0.6	-4.2	-1.8	2.0	-3.6	-3.0	2.1
<b>V</b>	-2.3	-0.5	1.6	-1.7	-2.3	-2.3	-1.7	-1.7	0.1	-1.7	-4.1	2.2
<b>A</b>	1.3	0.7	-1.8	1.6	1.9	-1.2	1.8	1.8	-2.4	-1.2	-2.4	2.3
<b>D</b>	0.4	-0.1	0.4	1.0	-1.9	2.1	-0.1	-1.9	2.0	-1.3	1.0	1.0
<b>E</b>	-1.6	0.2	1.5	-0.4	0.9	-1.0	2.5	2.2	-2.9	2.5	-2.9	-0.4
<b>G</b>	0.7	-1.7	0.7	2.5	2.2	2.2	2.2	-0.5	2.2	1.6	1.4	-1.1

## $Z(n)$ for Eur/US\$ for candle equal to 30 minutes



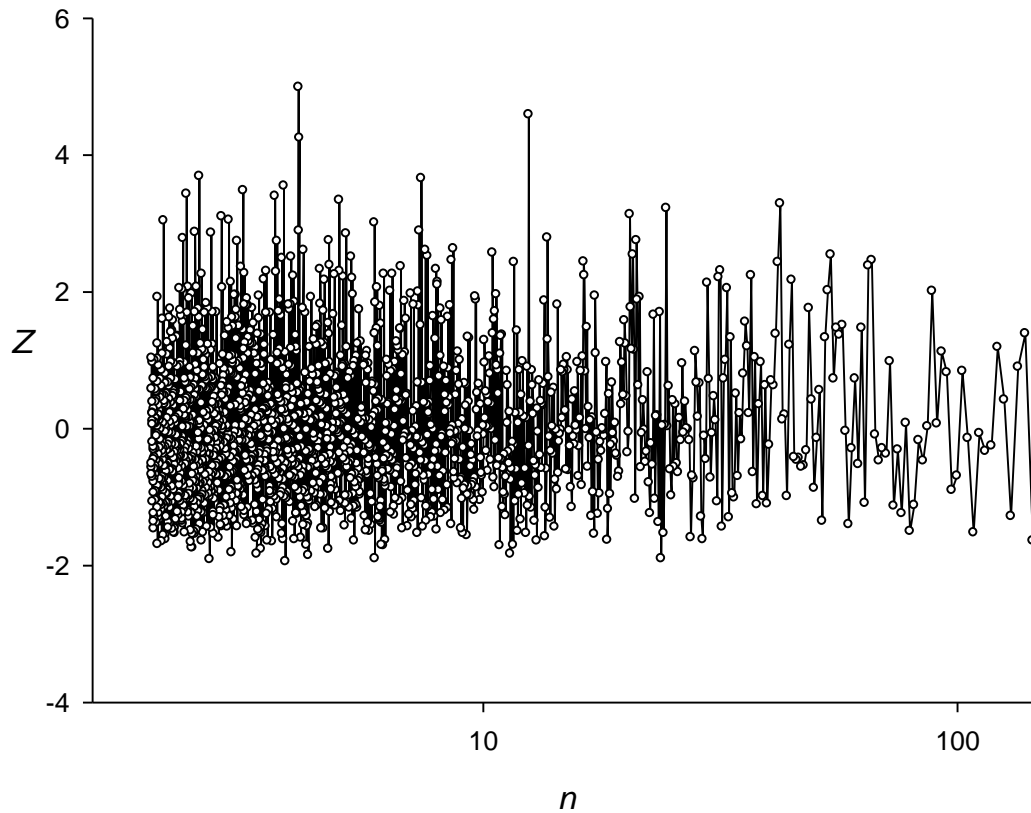
$n$  – length of period

## $Z(n)$ for Eur/US\$ for candle equal to 4 hours



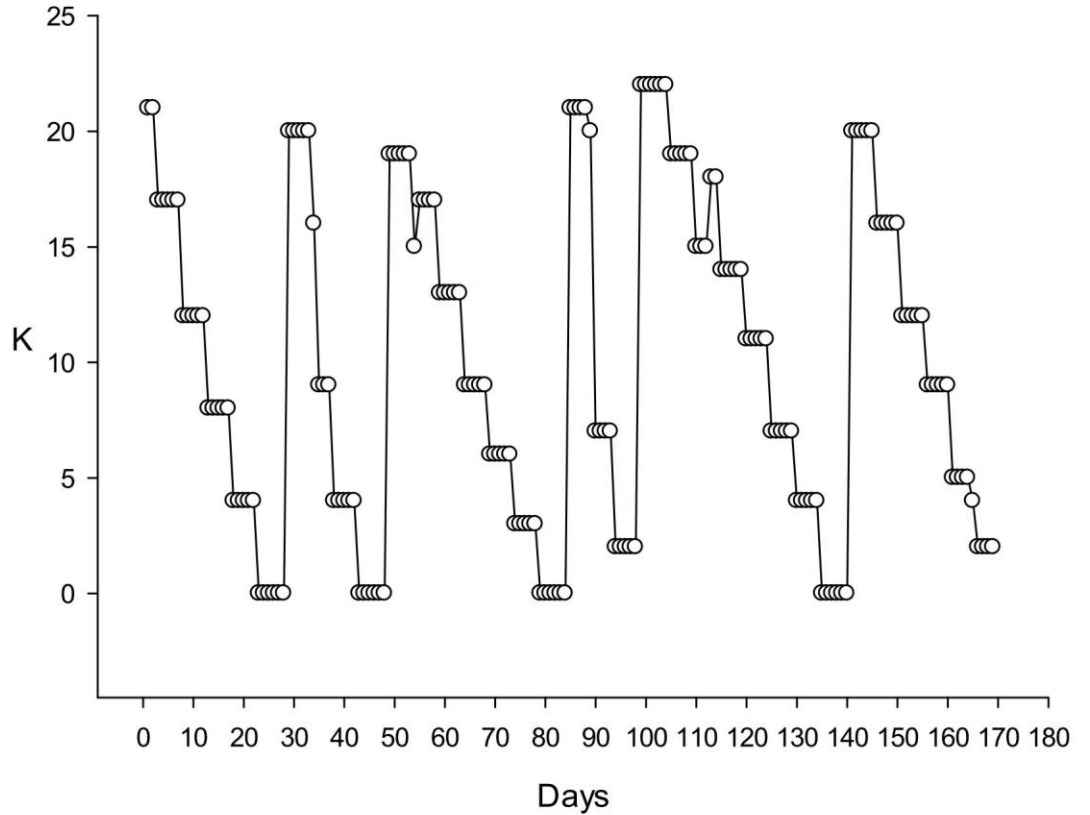
$n$  – length of period

# Fourier transform for Eur/US\$ for candle equal to 1 hour

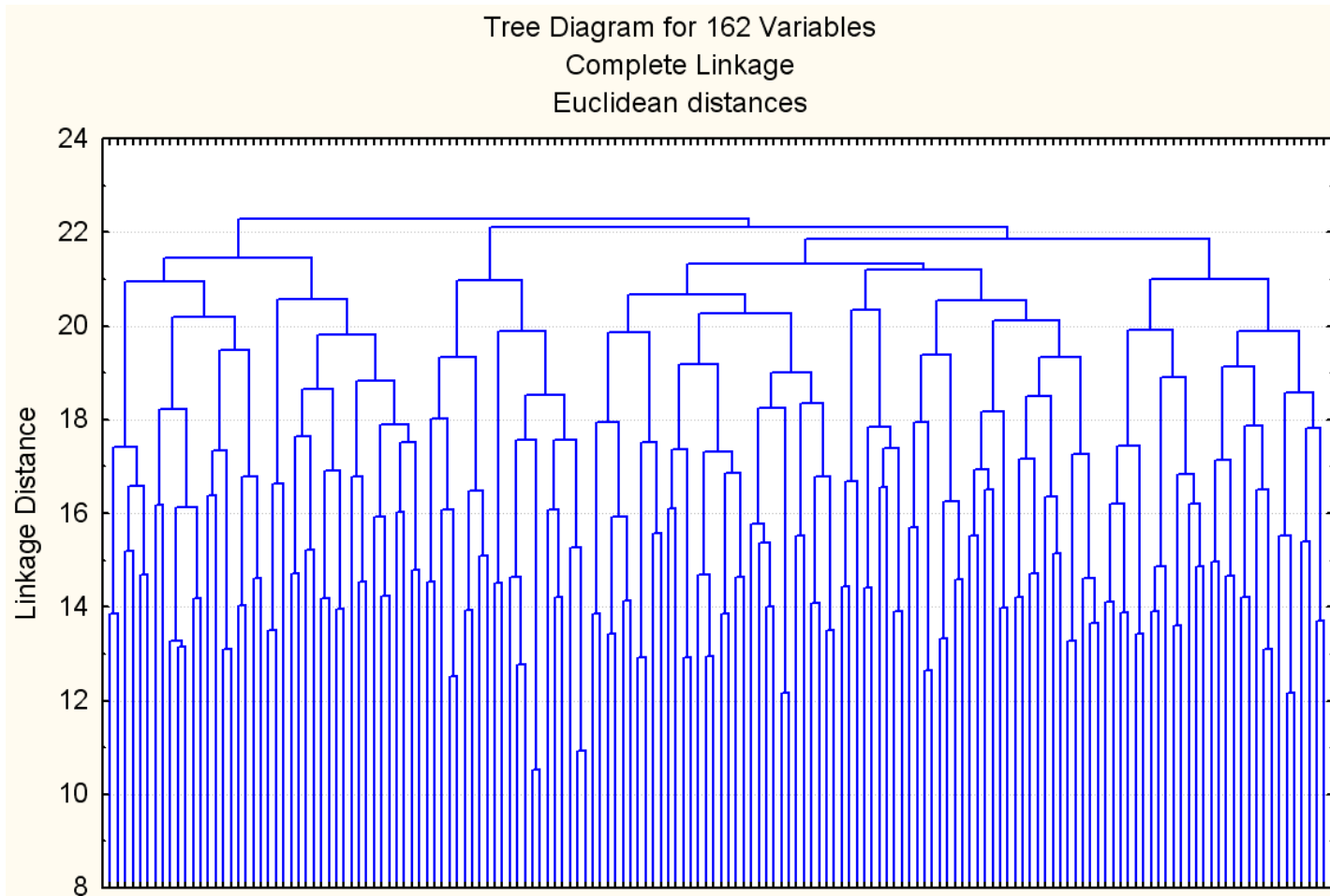


# Phase shift of period equal to 24 hours

**K- position of PWM which corresponds to the first candle of each day (0.00-1.00 hour)**



# Classification periods equal to 24 hours

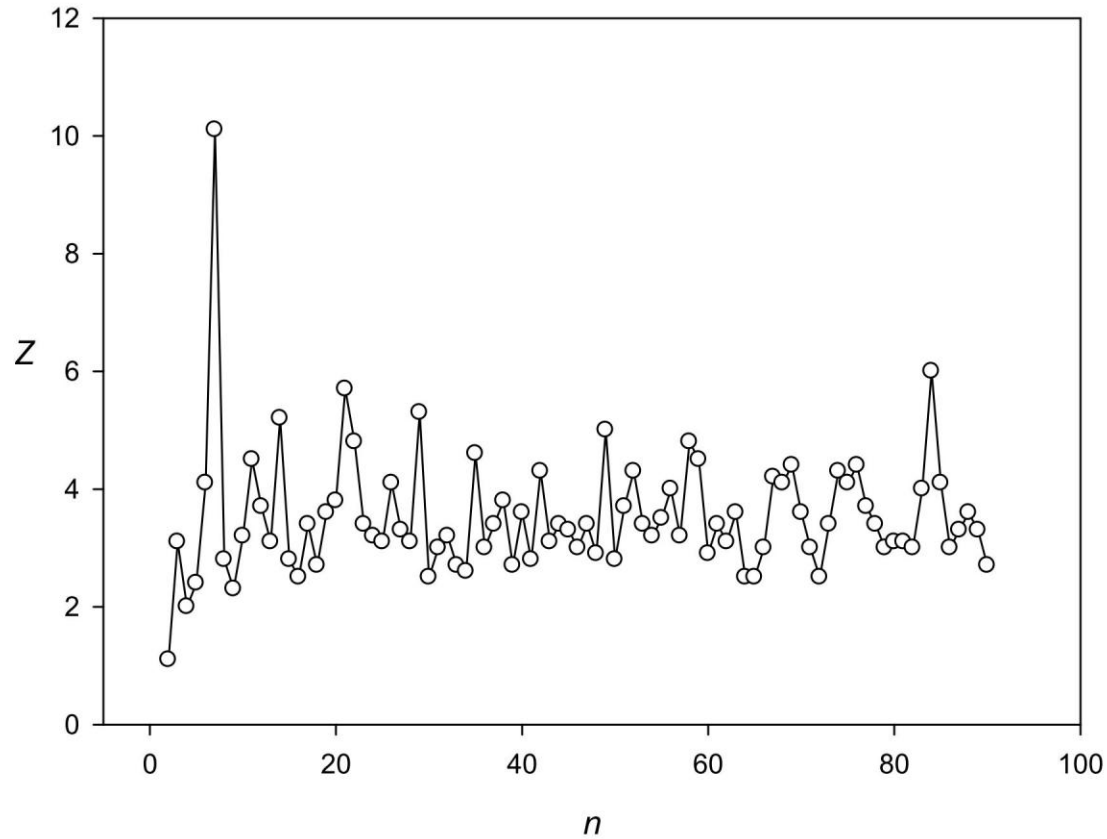




# Multiple alignment periods equal to 24 hours

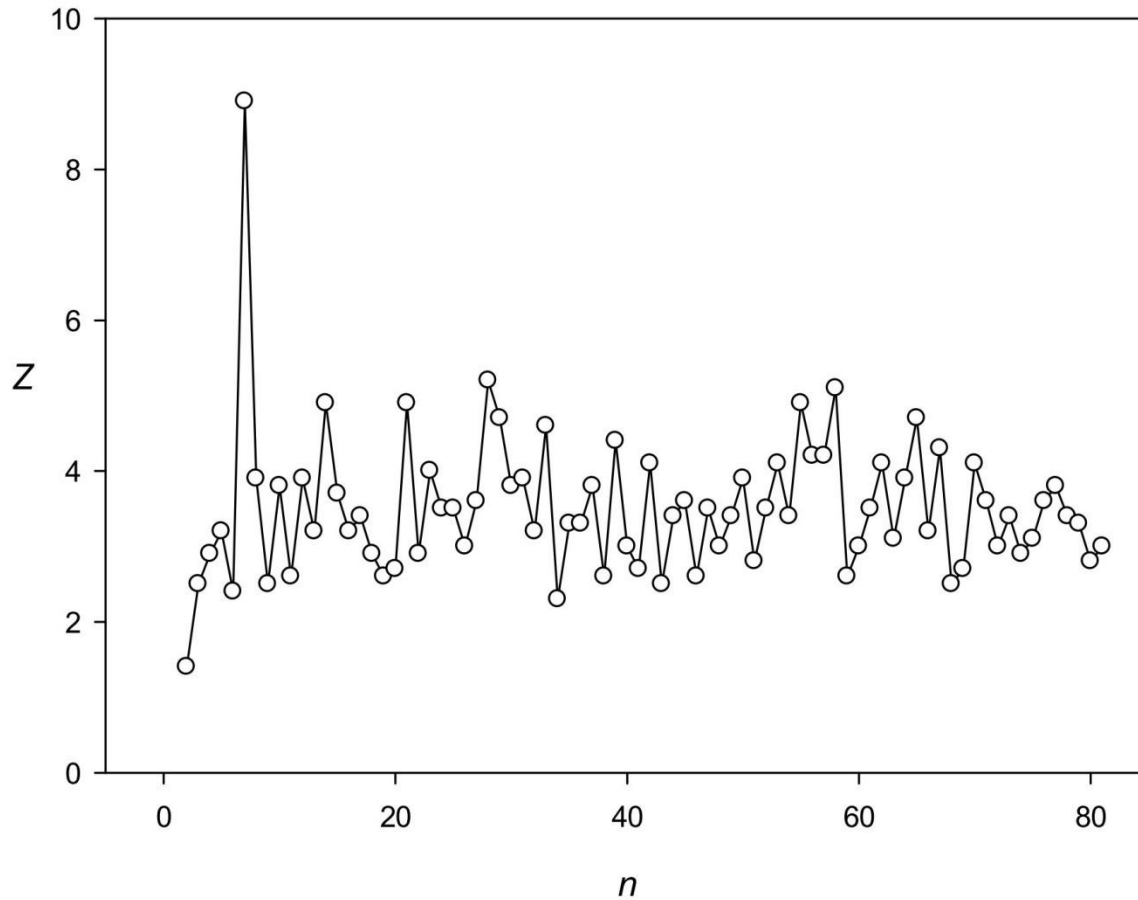
•	J=	1	WMLPA . . . . . VK . IA . MMKNAENWD . AWL . MRG
•	J=	2	PFQTQ . . . . . KR . GN . RNGLKLMGD . VWT . FQI
•	J=	3	PRRQW . . . . . KQ . MD . NWKWKMKND . RRM . MNG
•	J=	4	VTLVE . . . . . AE . LM . WALVVRMAW . FKA . NFD
•	J=	5	RMTAQ . . . . . MN . LP . KGDIEGNFE . PFL . RGL
•	J=	6	FWNIQ . . . . . NM . AR . DGAQKAMW . LLT . FIP
•	J=	7	VRWFA . . . . . GG . ET . NGGRGWMGV . RFT . MLI
•	J=	8	WWWAN . . . . . AK . WA . GKMIPIVIP . WPQ . DMW
