

Мутации и выравнивание

С.А. Спирин

27 апреля 2021

План

0. Введение: гомологичные белки
1. Источники гомологичных белков
 - Мутации:
 - Ошибки репликации.
 - Повреждения ДНК и их reparация.
 - Закрепление мутаций
2. Выравнивание:
 - последовательностей потомков относительно предка;
 - двух потомков одного предка.
3. Формализация: вес выравнивания
4. Программы парного выравнивания в EMBOSS
5. Редактор выравниваний Jalview

Последовательности миоглобинов человека, мыши и быка

>MYG_HUMAN

MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKA
SE DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH
PGDFGADAQGAMNKALELFRKDMASNYKELGFQG

>MYG_MOUSE

MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIGLFKTHPETLDKFDKFKNLKSEEDMKGSE
DLKKHGCTVLTALGTILKKKGQHAAEIQPLAQSHATKHKIPVKYLEFISEIIIIEVLKKRH
SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG

>MYG_BOVIN

MGLSDGEWQLVLNAWGKVEADVAGHGQEVLIRLFTGHPETLEKFDKFKHLKTEAEMKA
SE DLKKHGNTVLTALGGILKKKGHHEAEVKHLAESHANKHKIPVKYLEFISDAIIHV
LHAKH PSDFGADAQAAMSKAELFRNDMAAQYKVLGFHG

Напишем последовательности друг под другом, чтобы было видно сходство:

MYG_HUMAN	MGLSDGEWQLVLNVWGKVEADIPGHQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKA	60
MYG_MOUSE	MGLSDGEWQLVLNVWGKVEADLAGHGQEVLIGLFKTHPETLDKFDKFKNLKSEEDMKG	60
MYG_BOVIN	MGLSDGEWQLVLNAWGKVEADVAGHGQEVLIRLFTGHPETLEKFDKFKHLKTEAEMKA	60
	***** . * : * * * * * * : * * * * * * : * * : * : * * . *	
MYG_HUMAN	DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH	120
MYG_MOUSE	DLKKHGCTVLTALGTILKKKGQHAAEIQPLAQSHATKHKIPVKYLEFISEIIIIEVLKKRH	120
MYG_BOVIN	DLKKHGNTVLTALGGILKKKGHHEAEVKHLAESHANKHKIPVKYLEFISDAIIHVHLAKH	120
	***** * * * * * * : * * : * * : * * * * * * * * : * * : * : *	
MYG_HUMAN	PGDFGADAQGAMNKALELFRKDMASNYKELGFQG	154
MYG_MOUSE	SGDFGADAQGAMSKALELFRNDIAAKYKELGFQG	154
MYG_BOVIN	PSDFGADAQAAMSKAELFRNDMAAQYKVLGFHG	154
	***** . * * * * * * : * : * : * * : * * : *	

Видно, что большинство букв совпадает, но некоторые различаются.

Это последовательности **гомологичных** белков, что означает, что эти белки произошли от общего предка. За время, прошедшее от существования общего предка, некоторые буквы менялись, но большинство остались неизменными.

Последовательности миоглобинов человека и рыбы

>MYG_HUMAN

MGLSDGEWQLVNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKA
DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH
PGDFGADAQGAMNKALELFRKDMASNYKELGFQG

>MYG_DANRE

MADHDLVLKCWGAVEADYANGGEVLNRLFKEYPDTLKLFPKFSGISQGDLAGSPA
VAAH GATVLKKLGELLKAKGDHAALLKPLANTHANIHKVALNNFRLITEVLVKVMAE
KAGLDAA GQGALRRVMDAVIGDIDGYYKEIGFAG

Разная длина, как сравнивать?

Последовательности миоглобинов человека и рыбы

>MYG_HUMAN

MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFHLKSEDEMKA
DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH
PGDFGADAQGAMNKALELFRKDMASNYKELGFQG

>MYG_DANRE

MADHDLVLKCWGAVEADYANGGEVLNRLFKEYPDTLKLFPKFSGISQGDLAGSPAVA
AHGATVLKKLGELLKAKGDHAALLKPLANTHANIHKVALNNFRLITEVLVKVMAEKAGLDAA
GQGALRRVMDAVIGDIDGYYKEIGFAG

Разная длина, как сравнивать?

Ответ: **выравнивание**

MYG_HUMAN	MGLSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFHLKSEDEMKA	60
MYG_DANRE	----MADHDLVLKCWGAVEADYANGGEVLNRLFKEYPDTLKLFPKFSGISQG-DLAGSP	55
	. : :***: ** *** . :* *** *** :*:***: * ***. . . . : . *	

MYG_HUMAN	DLKKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKH	120
MYG_DANRE	AVAAHGATVLKKLGELLKAKGDHAALLKPLANTHANIHKVALNNFRLITEVLVKVMAEKA	115
	: *****. ** ;** **.* * :*****:***. **: . : .:***: : :*: . *	

MYG_HUMAN	PGDFGADAQGAMNKALELFRKDMASNYKELGFQG	154
MYG_DANRE	--GLDAAGQGALRRVMDAVIGDIDGYYKEIGFAG	147
	. . * . ***: *: . ***:*** *	

Источники разнообразия геномов

1. Ошибки репликации ДНК

- Регулярное переписывание информации как способ сохранения имеет определенные преимущества.
Рукопись «Слова о полку Игореве» сгорела, но информация сохранилась.
- При переписывании бывают ошибки – мутации.
- ДНК любого ныне живущего организма получилась переписыванием ДНК организма, жившего примерно 3,5 млрд лет тому назад.
Этот организм называют **LUCA** (Last Universal Common Ancestor).
При переписывании случались и добавления, при этом источники «новой» ДНК довольно загадочны.
- Текст ДНК, конечно, изменился до неузнаваемости.
Но родство (гомологичность) последовательностей некоторых белков во всех современных организмах устанавливается достаточно надежно

Источники разнообразия геномов

2. Повреждения ДНК и их репарация

- Имеется много источников повреждений ДНК
 - ультрафиолетовое излучение
 - различные химические вещества, содержащиеся в пище, воздухе, табачном дыме...
 - некоторые ферменты самого организма
- Повреждения ДНК контролируются клеткой и **репарируются**
Увы, не всегда правильно. Одно из следствий – онкологические заболевания.

