



**Три вещи, которые нужно уточнить,  
решая содержательную задачу**

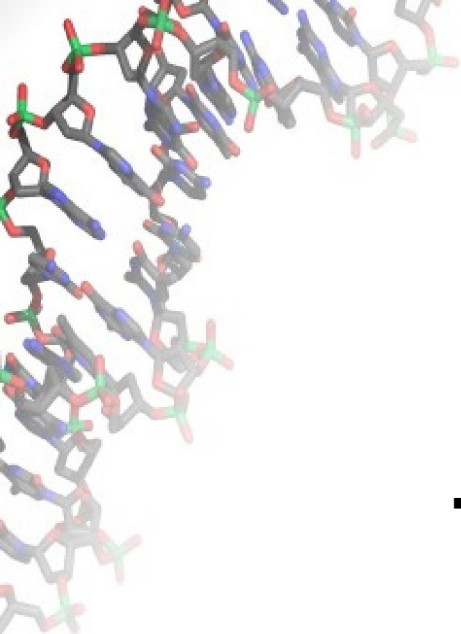
***1. ФОРМАЛЬНАЯ ПОСТАНОВКА ЗАДАЧИ.***

***2. АЛГОРИТМ РЕШЕНИЯ***

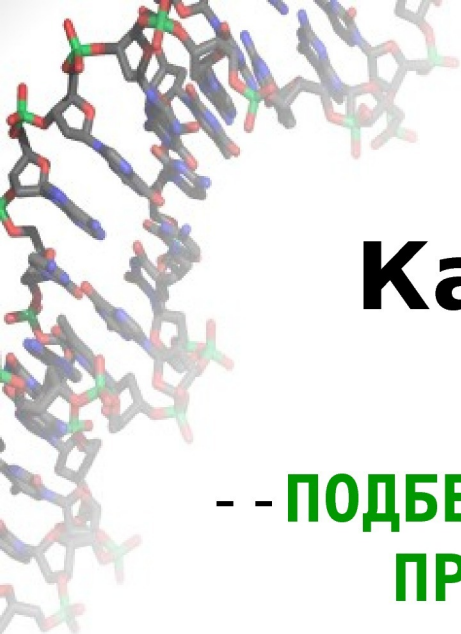
***3. ПРОВЕРКА - АДЕКВАТНО ЛИ ТО, ЧТО  
ПОЛУЧАЕТСЯ***

***(ИСХОДНОЙ СОДЕРЖАТЕЛЬНОЙ ЗАДАЧЕ).***

**Как правило, уточнение этих вещей  
происходит взаимосвязано и  
итеративно**



# **Часть 1. Точность и достоверность выравниваний.**



## Варианты выравниваний

# Какой вариант выбрать?

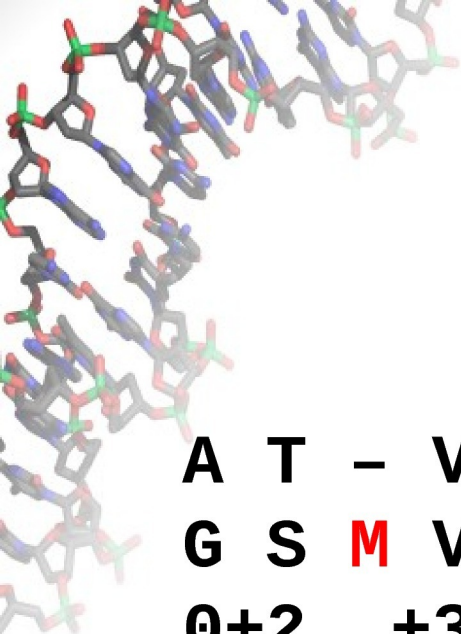
-- **ПОДБЕРЕЗОВИК**  
**ПРЕДОСИНОВИЧКИ**

**ПОДБЕРЕЗОВИК** --  
**ПРЕДОСИНОВИЧКИ**

**ПО -ДБЕРЕЗОВИК - -**  
**ПРЕДОСИН - ОВИЧКИ**

**ПО -ДБЕРЕЗОВИК - -**  
**ПРЕД - ОСИНОВИЧКИ**

**ПО -ДБЕРЕЗОВИ - К -**  
**ПРЕД - ОСИНОВИЧКИ**



# **Вес выравнивания**

**Пример:  $A=0$ ;  $B = 2$**

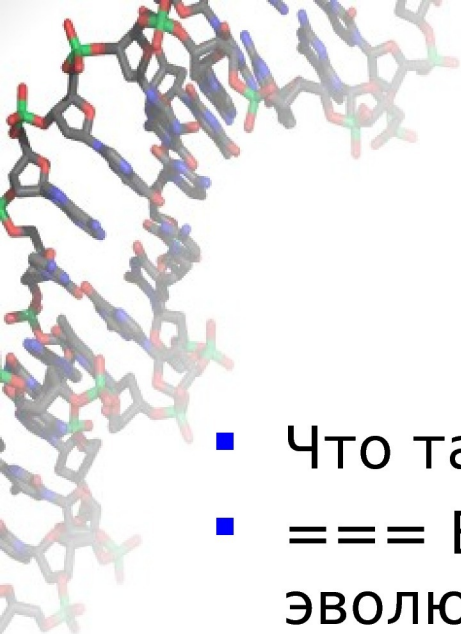
A	T	-	V	V	I	-	-	T	G	S
G	S	M	V	L	L	E	F	S	G	T
0+2			+3+2+3					+2+7+2=	21	
		-2				-4		=	-6	