•	J=	9	QTFWM . . . . . DE . RK . INFPEALIM . EFF . WML
•	J=	10	QAPLD . . . . . KF . NK . VWGANRPID . VPK . PQV
•	J=	11	LMLEW . . . . . FG . MV . KKQKGTRL . MDF . RWN
•	J=	12	VTFTE . . . . . KT . KG . NKFNKVRP . RDP . FTQ
•	J=	13	IPTTN . . . . . EE . KI . GKNKTGVLK . WWQ . FFL
•	J=	14	FVMPE . . . . . KG . EP . MKPNGVMWG . RTW . TPN
•	J=	15	MMWEE . . . . . IK . AD . MDTIFRGFT . LVT . VWA
•	J=	16	MIWFI . . . . . VE . TI . DVNVGKGQL . QMV . LCL
•	J=	17	QKCIA . . . . . AQ . LP . NKGEMQRII . VRP . PRQ
•	J=	18	LRLMR . . . . . WK . VM . VRGIGMDFE . MQA . TVQ
•	J=	19	FRRRV . . . . . GD . KL . LDRLGTLFF . ALF . PGI
•	J=	20	PWWTT . . . . . PE . TW . IGKLGKLN . PSS . LFF
•	J=	21	IFLAE . . . . . KQ . DM . KFQMGVIEF . LFF . I . E
•	J=	22	DEPTW . . . . . EL . VG . NIDFFTQPV . LDL . FSQ
•	J=	23	PTLVL . . . . . LG . NI . QGSKIFPAQ . QVW . NQP
•	J=	24	GNTLT . . . . . DN . SF . EMTGPSQND . FFQ . TLT
•	J=	25	WEWWT . . . . . SD . AQ . WFPGGNDTW . WPF . PLA
•	J=	26	PTIAF . . . . . MQ . TI . TWMQELNI . SNF . FTE
•	J=	27	ADAQS . . . . . TK . KV . NQDEWTKED . GWS . PFT
•	J=	28	WPTFS . . . . . VM . TD . NSAAGSKGD . QWP . ALQ
•	J=	29	NAFPF . . . . . VM . GE . TKNTNGGIP . TSW . LWT
•	J=	30	WPFFF . . . . . EA . PN . ENDAPKGQN . WTD . QFF
•	J=	31	LLLLS . . . . . AL . QS . FSNNKKVWV . EQL . STD
•	J=	32	DMSSQ . . . . . LN . KN . QDMGGTGNR . PSW . FLQ

The spectrum of  $Z(n)$  obtained for the stock of Bank of America.