Гомологичные последовательности

>First

```
CGTCCCCGGGTCTTGTACACACCGCCCGTCACACCACGAGAGTTGTAACACCCGAAGGCCGG  
TGGAGTAACCATTGGAGCTAGCCGTCGAAGGTGGG
```

>Second

```
CGTCCCCGGGCCTTGTACACACCGCCCGTCACACCATGGAAGTCTGCAATGCCCAAAGTCGG  
TGGGATAACCTTATAAGGAGTCAGCCGCCTAAGGCAGG
```

Выравнивание (демонстрирует сходство)

Негомологичные последовательности

>First

```
CGTCCCCGGGTCTTGTACACACCGCCCGTCACACCACGAGAGTTGTAACACCCGAAGCCGG  
TGGAGTAACCATTGGAGCTAGCCGTCGAAGGTGGG
```

>Third

```
CCTGCCTTAGGCGGCTGACTCCTATAAAGGTTATCCCACCGACTTGGGCATTGCAGACTTC  
CATGGTGTGACGGCGGTGTACAAGGCCGGAACG
```

Выравнивание (бессмысленное)

```
>First
CGTCCCCGGGTCTTGTACACACCGCCCGTCACACCACGAGAGAGTTGTAACACCCGAAGGCCGG
TGGAGTAACCATTGGAGCTAGCCGTCGAAGGTGGG
>Third
CCTGCCTTAGGC GGCTGACTCCTATAAAGGTTATCCCACCGACTTGGCATTGCAGACTTC
CATGGTGTGACGGGCGGTGTGTACAAGGCCGGAACG
```

First	1	-----	CGTTCCCCGGGT	11
			
Third	1	CCTGCCTTAGGC GGCTGACTCCTATAAAGGTTATCCCACCGACTTGGC		50
First	12	CTTGTACACACCGCCCGTCACACCACGAGAGAGTTGTAACACCCGAAGCCG		61
			
Third	51	ATTGCAGACTTCCATGGTGTGACGGGCGGTGTGTACAAGGCCGGAACG		100

Выравнивание последовательностей потомков относительно предка

Гомологичные нуклеотиды ставим друг под другом

ПРЕДОК 1 . TATGCGAATGCCCTGAA

сын 2 . TATGCAATGCCCTGAA замена

Выравнивание последовательностей потомков относительно предка

Гомологичные нуклеотиды ставим друг под другом

ПРЕДОК 1 . TATGCGAAT-GCCCTGAA

сын 2 . TATGCAAAAT-GCCCTGAA замена

внук 3 . TATGCAAAT-GCTCTGAA замена

правнук 4 . TATGCAAATCGCTCGGAA вставка 1 п.н.

Выравнивание последовательностей потомков относительно предка

Гомологичные нуклеотиды ставим друг под другом

- | | | | |
|---------------|----|--|----------------|
| ПРЕДОК | 1. | TATGCGAAT-GCCCTGAA | |
| сын | 2. | TATGCA A AAT-GCCCTGAA | замена |
| внук | 3. | TATGCA A AAT-GC T CTGAA | замена |
| правнук | 4. | TATGCA A AAT C GCT T CGGAA | вставка 1 п.н. |
| праправнук | 5. | TATGCA A AAA A CGCT T CGGAA | замена |
| прапраправнук | 6. | TATGCA A AAA- C GCT T CGGAA | делеция 1п.н. |

Выравнивание последовательностей потомков относительно предка

Гомологичные нуклеотиды ставим друг под другом

- | | | | |
|---------------|----|--|----------------|
| ПРЕДОК | 1. | TAT--GCGAAT-GCCCTGAA | |
| сын | 2. | TAT--GC A AAT-GCCCTGAA | замена |
| внук | 3. | TAT--GC A AAT-GCT T CTGAA | замена |
| правнук | 4. | TAT--GC A AAT C GCT C GGAA | вставка 1 п.н. |
| праправнук | 5. | TAT--GC A AA A C G CT C GGAA | замена |
| прапраправнук | 6. | TAT--GC A AA- C GCT C GGAA | делеция 1п.н. |
| | 7. | TAT--GC A T A- C GCT C GGAA | замена |
| | 8. | TAT--GC A T A- C GC---GAA | делеция 3 п.н. |
| | 9. | TAT A TGC A T A- C GC---GAA | вставка 2 п.н. |

Выравнивание геномов двух потомков общего предка микоплазм: *M. capricolum* и *M. mycoides*, (маленький фрагмент)