**Штраф за удаление фрагмента длины L:**

$$d(L) = A + B * L$$

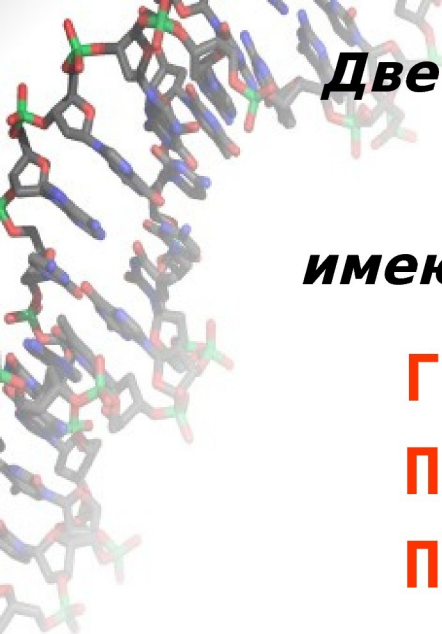
$$Score = \sum m(i,j) - \sum d(L_k)$$

$$Score = 21 - 6 = 15$$



## Правильно ли оптимальное выравнивание?

- Что такое «правильное выравнивание»?
- === Выравнивание, отражающее эволюционный процесс (не обязательно биологический).



**Две одинаковые буквы скорее имеют общего предка, чем две разные буквы**  
**Две буквы «одинаковой гласности» скорее имеют общего предка, чем две буквы «разные гласности»**

**Г)**

**П-ОДБЕРЕЗОВИК - -  
ПРЕД-ОСИНОВИЧКИ**

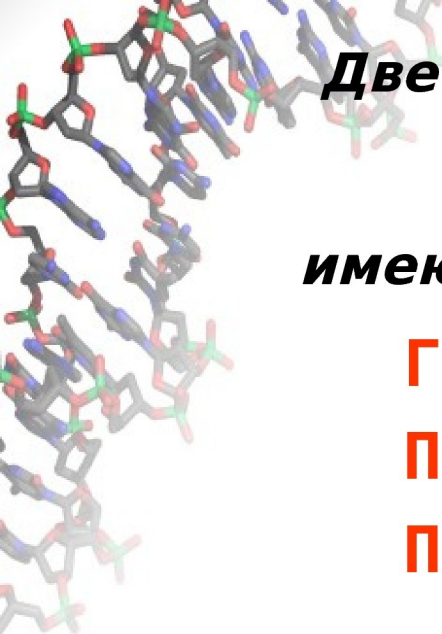
**Д)**

**ПО-ДБЕРЕЗОВИ-К-  
ПРЕД-ОСИНОВИЧКИ**

**При любом штрафе за удаление символа**

**Д) лучше, чем Г)**

**(удаление – посимвольное)**



**Две одинаковые буквы скорее имеют общего предка, чем две разные буквы**  
**Две буквы «одинаковой гласности» скорее имеют общего предка, чем две буквы «разные гласности»**

**Г)**

**П-ОДБЕРЕЗОВИК - -**

**ПРЕД-ОСИНОВИЧКИ**

**Д)**

**ПО-ДБЕРЕЗОВИ-К-**

**ПРЕД-ОСИНОВИЧКИ**

**НО :**

**Д) – «неправильно» («эволюционно»)**

**ПОД -БЕРЕЗ-ОВИ-К**

**ПРЕД-ОСИН –ОВИ-К**

**ПРЕД-ОСИН –ОВИ-Ч-ЕК**

**ПРЕД-ОСИН –ОВИ-Ч- К-И**



# Эталон: эволюционное выравнивание

ASVVLDFGTGT

ASVVLDFGTGT

ATVVI--TGS

AS-VVLDFGTGT

GSMVLLEFSGT

AS-VVLDFGTGT

AT-VVI--TGS

AS-VVLDFGTGT

GSMVLLLEFSGT

AS-VVLDFGTGT

AT-VVI--TGS

GSMVLLEFSGT





# Эталон: эволюционное выравнивание

**ASVVLDFGTGT**

**ASVVLDFTGT**

**AS\_VVLDFGTGT**

**ATVVI--TGS**

**GSMVLLEFSGT**

**AS-VVLDFGTGT**

**AS-VVLDFGTGT**

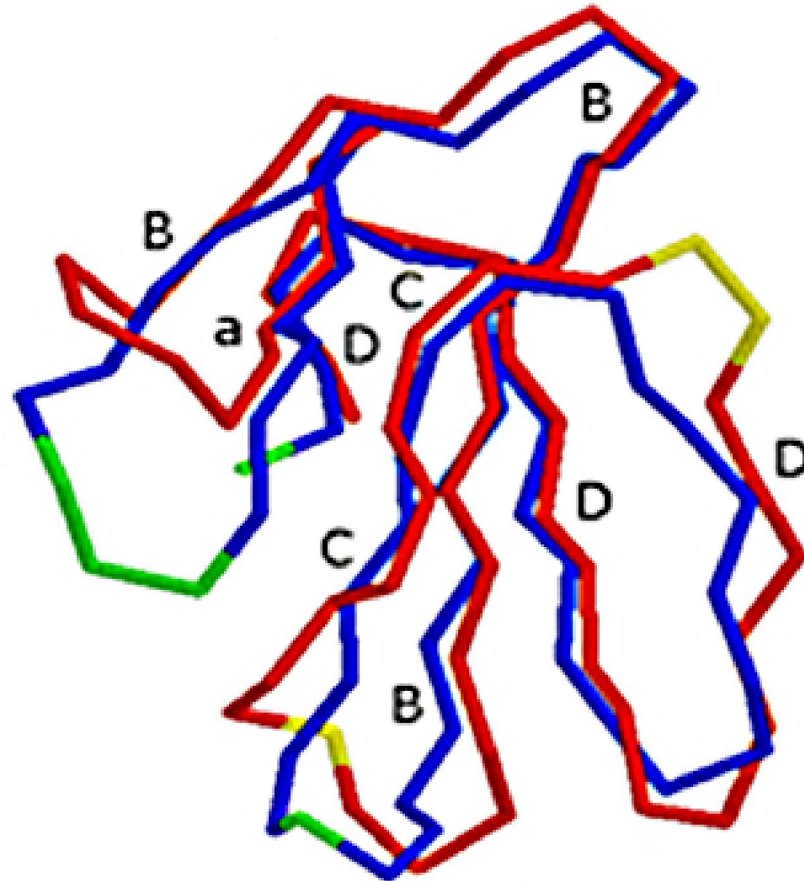
**AT-VVI--TGS**

**GSMVLLEFSGT**

**AT-VVI--TGS**

**GSMVLLEFSGT**

# Приближение: структурные выравнивания



# Структурное и алгоритмическое выравнивания

str) 40  
сопоставлений  
lkCnqli...PPFWKTCPKGKNLCYKmtmraapmvPVKRGcidv  
riCfnhqssqPQTTKTCSPGESSCYHkqwdsfrgtIIERGCg..  
\* \*\*\*\* \*  
1 16 6