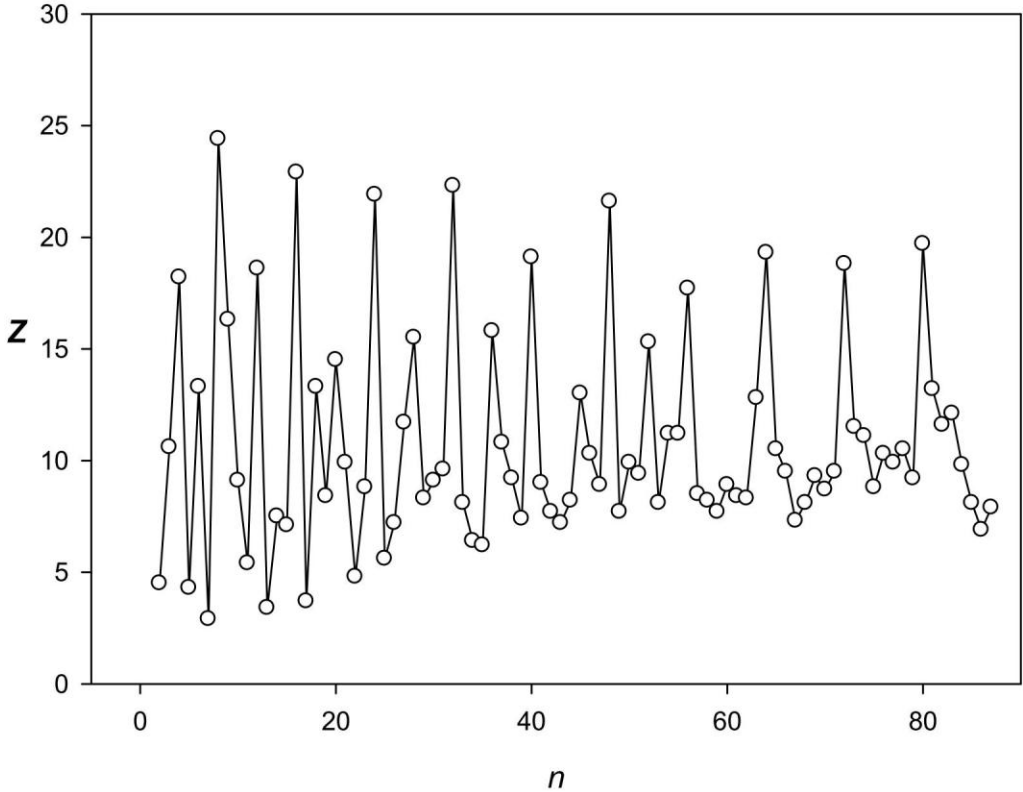
**Trading Time - 7  
hours every day**

# The spectrum of $Z(n)$ obtained for the stock of Microsoft corp.



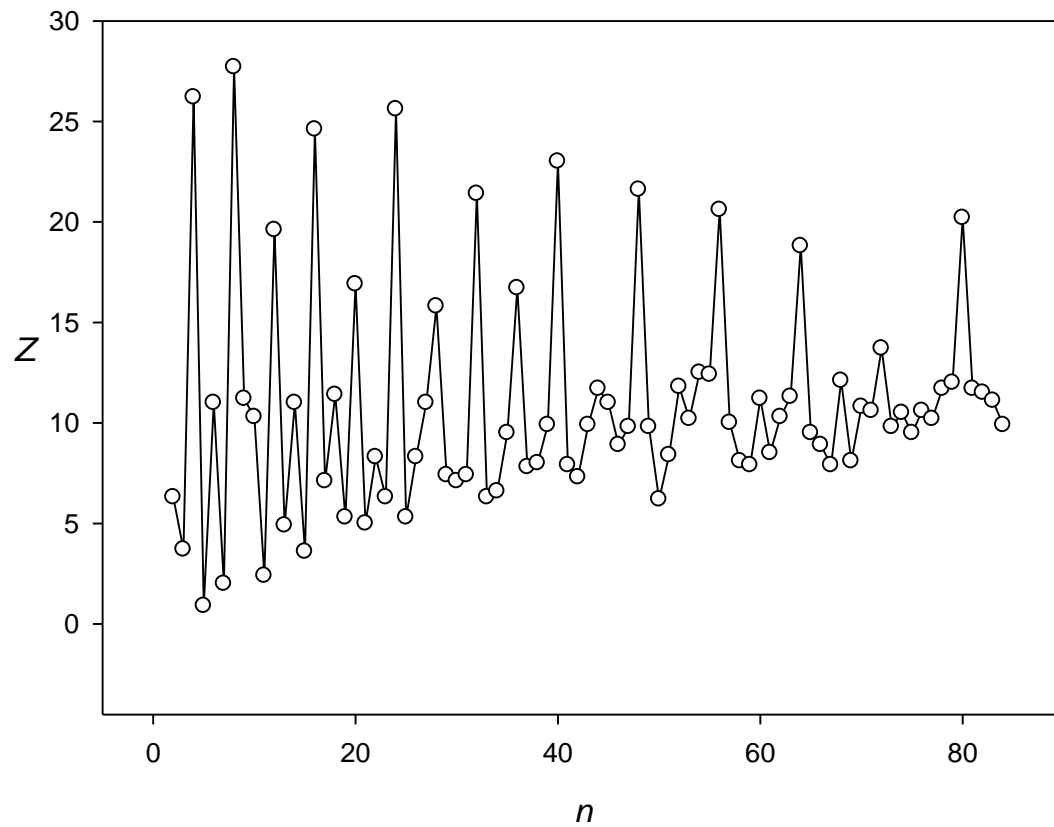
Trading Time - 7  
hours every day

**The spectrum of  $Z(n)$  obtained for the S&P500.**



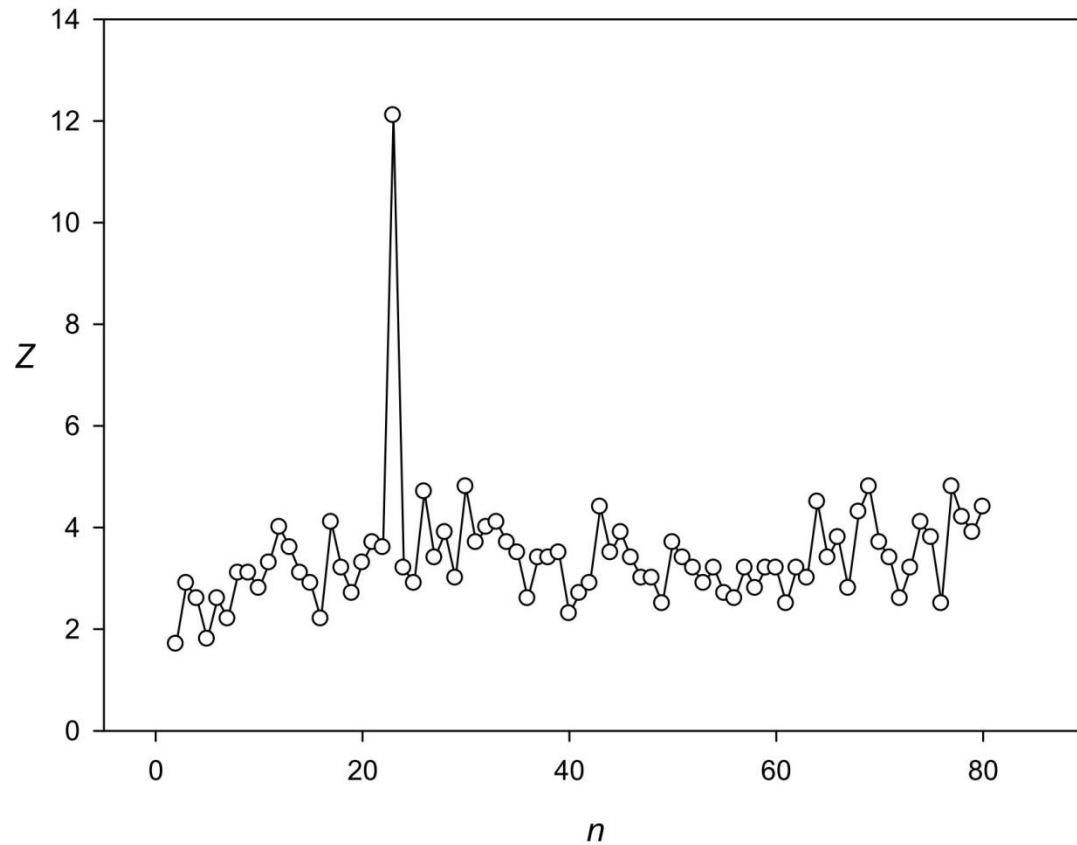
Trading Time - 8  
hours every day

# The spectrum of $Z(n)$ obtained for the NASDAQ.



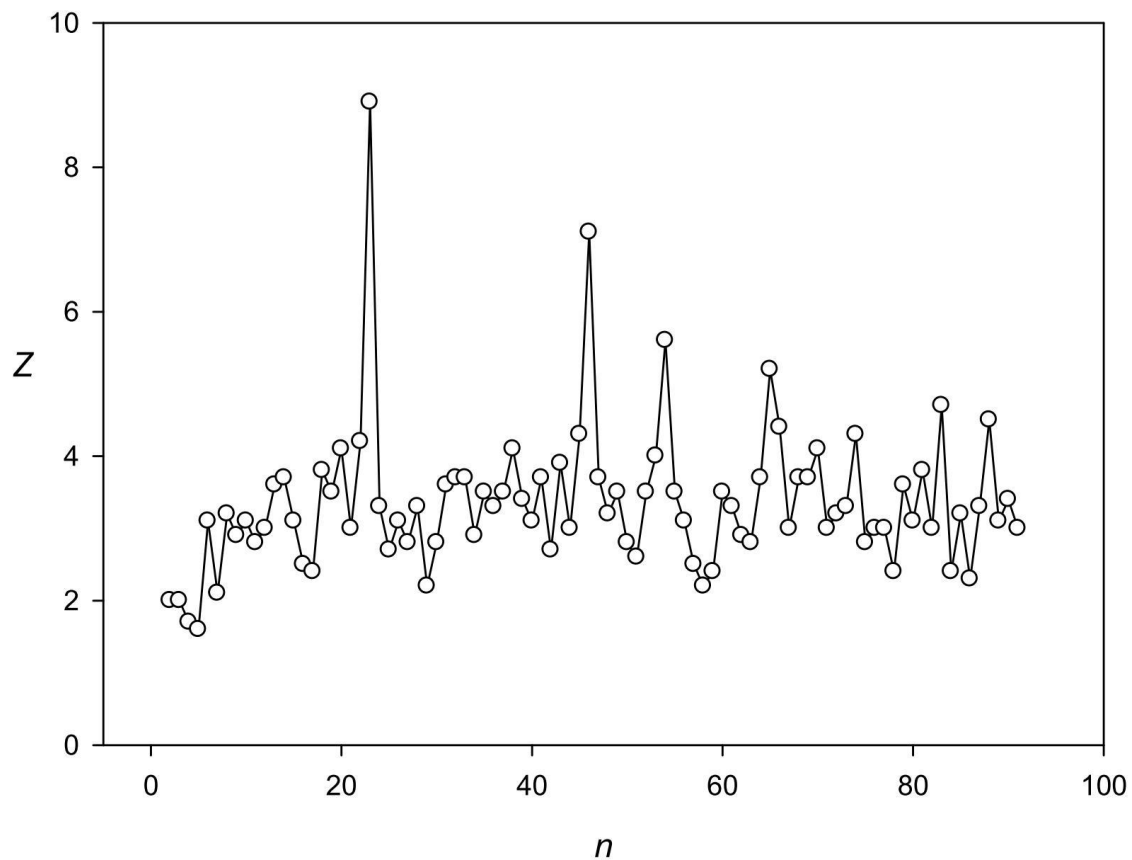
Trading Time - 8  
hours every day

# The spectrum of $Z(n)$ obtained for the Gold price.



Trading Time - 23 hours  
every day

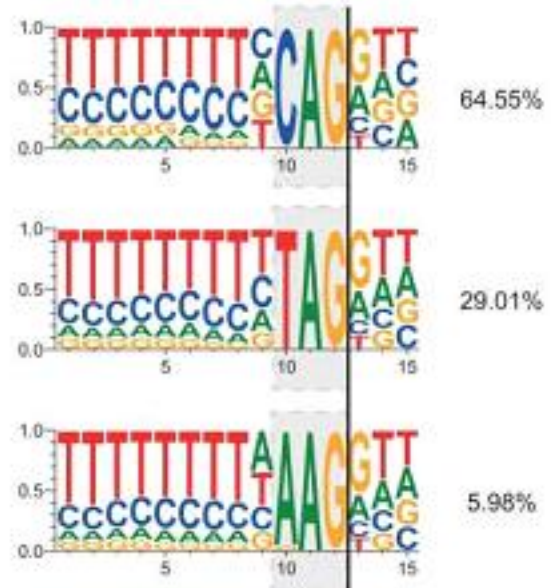
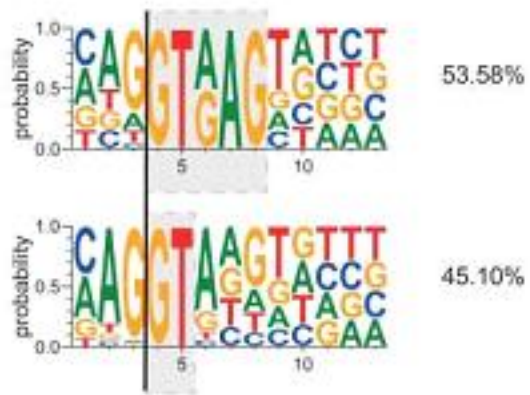
The spectrum of  $Z(n)$  obtained for the Silver price.