1 <i>M. mycoides</i>	1091	t a a - - - t t a a t t a t a a a t t t a t a a a a t t t t c a t t a a G T C T G A	1130
1 <i>M. capricolum</i>	1116	T A A T T T T T A A T T A T A A A A T T T A T A A A A T T T T C A T T A A G T C T A A	1158
1 <i>M. mycoides</i>	1131	T G T A T T C A C C T T T T T T T A A T A T A T A A A A C T C C A G A A A G A A A A T C	1173
1 <i>M. capriculum</i>	1159	T A T A T T C A C C T T T T T T A A C A T A T A T A A A A C T C C A G A A A G A A A A T C	1201
1 <i>M. mycoides</i>	1174	T T T A A A A A C G T T T A G C T T T A T T A T C A T C T A A G T T T T T A A A A T C T	1216
1 <i>M. capriculum</i>	1202	T T T A A A A A C G T T T A G C T T T A T T A T C A T C T A A G T T T T T A A A A T C T	1244
1 <i>M. mycoides</i>	1217	A C A A C A A C A A C A T T T T T G A T C T A A T A A A G T A T C T A C A A T T G A T T	1259
1 <i>M. capriculum</i>	1245	A T A A C A A C A A C A S A T T A T G T T C T A A T A A A G T A T C A A C A A T T G A T T	1287
1 <i>M. mycoides</i>	1260	G A A C T T C A G A A A A T T T C A T A G G A C T A A A T A C A T A A G T G T T A A T	1302
1 <i>M. capriculum</i>	1288	G A A T T T C A G A A A A T T T C A T A G G A C T A A A A A C A T A T G T A T T A A C	1330

В среднем 92% совпадающих букв на **гомологичных** участках

Гены и белки

Геном

$3 \cdot 10^9$ букв у человека,
 $\sim 10^6$ букв у бактерий

содержит

Гены

<2% генома у человека,
 $\sim 90\%$ у бактерий

кодируют

Белки

$\sim 25\ 000$ у человека,
600 – 6000 у бактерий

Генетический код

	T(U)	C	A	G	
T(U)	TTT Phe	TCT Ser	TAT Tyr	TGT Cys	
	TTC Phe	TCC Ser	TAC Tyr	TGC Cys	
	TTA Leu	TCA Ser	TAA Stop	TGA Stop	
	TTG Leu	TCG Ser	TAG Stop	TGG Trp	
C	CTT Leu	CCT Pro	CAT His	CGT Arg	
	CTC Leu	CCC Pro	CAC His	CGC Arg	
	CTA Leu	CCA Pro	CAA Gln	CGA Arg	
	CTG Leu	CCG Pro	CAG Gln	CGG Arg	
A	ATT Ile	ACT Thr	AAT Asn	AGT Ser	
	ATC Ile	ACC Thr	AAC Asn	AGC Ser	
	ATA Ile	ACA Thr	AAA Lys	AGA Arg	
	ATG Met	ACG Thr	AAG Lys	AGG Arg	
G	GTT Val	GCT Ala	GAT Asp	GGT Gly	
	GTC Val	GCC Ala	GAC Asp	GGC Gly	
	GTA Val	GCA Ala	GAA Glu	GGA Gly	
	GTG Val	GCG Ala	GAG Glu	GGG Gly	

Аминокислоты

A Ala	Alanine	Аланин
R Arg	Arginine	Аргинин
N Asn	Asparagine	Аспарагин
D Asp	Aspartic Acid	Аспарагиновая кислота
C Cys	Cysteine	Цистеин
Q Gln	Glutamine	Глютамин
E Glu	Glutamic Acid	Глутаминовая кислота
G Gly	Glycine	Глицин
H His	Histidine	Гистидин
I Ile	Isoleucine	Изолейцин
L Leu	Leucine	Лейцин
K Lys	Lysine	Лизин
M Met	Methionine	Метионин
F Phe	Phenylalanine	Фенилаланин
P Pro	Proline	Пролин
S Ser	Serine	Серин
T Thr	Threonine	Треонин
W Trp	Trryptophan	Триптофан
Y Tyr	Tyrosine	Тирозин
V Val	Valine	Валин
<i>"Stop"</i> в таблице кода означает стоп-кодон – сигнал окончания трансляции.		

Мутации

gatcaacactacttgacttcaagacttaccataaagaaaac



точечная замена

gatcaacactacttgacttcaa~~a~~acttaccataaagaaaac

gatcaacactacttgacttcaag~~a~~acttaccataaagaaaac



делеция

gatcaacactacttgacttcaacttaccataaagaaaac

gatcaacactacttgacttcaagacttaccataaagaaaac



инсерция
(вставка)

gatcaacactacttgacttcaaga~~t~~acttaccataaagaaaac

Классификация мутаций в кодирующих последовательностях ДНК

Синонимическая или молчащая: не меняет кодируемый аминокислотный остаток

Миссенс (missense): меняет остаток

Нонсенс (nonsense): заменяет кодон остатка на стоп-кодон

Сдвиг рамки (frameshift): вставка или делеция размера, не кратного трём.

Результатом являются совсем другие аминокислоты после мутации и, как правило, стоп-кодон сравнительно недалеко (в среднем через 21 триплет после мутации)

Точечные замены в гене

... AATCCGTCAAGTCTA...

... Asn Pro Ser Ser Leu ...

1) "молчащая"(синонимическая)мутация

... AATCCGTCGAGTCTA...

... Asn Pro Ser Ser Leu ...

2) замена остатка на близкий по свойствам

... AATCCGACAAGTCTA...

... Asn Pro **Thr** Ser Leu ...

3) замена остатка на остаток с иными свойствами

... AATCCGTCAAGACTA...

... Asn Pro Ser **Arg** Leu ...

Судьба мутации

Бактерия разделилась, и у одного из потомков произошла мутация.
(ошибка репликации, или повреждение ДНК и ошибка репарации).