AlgSW)  
1 16 6  
\* \*\*\*\* \*  
lk...C...nqliPPFWKTCPKGKNLCYK...mtmraapmvPVKRGcidv  
..riCfnhqssqPQTTKTCSPGESSCYHkqwdsfrgt...IIERGC..g  
35 сопоставлений

$S = 39$

$I = 23$

$A = 34$

*Точность*

$Acc = I/S =$

$23/39 = 0.59$

*Достоверность Conf =*

*Str)  $I/A = 23/34 = 0.68$*

1234567	89012345678901234567890123456789
lkCnqli...	PPFWKTCPKGKNLCYKmtmraapmvPVKRGcidv
riCfnhqssq	PQTTKTCSPGESSCYHkqw sdf rgtIIERGCg..
*	*****

1

16

6

AlgSW)

1

16

6

\*

\*\*\*\*\*

\*\*\*\*\*

lk..C...	nqliPPFWKTCPKGKNLCYK...	mtmraapmvPVKRGcidv
..riCfnhqssq	PQTTKTCSPGESSCYHkqw	sdf rgt...IIERGC..g
1	23456789012345678901	234567 890123 4

lkCnqli...PPFWKTCPKGNLCYKmtmraapmvPVKRGCi  
dv

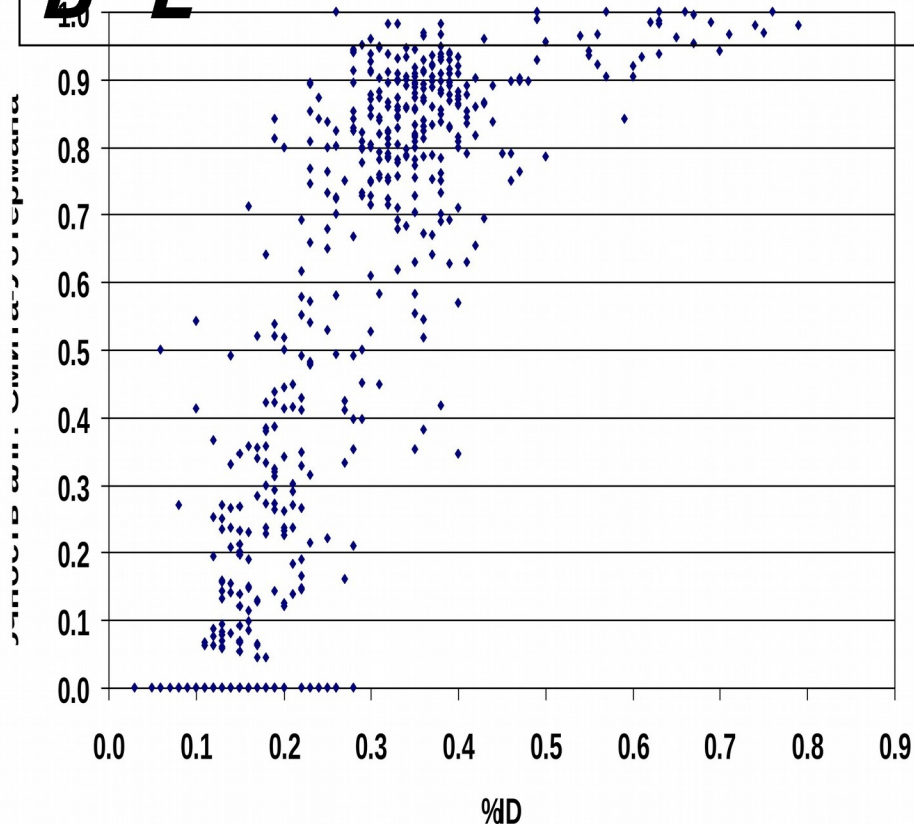
riCfnhqssqPQTTKTCSPGESSCYHkqwsdfrgtIIERGCg

%ID=  
= 11/39 =  
= **0.28**  
Acc = **0.59**

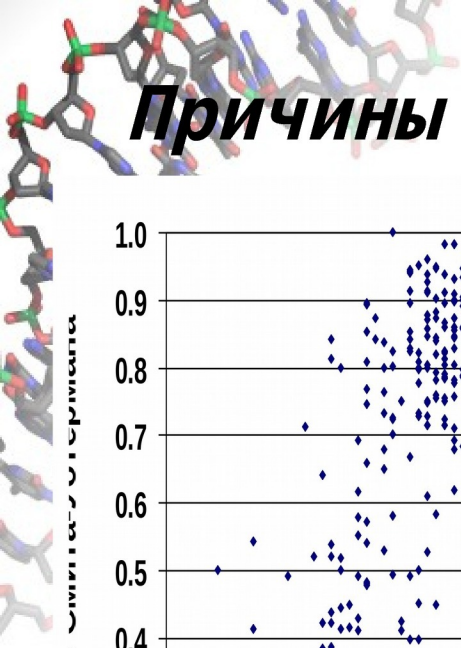
Алгоритм Смита-Уотермана (SW).  $d(L) = A + B * L$