Trading Time - 23 hours  
every day

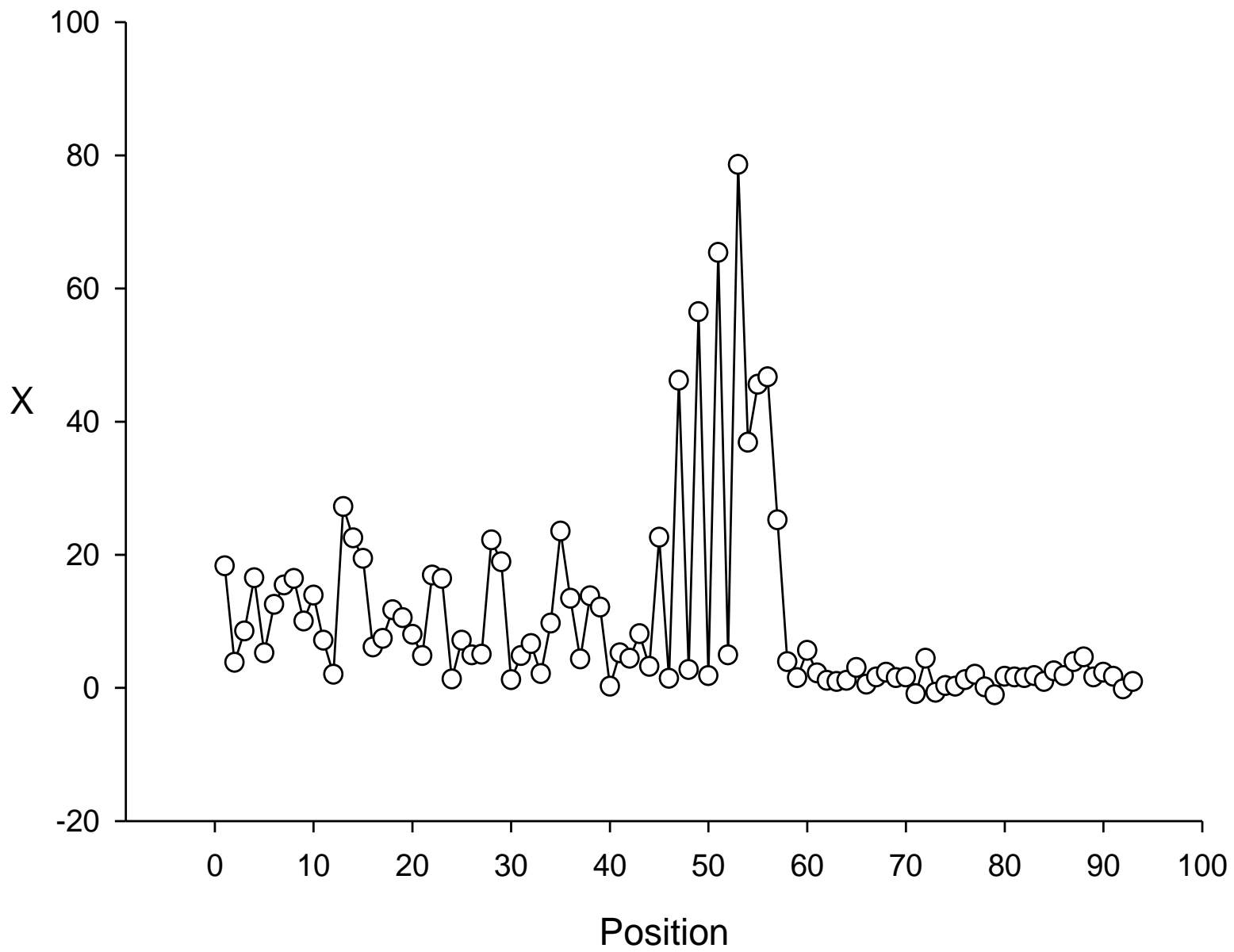


**A. GT-AG splice sites**

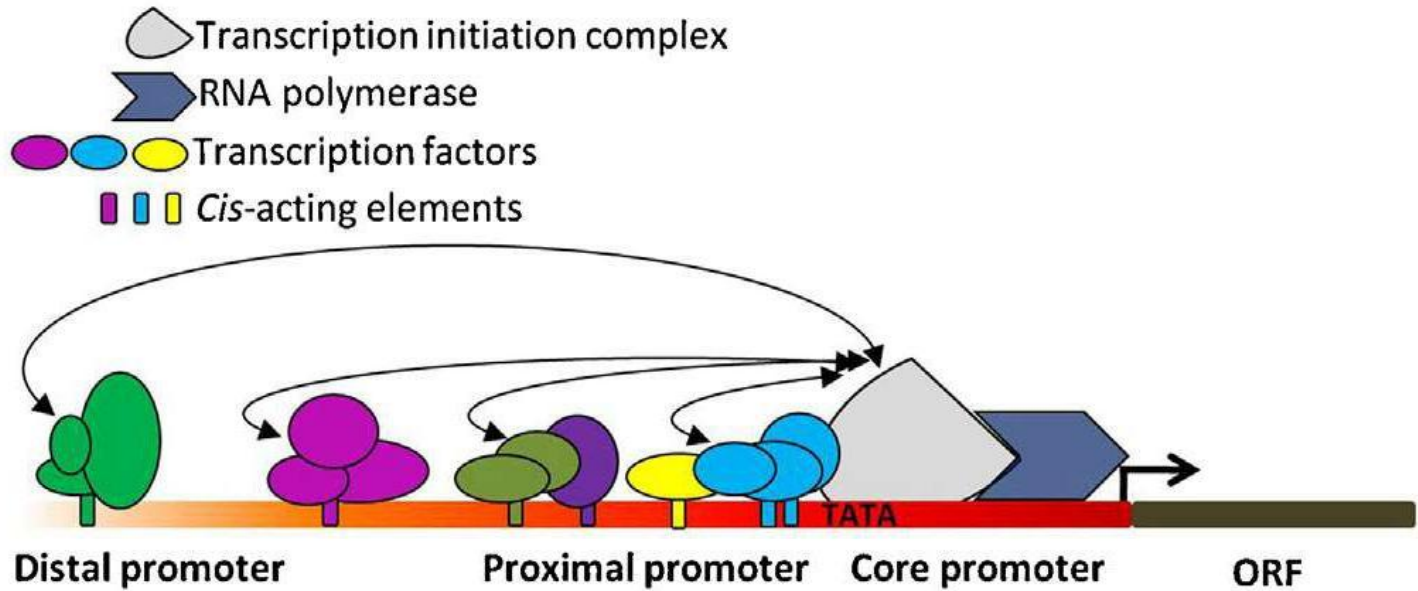




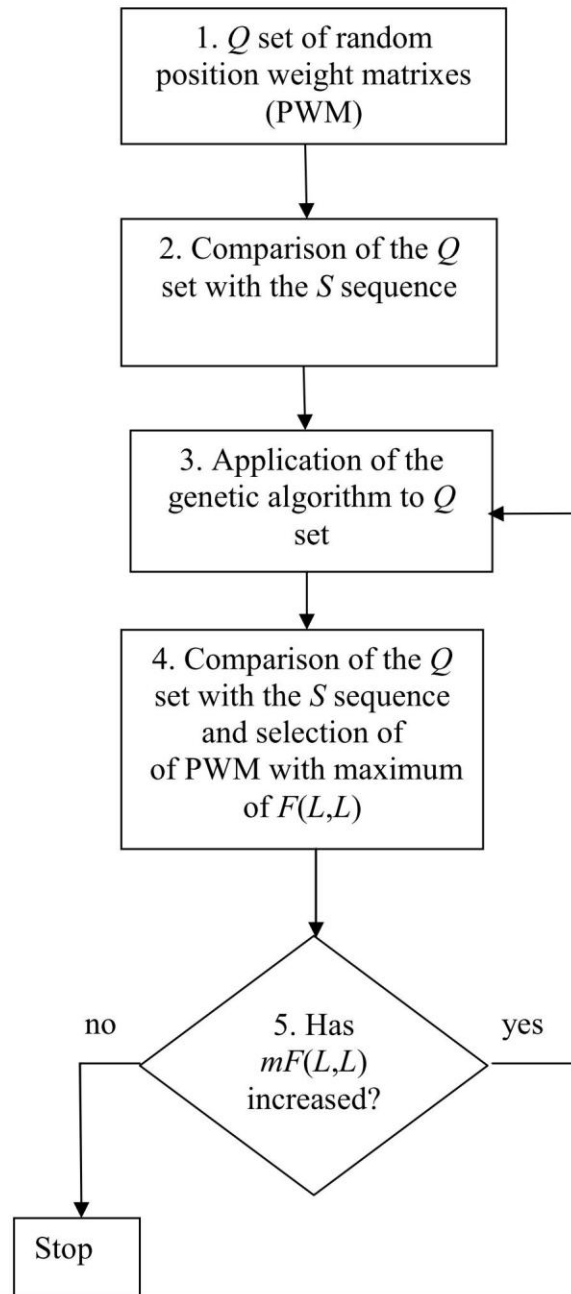
C.....TCCCAT.T.GG.T.G.G.C...A.G.C.C.AG.T.G.C.C.A..C.C...A.T.G.C..GC..G.C.T..C...A...GT..G..TAAGTATCATTCCCTCTCACTGTCTGGAGAGGACGAGAATT  
.....TGTCGGTA.T.CA.T.A.T.T...A.G.G.C.GC.T.G.T.A.T..G..A...C.A.A.T..CT..C.C.A..T...T...C...G..TAAGTACCTCTTGGTCATTGGACACATTGTAGATTAGTCCCT  
.....GCCGAGAA.C.TC.C.G.C.T...C.G.T.T.CT.G.T.G.C.G..T..T...C.T.C.C..TG..T.C.C..C...A...G...G..TAGGGAAGAGGGGCTGCCGGGGCGCTCTGCCGCCGTTTCT  
.....CAGCCAGCTG.A.AT.T.C.A.T...T.G.A.G.GA.C.C.A.G.....T...C.A.A.T..TC..T.A.A..A...C...T...G..TAAGCAATGTCACCTCAAAGATGCTTTTAGACTCTTCTCATAT  
.....CGTGATTGTC.G.GG.G.G.G.A...A.A.G.A.CA.T.C.C.A.G..G..G...C.T.C.C..T...T.G..C...A...G...G..TAACACATCTGTTTGGATAAATGGGTTCAAGGAGGACATT  
.....AGAGAATC.A.GA.G.A.C.A...G.C.C.T.TT.C.C.C.A.A..G..A...G.A.T.G..TT..G.G.C..A...A...G...G..TAAGTCAGACAAACAGCAAAATGACAAAACATGTTTTATGT  
.....GACACGGAC.C.TT.G.G.A.C...A.G.A.G.AA.G.T.A.A.T..C..G...C.T.G.T..TT..A.T.C..A...A...G...G...TGAGTACCCCTCTCCACCCCTGTGCGCAGAAATGTT  
.....GGCCCTTAC.A.TC.G.G.C.C...A.C.G.T.CA.T.G.G.A.A..G..G...G.T...C..AC..C.C.A..C...A...G...G..TATGTACATTGGCTTACCTTTAGCGTAATGGCTTGGAAAGT  
.....CATTGTC.ACTG.T.T.G.T...G.T.C.A.CC.T.G.C.G.C..T..G...C.T.G.G..AC..C.G.A..G...A...G...G..TGAGCTGAAAAGAATACCCTTTCTTTTTACAGAGAATAGAT  
.....TGACAAAA.A.TG.A.T.C.A...C.T.C.A.CC.A.A.A.A.T..T..C...A.C.C.A..AG..A.A.A..G...A...G...G..TAACCCCTGTGCCAAACACCAACCACCACTGTGGTCACAGT  
.....CAACCCCT.A.AA.C.C.A.A...C.G.A.A.GG.G.A.T.G.C..T..T...T.C.T.G..TT..G.T.C..A...A...A...G...G..TAAGGGTGTCTCAATTTGCCTTTCTCCCTCATGCGACACT  
.....ATGCCA.C.GG.A.C.C.C...C.G.A.C.TC.C.C.T.G.C.A..G...C.A.T.C..AA..A.C.C..ATC.A...G...G..TGGTGTAGTGGCTCCGGACACACAGCGGGAGTGGGCGGT  
.....AACAGCT.G.TA.C.A.C.T...G.C.T.C.TC.T.T.C.C.T..G..G...C.A.A.T..TG..A.C.A..G...T...G...G..TGAGTACTTGACAAAGACCATCAAGGGTATACCTTTCTGCTT  
.....GCTGTTC.T.GA.T.A.A.A...GTC.T.G.GA.A.G.A.T.C..T..C...C.A.A.G..AT..C.A.A..C...A...G...G...TAAGTCTGGCTAAAGCATTTCTACCTATTCTCCAGCTTTCT  
.....CTGCCCTCA.G.TC.C.T.G.C...T.C.C.T.CA.G.C.C.T.C..T..T...C.A.G.C..TT..A.G.C..T...T...G...G...TAAGTTGACCTAATCCAAATGCATGCACAAACAGGAATTTGT  
.....CTGTCC.C.CG.G.A.T.G...G.C.C.G.CC.T.C.G.C.T..A..T...C.A.T.CGTCC..A.T.G..A...A...G...G..TACTTGCCCTTGATGTGGGCGGGTACTGAGGGGACCACTGAT  
.....TTCATGAGC.T.GC.C.A.G.A...G.T.T.T.AA.C.C.T.C.A..T..C...C.A...C..TC..C.C.T..G...A...G...G..TAAGCCTTTGCTCGCAGTGGGGTGTGGTTTTATGCACACTC  
.....ATCAATAGC.T.CG.G.G.G.G...A.C.T.G.GA.T.T.G.C.T..T..T...T.G.G.C..TG..T..T..C...A...G...G..TTTGTCCCCCGCTGGGTGGTAGAGATGGACTGCCCATTAGT  
.....CGACGGAC.A.GT.A.C.A.T...C.G.T.G.AC.T.G.G.C.G..G..G...G.A.C.G..AC..G.G.C..A...A...G...G..TAGGCTCCTGTCCCGTCCCGTTGGCCTCTGTGCCTGGGGGT  
.....GTTCGAG.A.TC.T.C.T.T...G.C.A.A.CC.T.G.T.C.T..TT.G...G.A.C.G..CC..A.T.G..G...A...G...G..TGAGCCCGCAGCGCGGGCCGGATGGATGTTGCTTCCAATGT  
.....AGAAGAGT.A.TC.C.A.G.C...G.G.C.A.CC.T.G.G.A.C..G..A...C.C.T.G..TC..G.A.A..A...C...T...G..TACGTGTGGGTGAGGGCTGGGGTGGGGTGGGGTGCACGGT  
.....CTGCCATCT.A.CC.A.G.G.A...G.G.C.T.TT.G.G.A.G.C..G..G...C.T.A.A..AA..A...C...C...T...G..TGCATTTGAGTCCAAGGCTTCCGCTTTCCATGTGGAGT  