Что будет с потомством мутанта? Увидим ли мы эту мутацию, если отсеквенируем 1 000 000 бактерий этого штамма через 10 лет?

Потомство бактерии

В благоприятных условиях бактерия может делиться каждый час.

Сколько будет бактерий через 24 часа? А через год????

Потомство бактерии

В благоприятных условиях бактерия может делиться каждый час.

Сколько будет бактерий через 24 часа? А через год????

Ответ: примерно столько же, сколько сейчас.

Потомство бактерии

В благоприятных условиях бактерия может делиться каждый час.

Сколько будет бактерий через 24 часа? А через год????

Ответ: примерно столько же, сколько сейчас.

Численность подавляющего большинства популяций **постоянна** (по крайней мере на отрезках времени порядка лет) – погибает примерно столько же, сколько рождается.
Современная популяция человека – исключение!

Если члены популяции генетически идентичны, то вероятность оставить потомство для всех **одинакова** (точнее, зависит от только от внешних факторов).

Следствие: математическое ожидание числа потомков одной бактерии через достаточно большой промежуток времени равно 1.

Судьба нейтральной мутации

Предположим, что мутация **нейтральна** = никак не влияет на матожидание числа потомков (таких мутаций довольно много).

Мутация произошла и передаётся потомкам мутанта. Значит, в популяции появился новый **полиморфизм**. У данного варианта кода есть **частота** (сначала очень маленькая).

Судьба нейтральной мутации

Предположим, что мутация **нейтральна** = никак не влияет на матожидание числа потомков (таких мутаций довольно много).

Мутация произошла и передаётся потомкам мутанта. Значит, в популяции появился новый **полиморфизм**. У данного варианта генома есть **частота** (сначала очень маленькая).

Что произойдёт с частотой через пару суток?

Судьба нейтральной мутации

Предположим, что мутация **нейтральна** = никак не влияет на матожидание числа потомков (таких мутаций довольно много).

Мутация произошла и передаётся потомкам мутанта. Значит, в популяции появился новый **полиморфизм**. У данного варианта генома есть **частота** (сначала очень маленькая).

Что произойдёт с частотой через пару суток?

Ответ: частота либо немного возрастёт, либо немного упадёт. То и другое примерно равновероятно.

Случайное блуждание

Частота любого нейтрального полиморфизма постоянно колеблется случайным образом (это называется «генетический дрейф»).

Математическая модель такого процесса называется «случайное блуждание».

На тротуаре стоит пьяный и каждые 10 сек. делает шаг либо направо, либо налево, случайно выбирая направление. Как далеко он уйдёт за время T ?

Ответ: в среднем на расстояние, пропорциональное \sqrt{T} .

Случайное блуждание с поглощением

По длинной дамбе идёт пьяный и с каждым шагом отклоняется либо на полметра вправо, либо на полметра влево. Как скоро он свалится с дамбы?

Ответ: скоро...

Когда частота генетического варианта достигает 100% или 0%, процесс её изменения прекращается.

За исторически короткое время любой нейтральный вариант либо исчезает из популяции, либо закрепляется в ней!

Закрепление мутаций как результат генетического дрейфа

Вероятность закрепиться для новой нейтральной мутации очень мала, но не 0.

Организмов в популяции много, мутаций в них происходит тоже много (примерно 10^{-8} на п.н. на поколение – каждая сотая новорождённая бактерия несёт новую мутацию). Значительная доля мутаций нейтральна.

Итог: геномы независимых популяций начинают различаться, чем дальше, тем больше – в них независимо накапливаются нейтральные мутации.

А если мутация не нейтральна?

Каждому варианту генома можно сопоставить его «приспособленность» f = матожидание числа потомков организма с таким геномом (через какой-то фиксированный промежуток времени).

В подавляющем большинстве случаев новая мутация порождает либо нейтральный вариант ($f = 1$) либо вредный ($f < 1$).

Вредный вариант тоже начинает «блуждать», но вероятность «шага вверх» оказывается меньше вероятности «шага вниз». Это очень сильно уменьшает вероятность закрепления – тем сильнее, чем меньше f , и тем сильнее, чем больше популяция.

Явление невозможности закрепления вредной мутации называется **стабилизирующий отбор** или же **отрицательный отбор**.

Положительный отбор

Если вдруг $f > 1$, то вероятность закрепления мутации вырастает во много раз.
Процесс закрепления полезных мутаций называется **положительным отбором**.

Собственно, полезных мутаций так мало именно потому, что большинство возможных полезных мутаций уже закрепились.

Обычно полезные мутации начинают появляться в заметном количестве только при изменении условий жизни организмов – например при появлении нового источника пищи или новой опасности или попадании части популяции в другой климат...

Эволюция белков

Мутации возникают случайно.

Конкретная мутация может быть:

- летальной;
- вредной;
- слабовредной;
- нейтральной;
- полезной.

Мутация порождает **полиморфизм данного белка в популяции**.

Доля каждого варианта подвержена случайным изменениям (модель: «случайное блуждание с поглощением»).

За исторически короткое время один из вариантов (старый или новый) исчезает. В первом случае говорят, что мутация **закрепилась**.

Как правило, пространственная структура белка почти не меняется при эволюции его последовательности. В первом приближении верно утверждение: **гомологичные белки имеют почти одинаковые 3D-структуры**.