При ID < 0.3

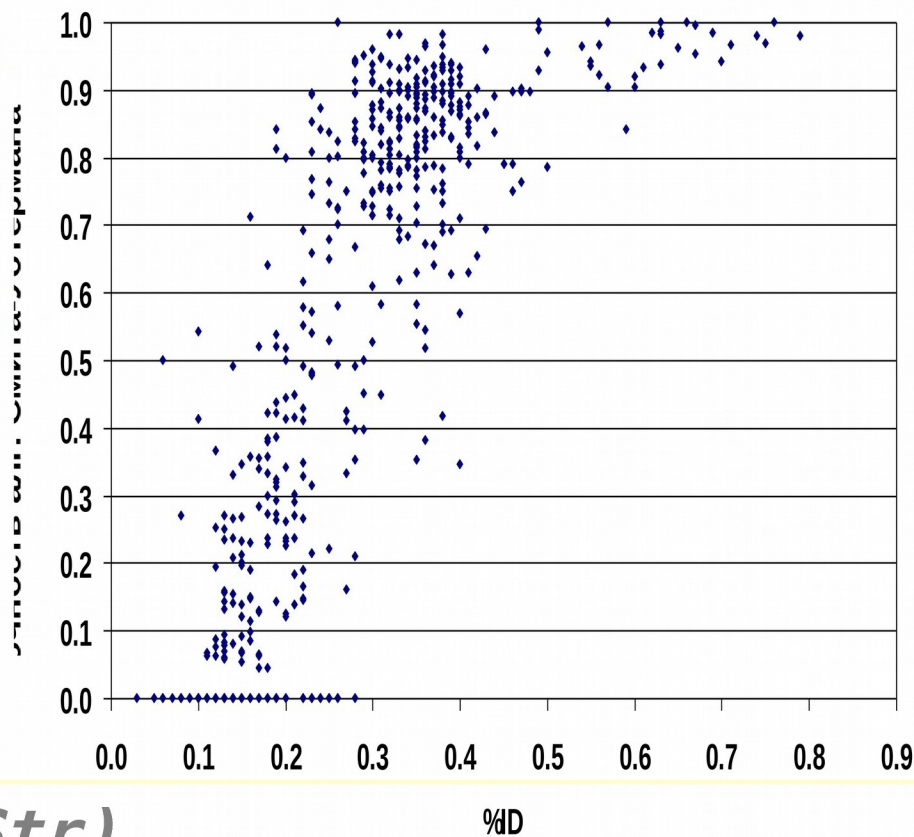
SW-выравнивание  
НЕ ПОХОЖЕ  
на структурное



%ID	SW точность (acc)
< 0,1	<b>0,037</b>
0,1-0,3	<b>0,306</b>
0,3-0,4	0,818
>0,4	0,893



# Причины плохого качества выравниваний SW



**Острова -  
безделеционные  
фрагменты  
выравниваний.**

**Вес острова - сумма  
весов  
сопоставлений**

*Str)*

lkCnqli...PPFWKTCPKGKNLCYKmtmraapmvPVKRGCidv  
riCfnhqssqPQTTKTCSPGESSCYHkqwsdfrgtIIERGCg..