.....CCGGTGATC.T.AC.A.C.A...T.G.G.A.G.A.C.A.A.G.C..C...C.A.T.C..G..T.C.A..C...C...T...G..TGTGTCCCGGAGCGGGAGGGCGCGCGGGTGGGGCGGT  
.....TCTCTCAGC.A.GC.T.T.C.C...A.A.G.A.GA.A.C.C.A.A..T..G...G.A.A.T..GG..A.G...A...A...G...G..TATGAGTTGTGTGTTTTTTGGCTGTATCCCTCTTTCATTTT  
.....CACGTTTC.C.CT.T.C.C.T...C.C.A.T.AT.T.T.A.C.T..G..G...A.C.A.G..AG..T.G.C..T...G...T...G..TGAGTACCAGCAGGGGAAGAGGATGCGCGGTGGGATGGGGTT  
.....GCCCTCTT.C.TC.A.G.T...A.T.G.T.GG.T.C.A.A.C..T..T...C.G.G.G..GC..C.G.G..G...G...G..TGAGTAGTGGCACTTCAAGTAAACGATGCTTTTCTTAGTGTGT  
.....GCTGCTGC.A.CA.A.C.T.G...C.G.T.G.CA.G.G.A.A.C..T..G...A.A.G.C..TC..T.T.A..G...A...A...G..TGAGTCCGCTGTTTTCCCAATTTACCCGTTTTCTTAGTGTTT  
.....CAGGCGGAG..A.GC.A.G.G.G...C.T.G.T.GG.T.C.T.A.C..T..C...C.A.A.C..AC...G.A..G...A...G...G..TGAGGTGCCGCCACCCAGGCCAAGGAGATGCCACCTTT  
.....ACAGAGAAT.G.GC.C.G.C.T...G.C.A.T.GG.T.T.A.C.C..A...C.A.G.A..CC..A.C.A..A...A...G...G...TGAGTAGGGACAGTGGAGGAGCTTAGCTTGTGGGGCTCCGT  
.....AAAGCCGT.A.TA.C.T.T.A...T.G.A.A.TT.T.A.A.A.G..T..G...G.A.A.A..AT..T.T.T..T...T...T...G..TAAGTATGATATATTTTGGGAGGGTGGGGTGGAAAATGT  
.....GAATCTGA.T.GA.G.C.A.T...T.T.T.T.CG.A.C.A.A.A..G..C...T.T.A.T..GT..A.T.G..A...A...G...G..TAGGTGGTCTGCAACCCATGGGCTCAAGCAGTCCCTCCCGCT  
.....GGGACTT.A.CC.T.T.C.T...C.A.T.A.CT.T.C.T.G.T.T..A..G...A.A.A.A..GC..C.A.C..C...A...G...G..TAAGAAGAGACCACCGAAGACCCCGGGCTGATTTCTCTCT  
.....AGCGTGGTG.A.CC.C.C.T.T...G...A.A.CA.T.C.A...A...A...G.A.T.G..CA..T.A.T..G...A...G...G..TAAACAGGGGAGCATGAGAGCAGGACCATGCTGGGAGGGT  
.....AAGCATTCTT.T.G...T.A.T...G.C.C.A.AT.T.A.A.G.A..A..T...C.A.T.A..GG..G.A.A..G...A...G...G..TAGCCAATTTCCATTTTATAAAAGGTCATTGGAAATATTTT  
.....ATCAGGAAA.C.GG.T.G.G.G...A.C.A.T.TT.A.C.T.G.A..A...A.A.G.A..TG..A...A...A...G...G..TACCAAACTTGTGTGAGACATGAGCCAATTTAGGAAGATGT  
.....GCTCTATG.A.GA.C.T.T.C...A.C.T.G.CT.G.G.A.G.A..A.C...A.A.G.A..GA..A.A.T..C...A...G...G..TCAGTCTGATTACTAAAGAAAGTGCATAGAAAATACTTTTTAT  
.....GCAGAAAT.G.CA.C.C.A.T...T.T.T.C.AT.A.G.T.G.G.A..A..T...G.A.A.A.G..AT..T.G.C..T...C...C...G..TGAGTACAGAGTATGACTTTTTATTTTATAGAACAGACAT  
.....AGTGTGA.G.GA.A.T.T.A...G.A.A.A.AA.T.C.A.T.C..T..G...A.T.A.CA.TT.A.A.C..A...A...G...G..TAAAGCAGCTATTCTCTACTTGTCTTTTACATGATATAT  
.....TATCCAG.A.T.GT.C.A.G.T...A.A.G.G.AT.A.A.A.A.G..A..T...T.T.C.C..CT..G.G.A..C...A...G...G...TAAGTAACCCAATTTTACTTGAAGCTAGAAGGCTGTTTTT  
.....GTACACC.T.AT.A.A.A.A...A.A.T.G.GA.A.A.G.A.A..G..A...A.T.A.A..AG..T.T.A..T...A...G...G..TAAACCTTTCCCTTTAAACCTATACTCTAAACTTATTTTAAAT  
.....AGCAGCTA..A.GG.A.A.G.A...A.A.C.A.AG.A.G.A.G.C..T..A...C.A.G.T..CG..T.T.A..G...C...G...G..TAATTATCCTCTTTCTTTTTAAAAATTAATAATACAGACAT  
.....TGACCTCA.A.CC.T.T.G.T...G.C.A.G.AA.C.T.T.C.C.A..A..G...A.T.G.A..AA..C.A.T..G...A...G...G..TAGGAATGCAGCAACACTCTTTGGAAAAGTCTATGTGTATGT  
.....TCAGTGA.T.GT.T.T.G.C...T.G.C.T.CT.A.A.A.A.G.C..A...G.C.A.A..AA..A.T.T.C...A...G...G..TACTTTTTAGAGCATTATATATCTAATAACATTTTTGTGTCT  
.....ATTTGGATC.A.CC.G.T.G.G...T.G.T.G.AA.T...A.A.C..T..C...T.T.A.C..AT..A.C.A..G...C...G...G..TAAGACTTTCCAGTAAAGCAAGAAAATGTTTTCTTTTAAAT  
.....CCATCATT.T.TG.A.C.C.A...G.T.G.C.CT.G.A.T.G.A..T..T...C.T.T.A..AT..A.G.T..C...C...A...G..TAAGTACCTAATTTATTGTTGTGGAATTTTAAAGAGAGCTT  
.....AGCAAGCTA.T.TT.T.A.G.C...T.A.C.A.GA.C.C.T.A.G..C..A...C.T.G.T..AC..A.T.T..A...A...G...G...TAAGTTTTAAGAAAATAACCTGAGAAAATAGGTATTTGAAGTCT