Множественное выравнивание белков

11 DRRE I RH IWDDWSSSFTDRRVA I VRAVFDDLFKHYPTSKALFERVK I DEPESGEF	66
8 DRHEVLDNWKG I WSAEFTGRRVA I GQA I FQELFALDPNAKGVFGRVNVDKPSEADW	63
8 DRREVQALWRS I WSAEDTGRRTL I GRLLFEEELFE IDGATKGLFKRVNVDDTHSPEE	63
7 QR I KVKKQQWAQVYSGES - - RTDFAIDVFNNFFRTNPD - RSLFNRVNGDNVYSPEF	59
9 QRLKVKQQWAKAYVGHE - - RVELGIALWKSMFAQDNDARDLFKRVHGEDVHSPAF	62
8 EGLKVKSEWGRAYGSGHD - - REAFSQA I WRATFAQVPESRSLFKRVHGDDTSHPAF	61
6 QRFKVKHQWAEAFGTSHH - - RLDFGLKLWNS I FRDAPE I RGLFKRVDGDNAYSAEF	59
7 QRLKVKRQWAEAYGSGND - - REEFGHF IWTHVFKDAPSARDLFKRVRGDN I HTPAF	60
67 KSHLVRVANGDLL I NLLDDTL VLQSHLGHLADQH I QRKGVTKEYFRG I GEAFA	120
64 KAHVIRV I NGLDLAVNL LEDPKALQEELKHLARQHRERSGVKAVYFDEMEKALL	117
64 FAHVLRVVNGLDTL I GVLGDSDTLNSL IDHHLAEQHKARAGFKTVYFKEFGKALN	117
60 KAHMVRVFAGFD I L I SVLDDKPVL DQALAHYAAFHKQF - GT I P - - FKAFGQTMF	110
63 EAHAMARVFNGLDRV I SSLTDEPVLAQLEHLRQQH I KL - G I TGHMFNL MRTGLA	115
62 I AHAERVLGGLD I A I STL DQPATLKEELDHLQVQHEGR - K I PDNYFDAFKTA I L	114
60 EAHAERVLGGLDMT I SLLDDQAAFDAQLAHLKSQHAER - N I KADYYGVFVNELL	112
61 RAHATRVLGGLDMC I ALLDDEGV LNTQLAHLASQHSSR - GVSAQYDVVEHSVM	113

Мы видим только закрепившиеся мутации!

Гомология – общность происхождения

- При репликации почти всегда каждый нуклеотид потомка происходит от определенного нуклеотида предка.
- В выравнивании гомологичных последовательностей у разных потомков одного и того же предка гомологичные нуклеотиды должны стоять в одной колонке.
- Как правило, нам известны геномы только современных организмов, и потому у нас нет способа проверить, какие нуклеотиды гомологичны.
- Гомологичность последовательностей часто можно установить анализом их выравнивания.
- Проблема построения выравнивания обсуждается ниже.

Выравнивание последовательностей касается всех студентов МГУ!

Положение об обеспечении самостоятельности выполнения письменных работ в МГУ имени М.В.Ломоносова на основе системы «Антиплагиат»

Самостоятельное выполнение письменных работ обучающимися в МГУ имени М.В.Ломоносова (далее – МГУ) является необходимым условием эффективности этих работ как элементов учебного процесса, развития у обучающихся навыков научной работы.

К обучающимся в Университете относятся студенты, аспиранты, докторанты, слушатели и соискатели (ст.ст. 123-128 Устава МГУ).

Для данных двух последовательностей существует много разных выравниваний

TGGAGTAACCAT-
TGGGATAACCTTG

TGGAGTAACCAT-----
-----TGGGATAACCTTG

-TGGAGTAACCAT
TGGGATAACCTTG

TGGA--GTAACCAT--
TGGGATAA---CCTTG

Всего для двух последовательностей одинаковой длины p имеется 2^n разных выравниваний

Биоинформатическая задача: выбрать среди множества выравниваний правильное

Алгоритм выравнивания решает математическую задачу, а не биологическую

Математическая задача разбивается на две:

- Любому выравниванию сопоставить число – его **вес**
- Для данных последовательностей построить выравнивание с наибольшим весом

Три понимания «правильного» выравнивания

1

Оптимальное выравнивание: наилучшее по весу

Его ищут программы.

Оптимальное выравнивание существует для любого набора последовательностей, даже негомологичных!

2

Эволюционное выравнивание: запись, отражающая ход эволюции

Не поддается достоверной реконструкции в большинстве реальных случаев; может отличаться от оптимального выравнивания.

Алгоритм вычисления веса стараются выбрать так, чтобы можно было ожидать, что эволюционное выравнивание будет среди нескольких оптимальных.

Для негомологичных последовательностей эволюционного выравнивания не существует!

3

Функциональное выравнивание: сопоставление функционально

идентичных частей белков или нуклеиновых кислот

Объясняет сохранение в эволюции одних частей белка и варьирование других.

Поскольку функция и 3D-структура белка очень тесно связаны, функционально выровненные аминокислотные остатки должны иметь примерно одинаковое расположение в пространстве.

Вес парного выравнивания

Простейший вариант

За каждую колонку с совпадающими буквами прибавляем число A

За каждую колонку с разными буквами вычитаем число B

За каждую «чёрточку» (гэп) вычитаем число C

Вес парного выравнивания

Простейший вариант

За каждую колонку с совпадающими буквами прибавляем число A

За каждую колонку с разными буквами вычитаем число B

За каждую «чёрточку» (гэп) вычитаем число C

First TGGAGTAACCAT--TAGGAGCTAGCCG

 ||||..|||||.|| ||||||..|||||

Second TGGGATAACCTTGATAGGAGTCAGCCG

Здесь 20 совпадений, 5 несовпадений, два гэпа, значит вес $20A - 5B - 2C$

Например, при $A = 5$, $B = 4$, $C = 6$ вес равен 68.