^^^

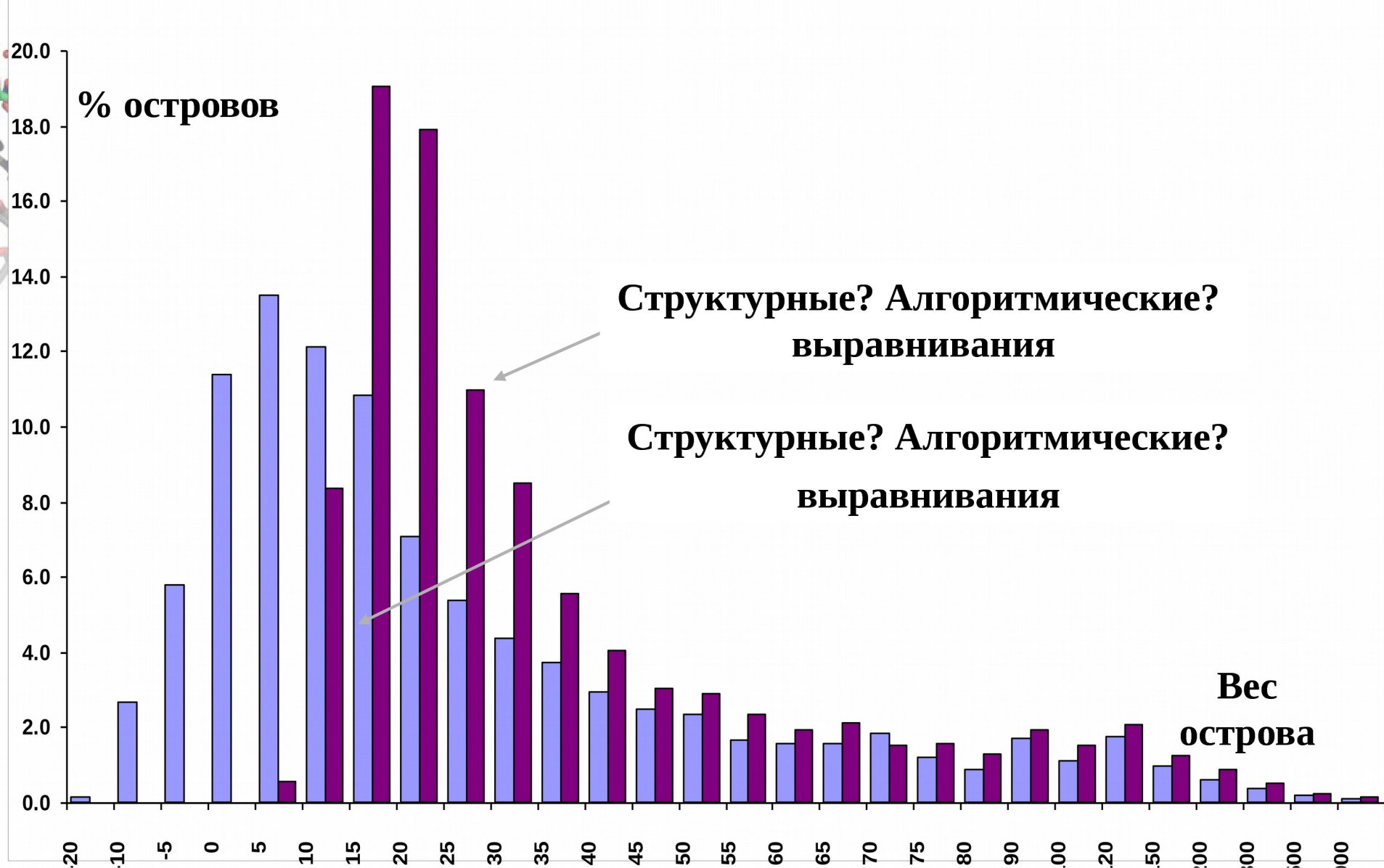
^^^

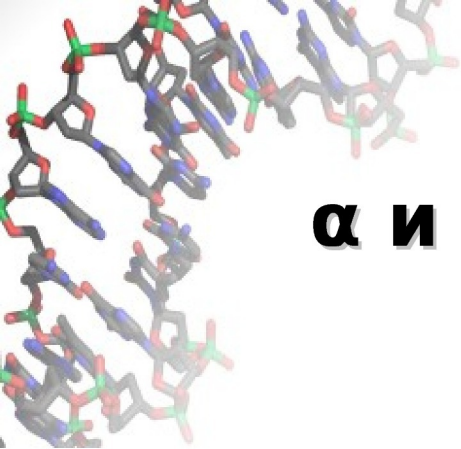
Остров 1

Остров 2

# Причины плохого качества выравниваний SW

## Гистограммы весов островов





# Пример. α и β цепи гемоглобина лошади

	10	20	30	40	50	60	70	80			
1	<u>234567890</u>	<u>1234567890</u>	<u>12345678901234</u>	<u>56789012</u>		6789012456789012345678901					
	V-LSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPHF-DLSHGSAQ-----				VKAHGKKVGDALTLAVGHLDDLPGA						
	<u>VqLSgeeKaaVlAlWdKVnee--</u>		<u>EvGgEALgRllvypwTqrFFdsFgDLSnpgAvmgnpkVKAHGKKVlhsfgegVhHLDnLkGt</u>								
1	23	<u>4</u>	5 6 7 8 9	0	1 2 3 4 5	6	7 8 9 0 1	2	3 4 5 6 7 8 9 0	1	2 3 4 5 6
				10			20		30		

$$\%ID = 36/81 \approx 0.44$$



# Пример. $\alpha$ и $\beta$ цепи гемоглобина лошади

Оптимальные  
выравнивания,  
содержащие  
1, 2 и 3  
удаленных  
фрагмента

**S** - сумма весов  
сопоставлений

$S(1)=182$ ;

VLSAADKTNVKAAWSKVGGHAGEYGAERALERMFLGFPTTKTYFPHF  
VqIsgeekaav1A1wdkvneeevgGealgr11vvypwTqrffdsfg

|=====>

DLSHGSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA

DLsnpgAvmgnpkVKAHGKKVlHsfgegVhHLDnLkGt

-----  
 $S(2)=223$ ;  $S(2) - S(1) = 41$

VLSAADKTNVKAAWSKVGGHAGEYGAERALERMFLGFPTTKTYFPHF  
VqIsgeekaav1A1wdkvnee-EvGgEALgR11vvypwTqrFFdsF

|=====

DLSHGSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA

gdlsnpgavmgnpkVKAHGKKVlHsfgegVhHLDnLkGt

=====>

-----  
 $S(3)=272$ ;  $S(3) - S(2) = 49$

|=====>

V-LSAADKTNVKAAWSKVGGHAGEYGAERALERMFLGFPTTKTYFPHF

VqLSgeeKaaV1A1wdKVnee--EvGgEALgR11vvypwTqrFFdsF

DLSHGSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA

gdlsnpgavmgnpkVKAHGKKVlHsfgegVhHLDnLkGt

S(3)=272; S(3) - S(2) = 49

|=====>

V-LSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPFH  
VqLSgeeKaaVlAlWdKVnee--EvGgEALgRllvvypwTqrFFdsF

DL SHGSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA  
gdlsnpgavmgnpkVKAHGKKVlhsfgegVhHLDnLkGt

S(4)=303; S(4) - S(3) = 31

V-LSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPFH  
VqLSgeeKaaVlAlWdKVnee--EvGgEALgRllvvypwTqrFFdsF

|=====>

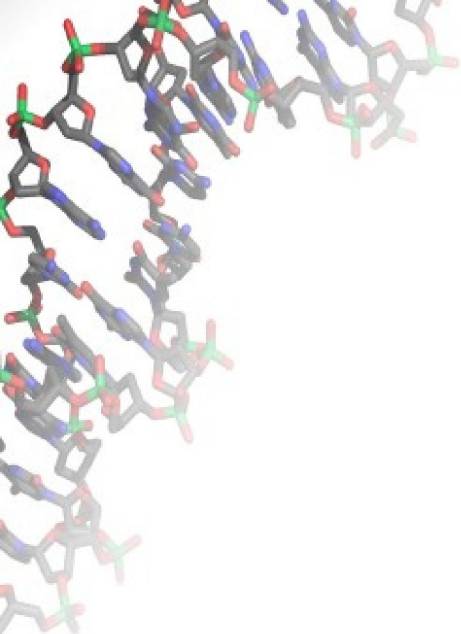
-DL SHGSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA  
gDLSnpgAvmgnpkVKAHGKKVlhsfgegVhHLDnLkGt