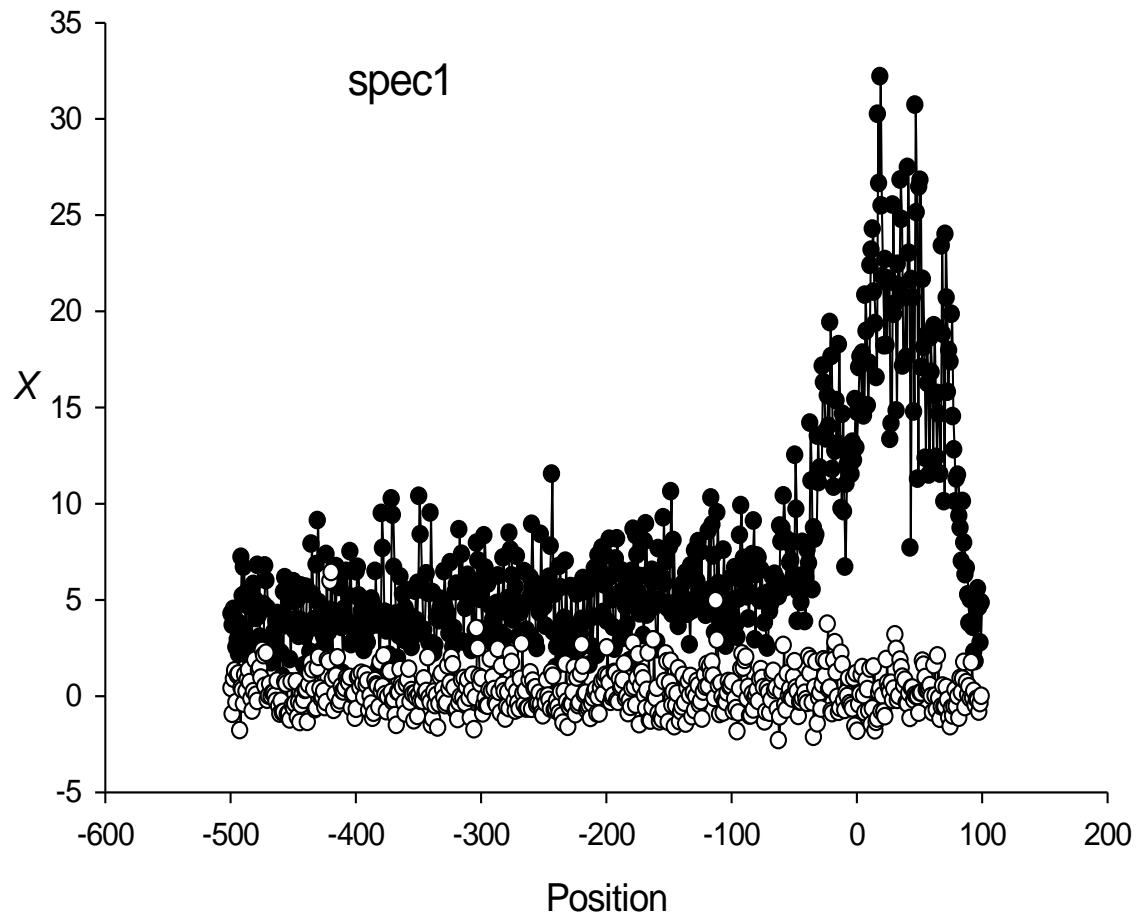
CONSENSUS=TTACCGCATCACCCAGCGACTGCAACATCAACTTGC  
ACCTTACGGTAACGGGGTAAGTTGCAAAAATCAGAAGGAGAAAA  
TTAGTAATTATGCA