Проверьте, что ни у какого выравнивания этих последовательностей вес не будет большим. Это выравнивание является оптимальным при данных параметрах A , B , C .

Вес парного выравнивания – белки

Какое выравнивание имеет большие шансы оказаться правильным?

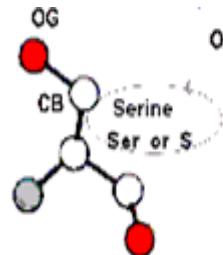
AFTGAHAYL

AYS---AYM

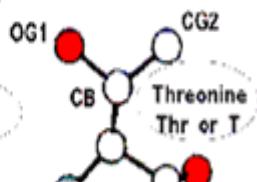
AFTGAHAYL

AY---SAYM

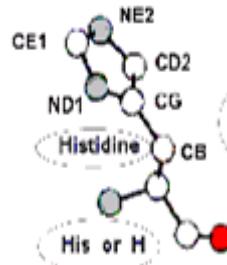
Вес парного выравнивания – белки



Серин S



Треонин T



Гистидин H

Мутация серина в треонин закрепляется с гораздо большей вероятностью по сравнению с мутацией серина в гистидин.

Поэтому если в одной колонке выравнивания оказались буквы S и T, это скорее аргумент за данное выравнивание, чем против него. Значит, за такую колонку лучше увеличивать вес, чем уменьшать.

Матрица весов аминокислотных замен BLOSUM 62

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-3	0	-2	-2	0	6														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	2	5											
Q	-3	0	-1	-1	-1	-2	0	0	2	5										
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8									
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5								
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5						
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	1	4						
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4				
V	-1	-2	0	-2	0	-3	-3	-3	-2	-3	-3	-2	1	3	1	4				
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7	
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

Треугольная (симметричная)
матрица

Из работы (Henikoff&Henikoff, 1992, PNAS)

Матрица BLOSUM62

```

# Matrix made by matblas from blosum62.iij
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units
# Blocks Database = /data/blocks_5.0/blocks.dat
# Cluster Percentage: >= 62
# Entropy = 0.6979, Expected = -0.5209
A R N D C Q E G H I L K M F P S T W Y V B Z X *
A 4 -1 -2 -2 0 -1 -1 0 -2 -1 -1 -1 -1 -2 -1 1 0 -3 -2 0 -2 -1 0 -4
R -1 5 0 -2 -3 1 0 -2 0 -3 -2 2 -1 -3 -2 -1 -1 -3 -2 -3 -1 0 -1 -4
N -2 0 6 1 -3 0 0 0 1 -3 -3 0 -2 -3 -2 1 0 -4 -2 -3 3 0 -1 -4
D -2 -2 1 6 -3 0 2 -1 -1 -3 -4 -1 -3 -3 -1 0 -1 -4 -3 -3 4 1 -1 -4
C 0 -3 -3 -3 9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
Q -1 1 0 0 -3 5 2 -2 0 -3 -2 1 0 -3 -1 0 -1 -2 -1 -2 0 3 -1 -4
E -1 0 0 2 -4 2 5 -2 0 -3 -3 1 -2 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
G 0 -2 0 -1 -3 -2 -2 6 -2 -4 -4 -2 -3 -3 -2 0 -2 -2 -3 -3 -1 -2 -1 -4
H -2 0 1 -1 -3 0 0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2 -3 0 0 -1 -4
I -1 -3 -3 -3 -1 -3 -3 -4 -3 4 2 -3 1 0 -3 -2 -1 -3 -1 3 -3 -3 -1 -4
L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -2 -1 1 -4 -3 -1 -4
K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 5 -1 -3 -1 0 -1 -3 -2 -2 0 1 -1 -4
M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5 0 -2 -1 -1 -1 -1 1 -3 -1 -1 -4
F -2 -3 -3 -3 -2 -3 -3 -1 0 0 -3 0 6 -4 -2 -2 1 3 -1 -3 -3 -1 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 7 -1 -1 -4 -3 -2 -2 -1 -2 -4
S 1 -1 1 0 -1 0 0 -1 -2 -2 0 -1 -2 -1 4 1 -3 -2 -2 0 0 0 -4
T 0 -1 0 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 5 -2 -2 0 -1 -1 0 -4
W -3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1 1 -4 -3 -2 11 2 -3 -4 -3 -2 -4
Y -2 -2 -2 -3 -2 -1 -2 -3 2 -1 -1 -2 -1 3 -3 -2 -2 2 7 -1 -3 -2 -1 -4
V 0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4 -3 -2 -1 -4
B -2 -1 3 4 -3 0 1 -1 0 -3 -4 0 -3 -3 -2 0 -1 -4 -3 -3 4 1 -1 -4
Z -1 0 0 1 -3 3 4 -2 0 -3 -3 1 -1 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
X 0 -1 -1 -1 -2 -1 -1 -1 -1 -1 -1 -1 -1 -2 0 0 -2 -1 -1 -1 -1 -1 -1 -4
* -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 1

```

Вес парного выравнивания – белки

Посчитаем веса выравниваний, используя матрицу BLOSUM62

AFTGAHAYL
AYS---AYM

AFTGAHAYL
AY---SAYM

Обозначим значения матрицы на пересечении строки A и столбца B через $M(A,B)$

Тогда вес левого выравнивания равен:

$M(A,A) + M(F,Y) + M(T,S) + M(A,A) + M(Y,Y) + M(L,M)$ – штраф за гэпы,

а правого:

$M(A,A) + M(F,Y) + M(H,S) + M(A,A) + M(Y,Y) + M(L,M)$ – штраф за гэпы.