S(5)=310; S(5) - S(4) = 7

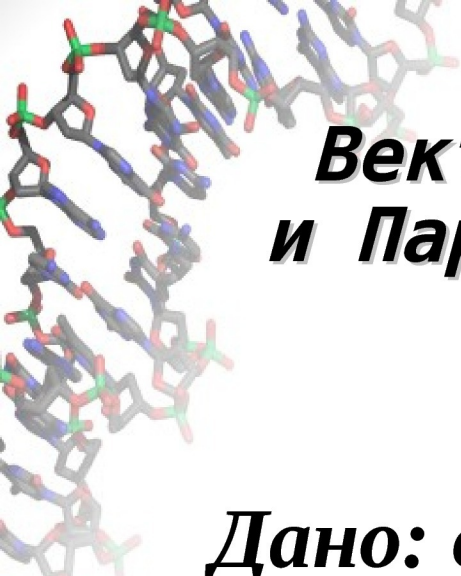
V-LSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPFH  
VqLSgeeKaaVlAlWdKVnee--EvGgEALgRllvvypwTqrFFdsF

|==>

-DL SH-GSAQ-----VKAHGKKVGDALTLAVGHLDDLPGA  
gDLSnpgAvmgnpkVKAHGKKVlhsfgegVhHLDnLkGt



**КАК  
УЛУЧШИТЬ  
КАЧЕСТВО  
ВЫРАВНИВАНИЙ?**

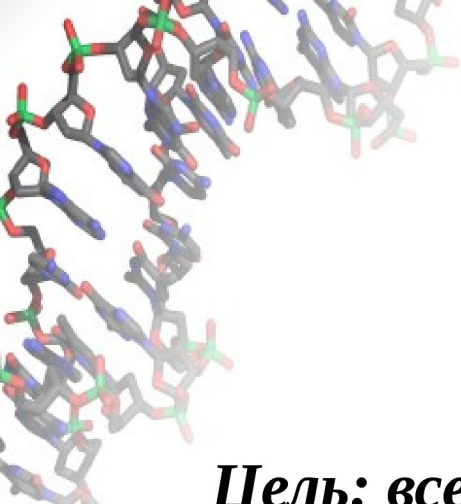


# **Векторные веса выравниваний и Парето-оптимальные выравнивания**

$$W(G) = c \cdot p - f \cdot r - d \cdot s$$

**Дано: вес выравнивания, зависящий от характеристик выравнивания и числовых параметров**

**Цель: все выравнивания, которые могут быть оптимальными при заданном виде весовой функции и каких-либо значениях числовых параметров.**



$$W(G) = c \cdot p - f \cdot r - d \cdot s$$

**Цель:** все выравнивания, которые могут быть оптимальными при заданном виде весовой функции и произвольных числовых параметрах.

**Идея:**

- использовать векторные веса:

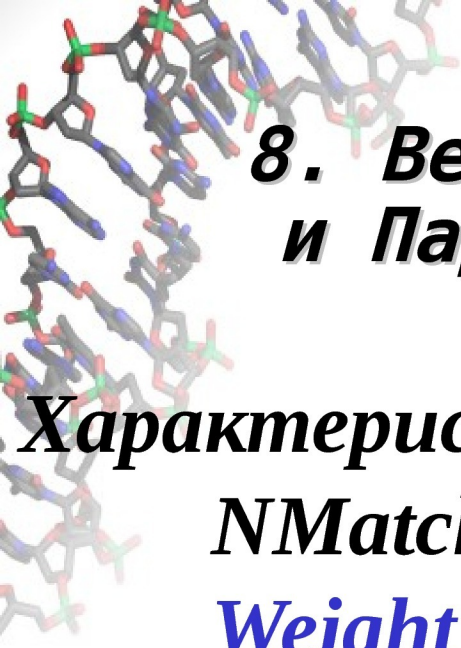
$$V(G) = \langle p, r, s \rangle$$

- искать Парето-оптимальные выравнивания



**КАК УЛУЧШИТЬ КАЧЕСТВО  
ВЫРАВНИВАНИЙ?**

**ИСПОЛЬЗОВАТЬ  
СПЕЦИФИКУ ЗАДАЧИ!**



## **8. Векторные веса выравниваний и Парето-оптимальные выравнивания**

**Характеристики:**

$NMatch(G)$ ,

$WeightMatch(G)$

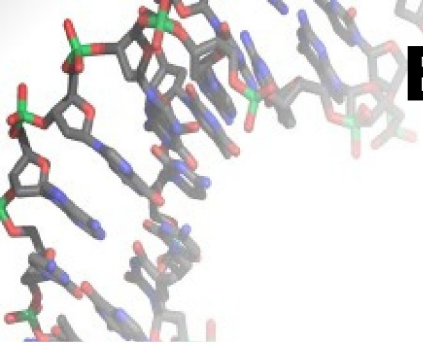
$-NDel(G)$ ,

$-NGap(G)$

**Алгоритм на основе обобщенных  
статистических сумм.**

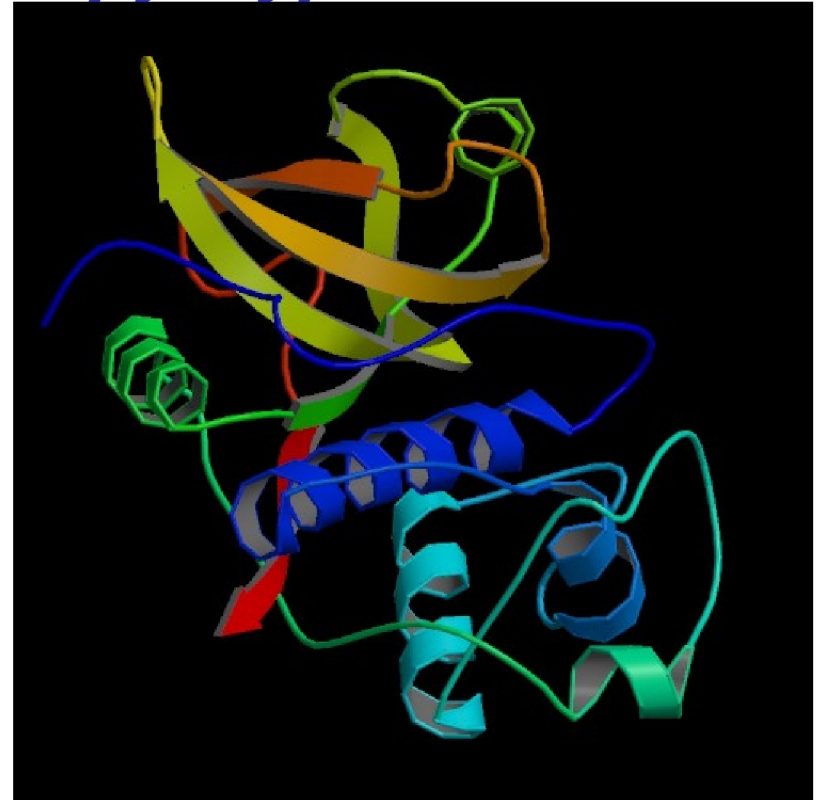
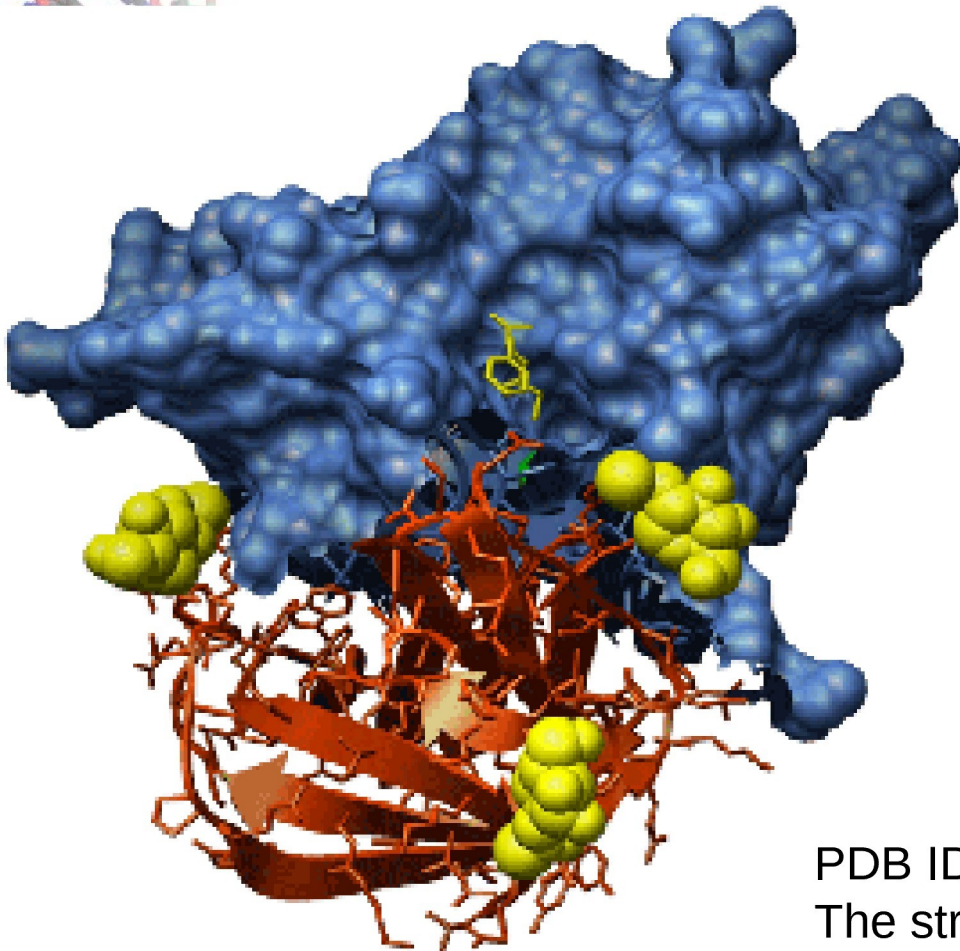
**Время для 2-компонентных весов:**

$$\sim \min(L_1, L_2)L_1L_2$$



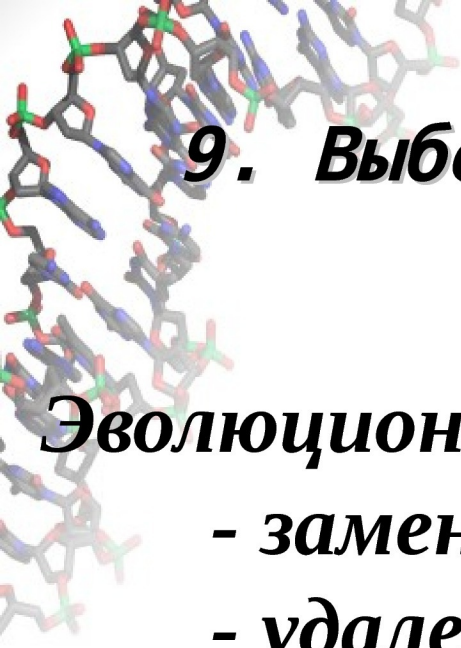
**Белки: 1 нить;  $L \sim 10^2 - 10^3$   
аминокислоты (20)**

*вторичная структура*



PDB ID: **2act** E.N. Baker, E.J. Dodson (1980):  
The structure of actinidin at 1.7 Ångstroms





## **9. Выбор биологически корректного выравнивания**

**Эволюционные события:**

- замены символов;**
- удаления/вставки фрагментов**

**Вес:**

**$\langle \text{WeightMatch}(G), -N\text{Gap}(G) \rangle$**



## **2. Улучшение качества выравниваний – учет вторичной структуры.**

$$W(i, j) = \max \begin{cases} W(i-1, j-1) + M(a_i, b_j) + SBON \times Q(i, j) \\ W(i-1, j) - GOP - GEP \\ WA(i-1, j) - GEP \\ W(i, j-1) - GOP - GEP \\ WB(i, j-1) - GEP \\ 0 \end{cases}$$