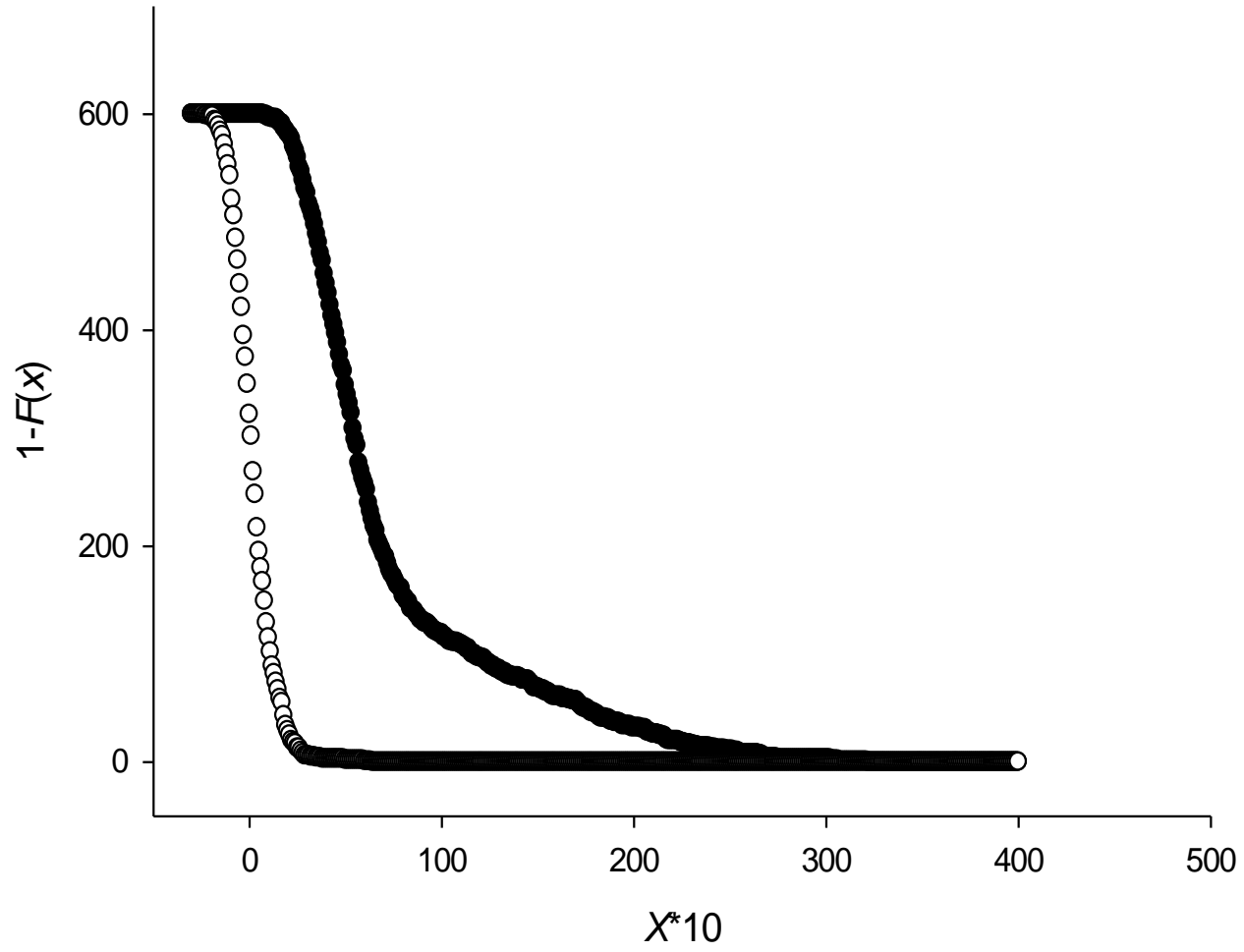


**Fig. 1.** Simplified model of transcriptional regulation of protein-encoding genes.

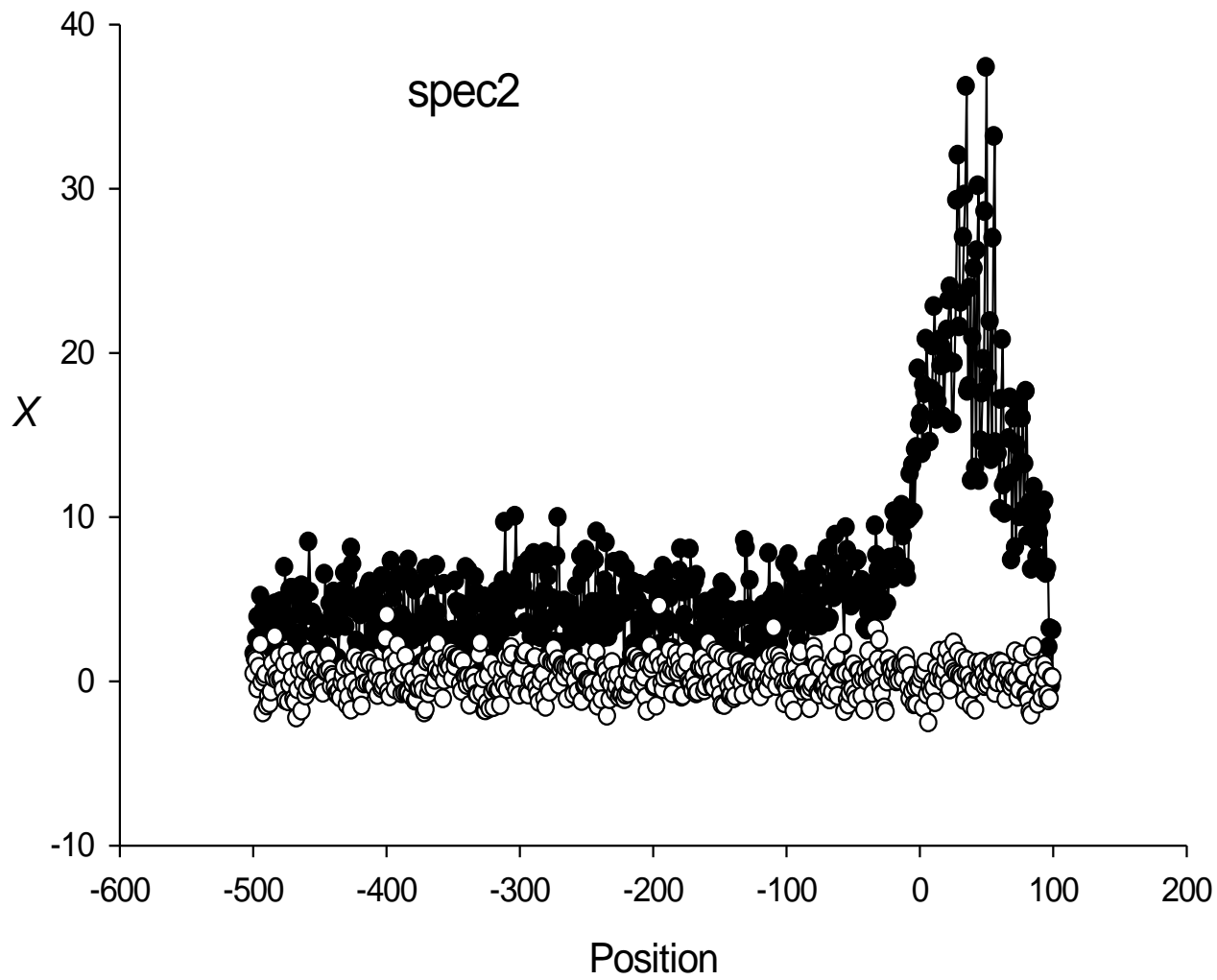


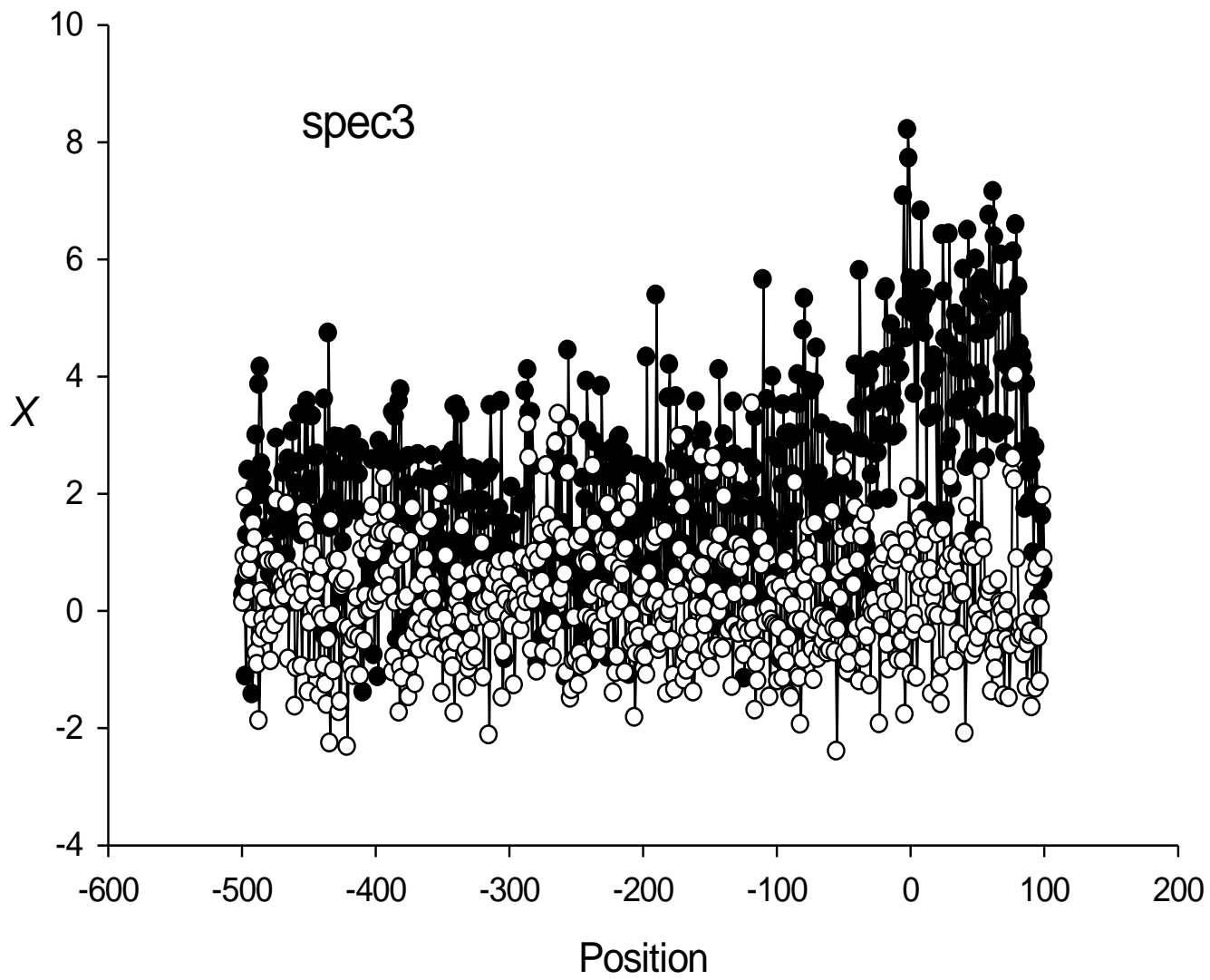
N	Volume of promoter group	Volume of random group
1	7247	32
2	1773	28
3	633	21
4	201	24
5	279	18
6	182	24
7	152	28
8	152	33
9	141	30
10	138	23
11	138	35
12	103	18
13	85	20
14	124	22
15	95	25
16	77	26
17	95	21
18	118	23
19	44	29
20	90	31
21	81	30
22	23	25
23	54	32
24	81	24
25	24	18