Штрафы за гэпы одинаковы, значит веса различаются слагаемым $M(T,S)$ слева против $M(H,S)$ справа. Но $M(T,S) = 1$, а $M(H,S) = -1$, поэтому вес левого выравнивания больше на 2.

Вес парного выравнивания: аффинные штрафы за гэпы

First TGGAGTAACCAT--TTGGAGCTAGCCG
 |||..|||||. | .||||..|||||

Выравнивание 1

Second TGGGATAAACCTTATAGGAGTCAGCCG

First TGGAGTAACCAT-TT-GGAGCTAGCCG
 |||..|||||. | .| ||||..|||||

Выравнивание 2

Second TGGGATAAACCTTATAGGAGTCAGCCG

Вес парного выравнивания: аффинные штрафы за гэпы

First TGGAGTAACCAT--TTGGAGCTAGCCG
 |||..|||||. | .||||..|||||

Выравнивание 1

Second TGGGATAAACCTTATAGGAGTCAGCCG

First TGGAGTAACCAT-TT-GGAGCTAGCCG
 |||..|||||. | .| ||||..|||||

Выравнивание 2

Second TGGGATAAACCTTATAGGAGTCAGCCG

Выравнивание 1 биологически более вероятно, чем выравнивание 2
(потому что одна делеция в две буквы случается чаще, чем две делеции в одну букву)

Чтобы выравнивание 1 имело больший вес, чем выравнивание 2, штрафы за гэпы делают зависимым от числа подряд идущих гэпов. Стандартный способ: за первый гэп вычитается «штраф за открытие», за каждый последующий — меньший «штраф за удлинение»

Терминология: гэпы и индели

Один знак "-", означающий отсутствие в данной последовательности **одной** буквы, гомологичной другим буквам данного столбца, мы будем называть «гэп»

Совокупность нескольких подряд идущих гэпов мы будем называть «индель», от **инсерция/делеция**.

First TGGAGTAACCAT--TTGGAGCTAGCCG
 |||..|||||. | .||||..|||||
Second TGGGATAACCTTATAGGAGTCAGCCG

Тут два гэпа и один индель

К сожалению, терминология не вполне устоялась. В литературе и описаниях программ вы можете встретить употребление термина «гэп» для обозначения инделя.

Парное выравнивание: локальное и глобальное

На самом деле это две формализации одной и той же задачи: даны последовательности двух белков, найти гомологичные аминокислотные остатки.

Задача глобального выравнивания: найти выравнивание с наибольшим весом. При вычислении весов учитываются: матрица замен (BLOSUM62) и аффинные штрафы за гэпы.

Задача локального выравнивания: найти

- участок в первой последовательности;
- участок во второй последовательности;
- выравнивание выбранных участков;

так, чтобы вес выравнивания был наибольшим.

Разница в том, что теперь выбираем не только как выравнивать, но и что.

Для самостоятельного обдумывания

1. (формальный вопрос). Если не штрафовать гэпы до первого сопоставления и после последнего, то глобальное выравнивание сведётся к локальному. Иначе говоря, задачу локального выравнивания можно сформулировать так: найти выравнивание с наибольшим весом, но при этом вес считаем хитрее: штрафуем только те гэпы, которые оказались между сопоставлениями букв.
2. (содержательный вопрос). Почему локальное выравнивание довольно часто имеет больший биологический смысл по сравнению с глобальным?

Форматы хранения выравниваний

Fasta-формат

```
>CHICK
MVGSSSEAGGEAWRGRRYYRLEEVQKHNNSQSTWIIVHHRIYDITKFLDEHPGEEVLREQAGGDATENFEDVG
HSTDARALSETFIIGELH-PDDRPKLQK--PAETLITTQSNSSWSN---WVIP-AIAAIIVALMYRSYMS
E-
>HUMAN
---MAEQSDEAVK--YYTLEEIQKHNHSKSTWLILHHKVYDLTKFLEEHPGEEVLREQAGGDATENFEDVG
HSTDAREMSKTFIIGELH-PDDRPKLNK--PPETLITTIDSSSSWWTN---WVIP-AISAVAVALMYRLYMA
ED
>CUSRE
-----MGGSKV----YSLAEVSEHSQPNDCWLVIGGKVYDVTKFLDDHPGGADVLLSSTAKDATDDFEDIG
HSSSARAMMDEMCGDID-SSTIPTKTSYTPPKQPLYNQDKTPQFIIKLLQFLVPLIILGVAVGIRFYKKQS
SD
```

Aln-формат (он же Clustal)

CHICK	MVGSSSEAGGEAWRGRRYYRLEEVQKHNNSQSTWIIVHHRIYDITKFLDEHPGEEVLREQA
HUMAN	---MAEQSDEAVK--YYTLEEIQKHNHSKSTWLILHHKVYDLTKFLEEHPGEEVLREQA
CUSRE	-----MGGSKV----YSLAEVSEHSQPNDCWLVIGGKVYDVTKFLDDHPGGADVLLSST
CHICK	GGDATENFEDVGHSTDARALSETFIIGELH-PDDRPKLQK--PAETLITTQSNSSWSN
HUMAN	GGDATENFEDVGHSTDAREMSKTFIIGELH-PDDRPKLNK--PPETLITTIDSSSSWWTN
CUSRE	AKDATDDFEDIGHSSSARAMMDEMCGDID-SSTIPTKTSYTPPKQPLYNQDKTPQFIIK
CHICK	---WVIP-AIAAIIVALMYRSYMSE-
HUMAN	---WVIP-AISAVAVALMYRLYMAED
CUSRE	LLQFLVPLIILGVAVGIRFYKKQSSD