**$Q(i, j)$  – мера сходства вторичной  
структуры в  $i$ -й и  $j$ -й позициях**

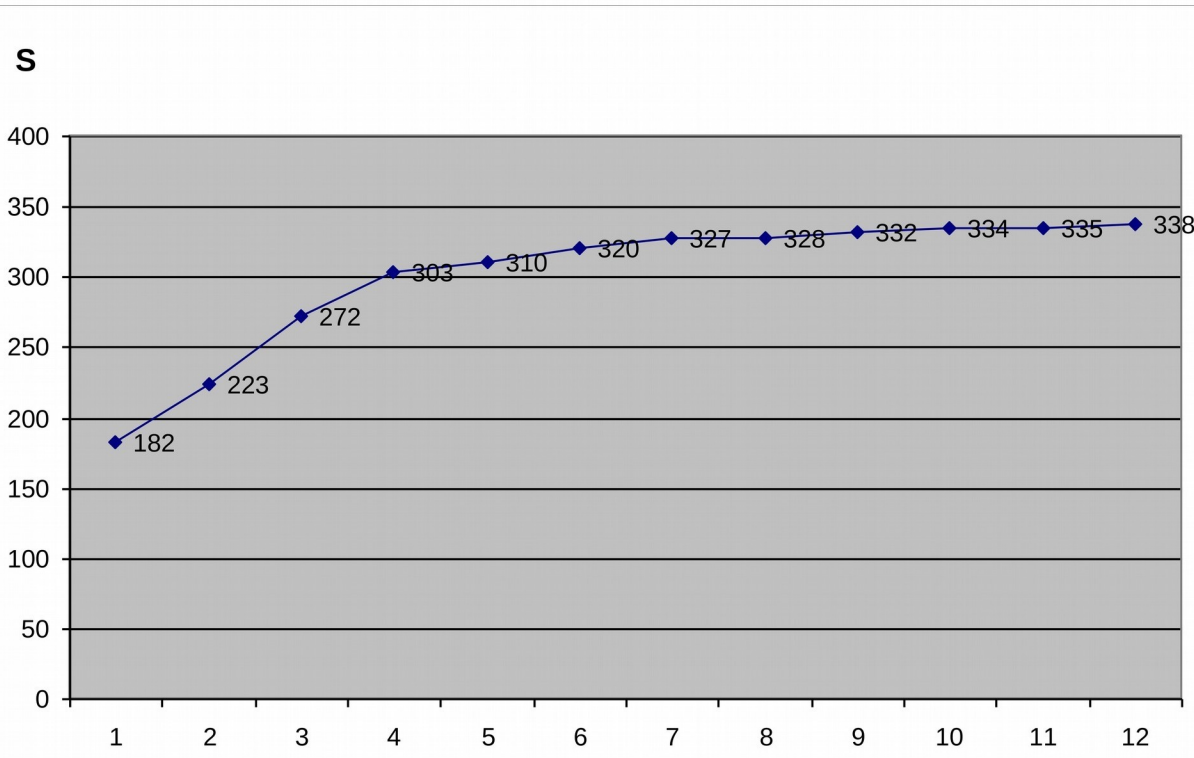
**$SBON$  – вес учета вторичной структуры**

# Выбор биологически корректного выравнивания.

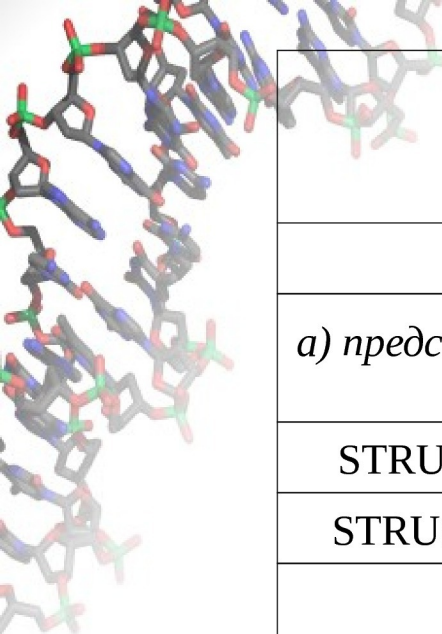
Ось X:  
*NGap, g*

Ось Y:  
*WeightMatch, S(g)*

$$Dcr \leq \log(L)/\log(1/(1-p)) - 1$$



$$D(g) = S(g) - S(g-1)$$



Метод	bonu s	GO P	GE P	Acc	Con f	Acc, ID < 30%	Conf, ID < 30%
SW	-	7	1	0.53	0.59	0.353	0.429
<i>а) предсказание вторичной структуры по последовательности</i>							
STRUSWER_SIN_S	2	10	1	0.58	<b>0.62</b>	0.428	<b>0.482</b>
STRUSWER_SIN_%	7	8	2	<b>0.6</b>	0.62	<b>0.461</b>	0.477
WFMFL_SIN	матр ица	13	1	0.4	0.49	0.263	0.346
<i>б) предсказание вторичной структуры с привлечением данных о гомологичных белках</i>							
STRUSWER_PSI_S	8	9	1	0.66	0.68	0.546	0.573
STRUSWER_PSI_%	17	6	2	<b>0.68</b>	<b>0.7</b>	<b>0.579</b>	<b>0.589</b>
WFMFL_PSI	матр ица	16	1	0.63	0.67	0.503	0.56
<i>в) экспериментально известная структура</i>							
STRUSWER_EXP	8	10	1	<b>0.68</b>	<b>0.7</b>	<b>0.577</b>	0.601
WFMFL_EXP	матр ица	15	1	0.64	0.7	0.527	<b>0.602</b>