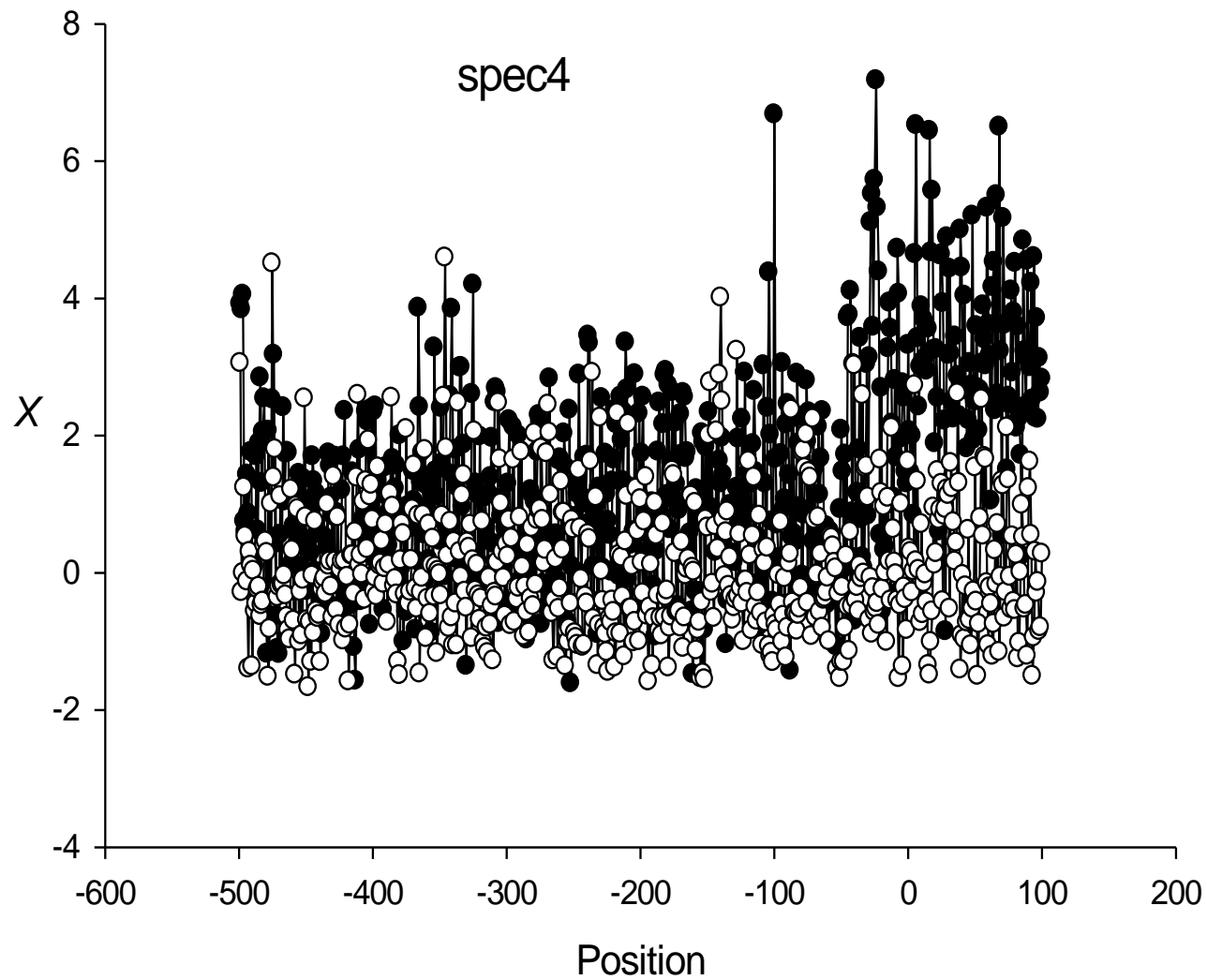


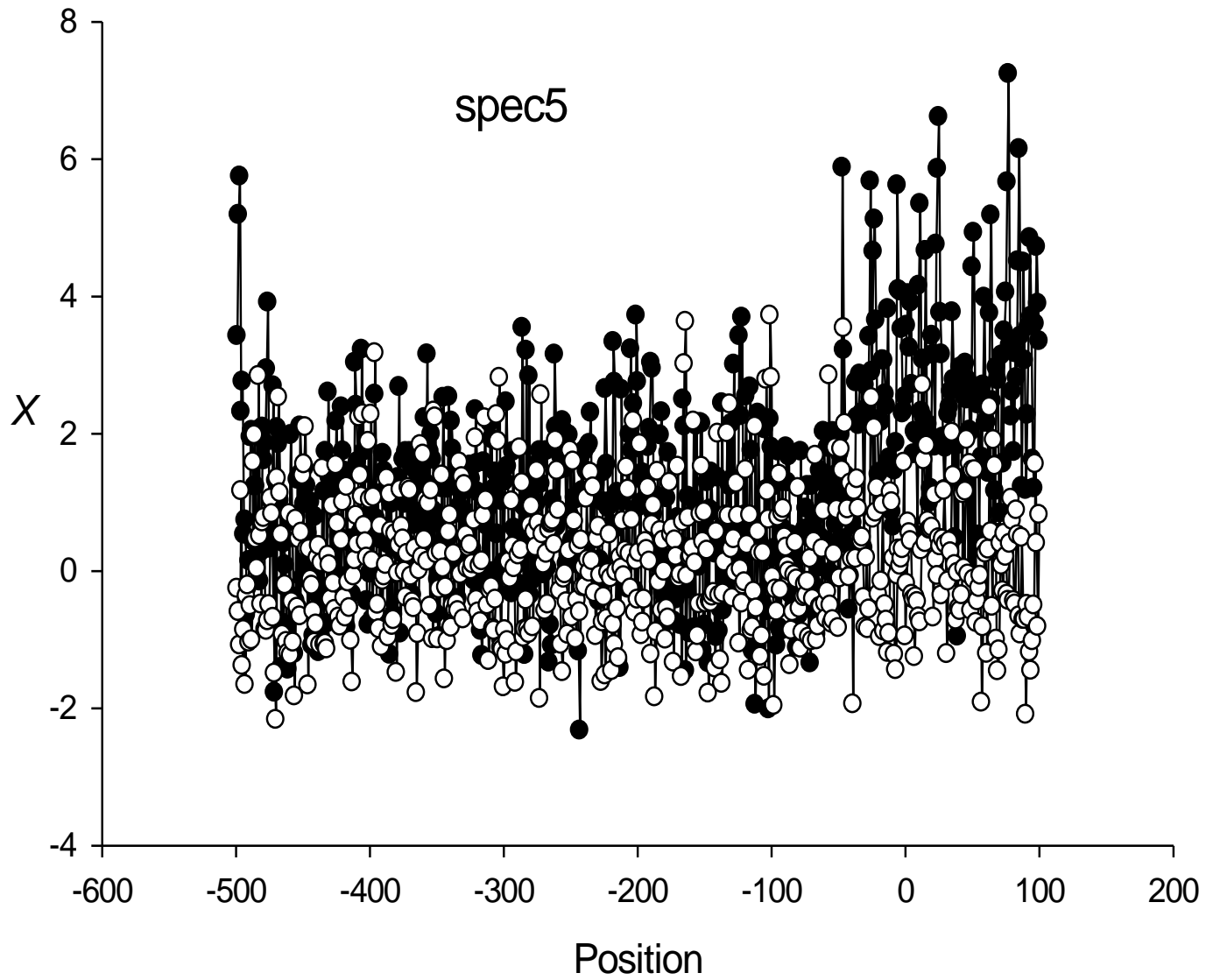


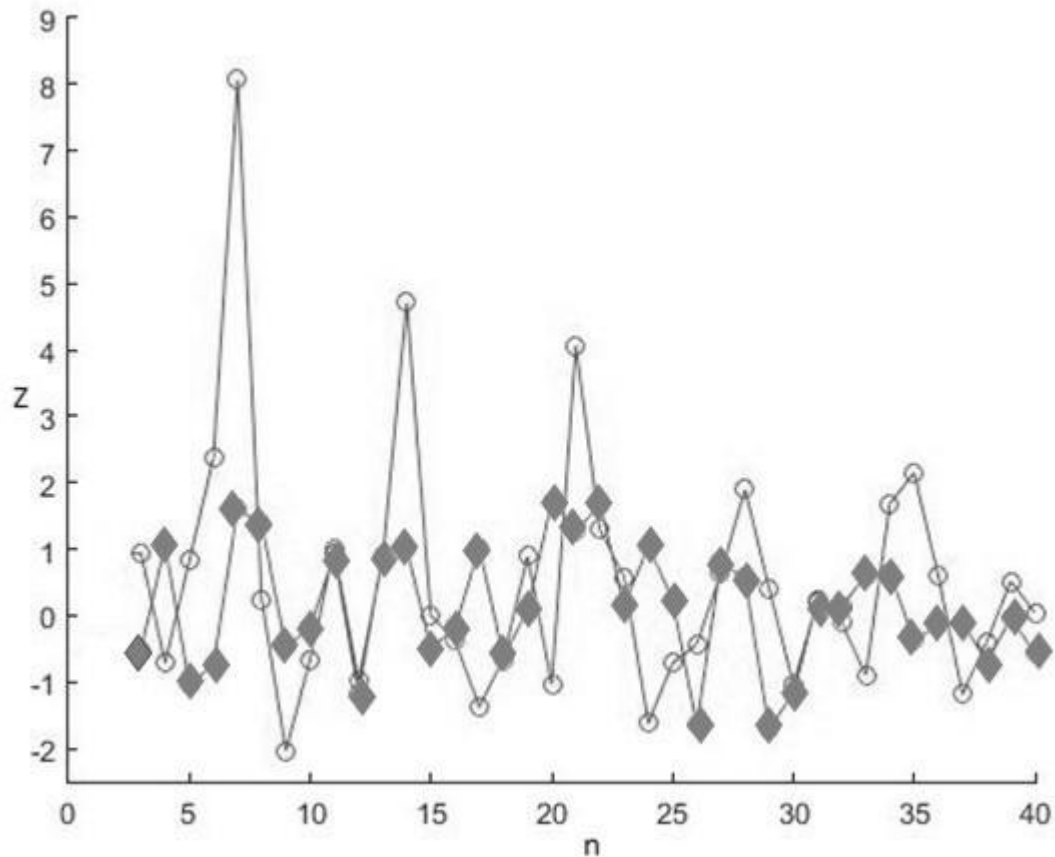




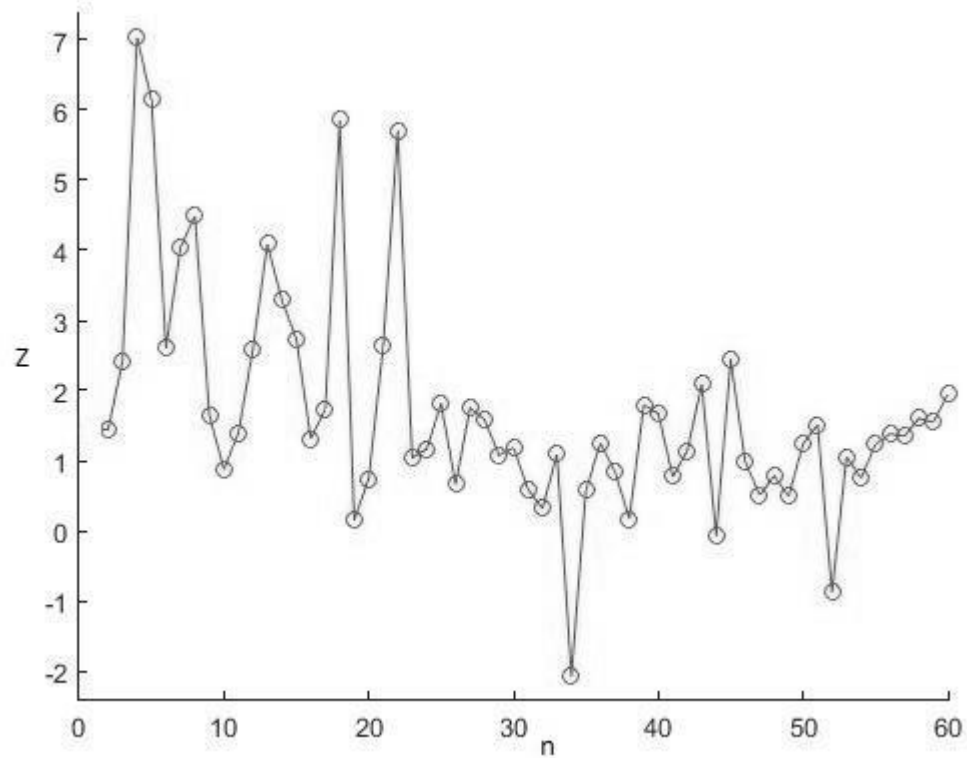




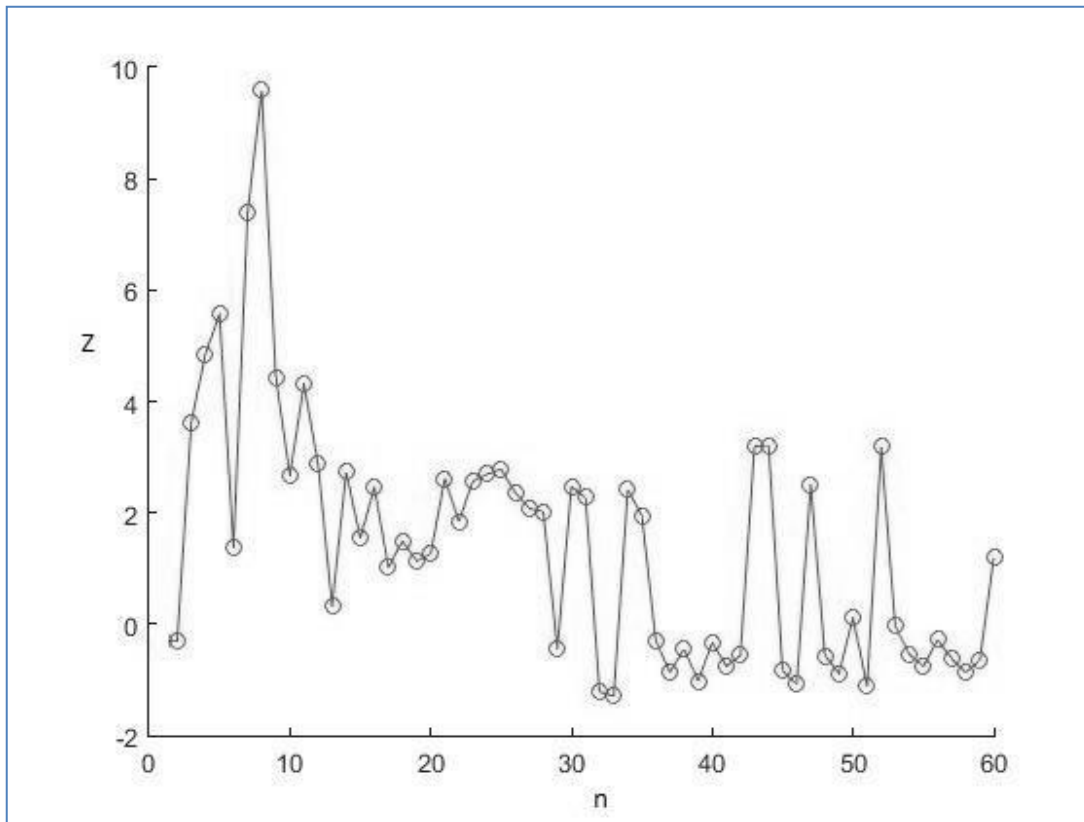




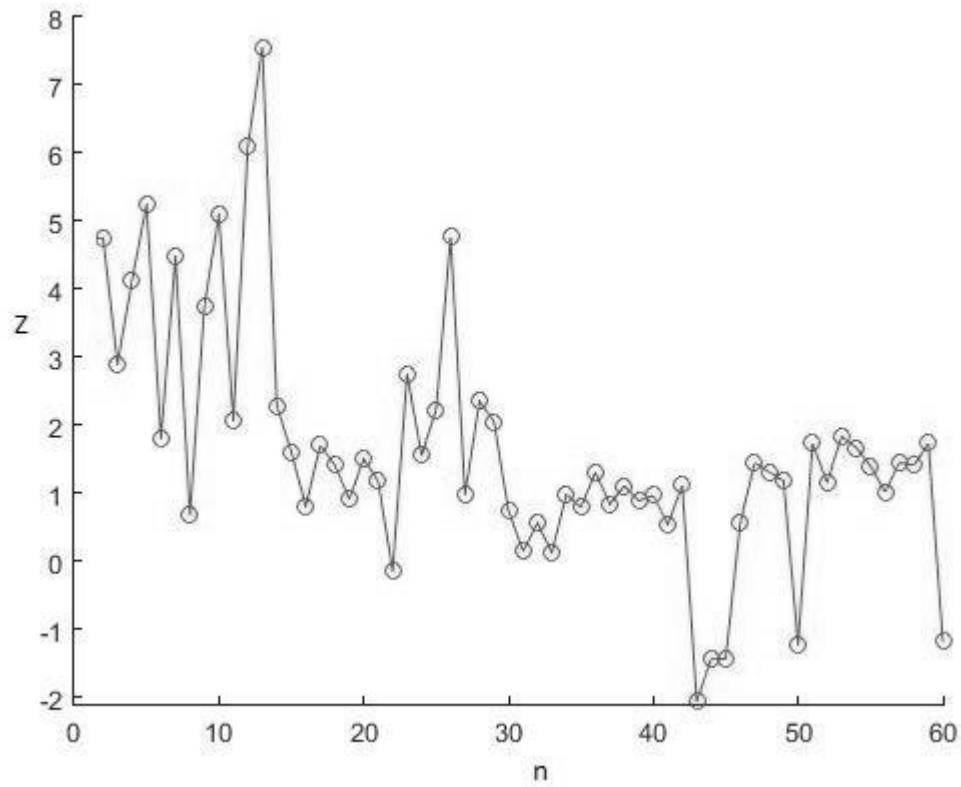
Graph  $Z(n)$  for an artificial sequence containing 60% random base substitutions with the addition of 8 inserts and 12 deletions (circles) and  $Z(n)$  for a random sequence (rhombus).



Graph  $Z(n)$  for the poem A.S. Pushkin "I remember a wonderful moment...".



Graph  $Z(n)$  for the William Blake's poem "Spring".



Graph  $Z(n)$  for the poem by Robert Frost "Fire and ice".



Lengths of fuzzy periods found in 95 works of poetry of different authors.

Author	Number of analyzed works of poetry	Lengths of fuzzy periods found in the works of poetry
Pushkin A.	15	2,4,7,10
Yesenin S.	10	2,6,8,11
Blok A.	5	3,5,10,11
Tutchev F.	5	2,7,8,15
Fet A.	5	2,4,6,9
Mayakovsky V.	7	2,3,5,8
Shakespeare W.	5	10,15,18,37
Byron D.	5	8,16,28
Frost P.	20	5,8,13,15,20
Blake W.	18	4,6,8,10,11,23