Программы

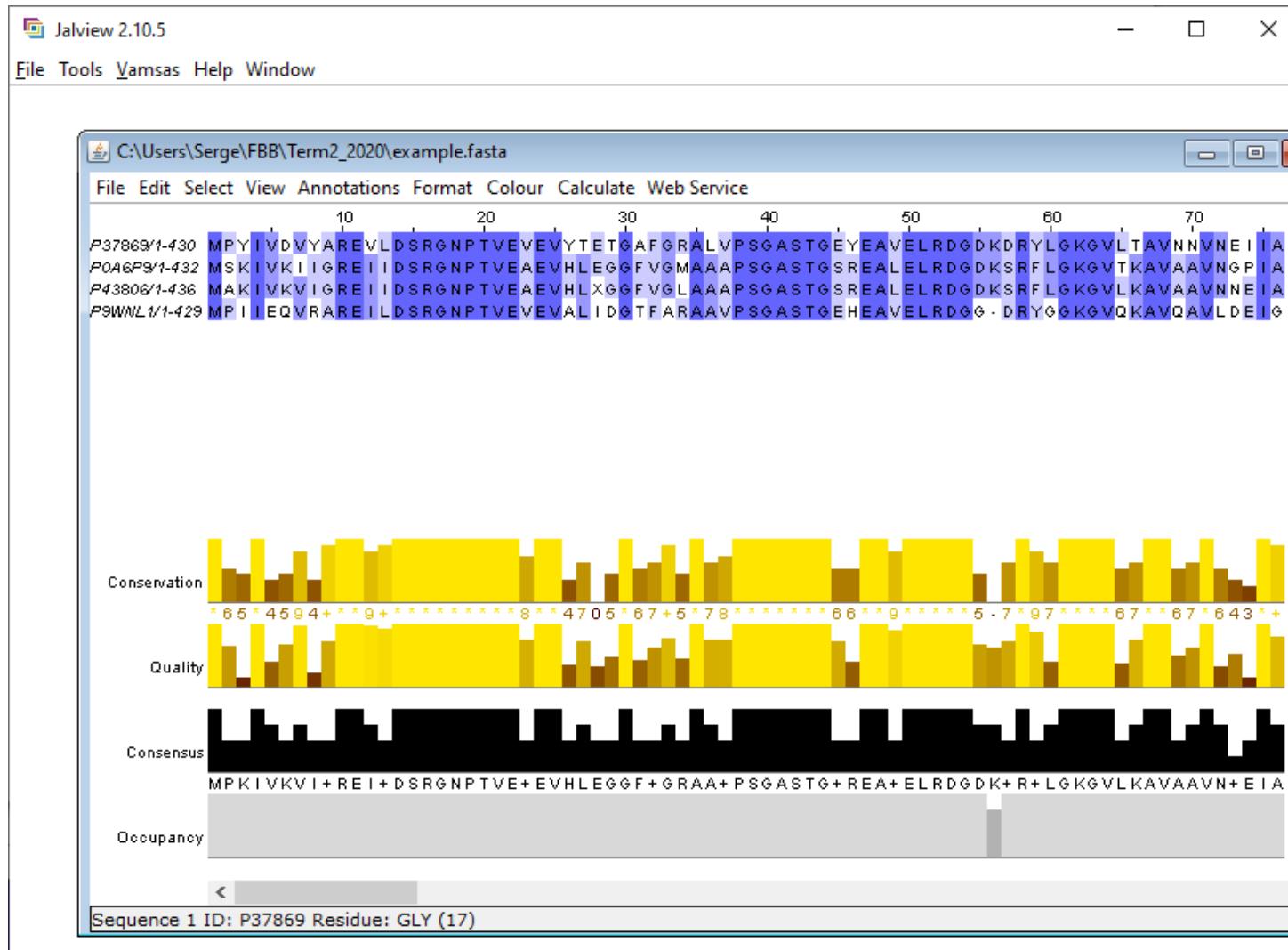
Парное глобальное выравнивание: needle, stretcher

Парное локальное выравнивание: water, matcher, blast

Множественное выравнивание: Muscle, MAFFT, Pride, ...

Редакторы выравниваний: JalView, GeneDoc, ...

Jalview



Словарик

Alignment	Выравнивание
Gap	Гэп
Indel	Индель
Gap penalty	Штраф за гэпы
Gap opening penalty	Штраф за открытие гэпа
Gap extension penalty	Штраф за удлинение гэпа
Score	Вес выравнивания
Scoring matrix	Матрица замен аминокислот

Вопросы и ответы

Что такое гомология?

Ответ: общность происхождения

(*НЕПРАВИЛЬНО говорить «последовательности гомологичны на 56%. Последовательности либо гомологичны, либо нет*)

Как определить, гомологичны ли белки?

Ответ: в большинстве случаев единственный способ — выровнять их последовательности и посмотреть на процент совпадающих букв. Если он достаточно велик, то белки, вероятно, гомологичны. Если нет, то всякое может быть. Если для обоих белков известны пространственные структуры, то есть гораздо более чувствительный способ: сравнить ход полипептидной цепи.

Какой процент идентичности служит надёжным признаком гомологии?

Ответ: для белков обычно более 20–25% на достаточно длинном участке (более точный ответ будет дан на следующей лекции)