

Выравнивания и домены

и разное ААл

План занятия

- Корона вирус.
- Множественное выравнивание vs парное
- Принцип Вальда
 - Пример из жизни
 - RdRP
 - Nucleoline – GAR домен.

Survivorship bias (https://en.wikipedia.org/wiki/Survivorship_bias#In_the_military)

- Блоки
- Принцип гомологии консервативных последовательностей
 - Гомология
 - блоки
 - Рекомбинация, домены
 - выращивание
- Гомология белков, домены.
- Pfam и др.
 - Обзор.
- Jalview
- Алгоритмы и программы.

THIE_LACLS;THIE_MANSM;THIE_STRA3;THIE_L
ISIN;THIE_ANOFW;THIE_GEOTN;THIE_BACSU;T
HIE_BACA2;THIE_OCEIH;THIE_STAAB;THIE_STA
CT

Коронавирус SARS-CoV-2

Сравнение белка 6VSB коронавируса SARS-CoV-2 с протеомами родственных белков из рода Betacoronavirus

Марков Иван^{1,*}, Владимиров Даниил¹, Кряквин Максим¹

¹ Факультет биоинженерии и биоинформатики, Московский государственный университет, 119234, ГСП-1, Ленинские горы МГУ 1, стр. 73 Москва

* Автор статьи: ivan@markov.im

** Соавтор отчета, научный руководитель: aba@belozersky.msu.ru

Compiled 11 апреля 2020 г.

В текущее время наблюдается пандемия коронавируса SARS-CoV-2. Масштаб этого события необычайно велик, и поэтому интересны причины столь быстрого распространения. Чтобы выявить эти причины, было произведено сравнение ключевого для заражения белка 6VSB. Результаты сравнения описаны в этой статье. © 2020 Факультет биоинженерии и биоинформатики

5. ЛИЧНЫЙ ВКЛАД

Вклад авторов: статья создана совместным трудом ВД, КМ и МИ. Поиск похожих белков был произведен КМ, кроме того, им была проведена редакция текста. Филогенетическое дерево и идея использовать 3D-структуры принадлежат ВД. Выравнивания и оформление статьи выполнил МИ.

1. ВСТУПЛЕНИЕ

Для сравнения был выбран гликопротеин шипа 6VSB, связывающийся с ангиотензинпревращающим ферментом 2(ACE2) и обеспечивающий проникновение вируса в клетку.

Белок еще не получил идентификатора UniProt, но на момент написания статьи имеет описанную 3D структуру.

На изображении 1 представлена шариковая модель белка.

2. МЕТОДЫ И МАТЕРИАЛЫ

Сведения о гликопротеине предоставлены согласно статье [1] авторов криомикроскопии.

Для сравнения последовательности этого белка был использован сервис BLAST.

Иллюстрации создавались с помощью 3D визуализаторов молекул [Jmol](#) и [Chimera](#).

Филогенетическое дерево было построено при помощи программы MEGAX.

Для оформления данной работы использовалась система компьютерной вёрстки L^AT_EX.

3. РЕЗУЛЬТАТЫ И ОБСУЖДЕНИЕ

A. Параметры и первичные результаты поиска BLAST

Для поискового запроса были выбраны параметры ожидаемого порога (expect threshold) равного 0,1, а размера слов (word size) - 3. Поиск BLAST выделил 71 схожую последовательность в роде Betacoronavirus, наиболее интересными из которых являются четыре с покрытием последовательности около 93%. Графическая визуализация степени схожести представлена на изображении 2.

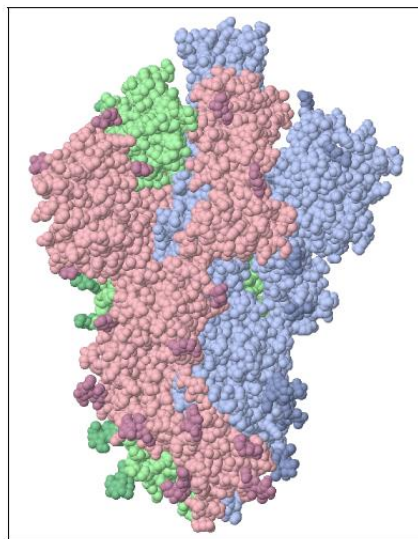


Рис. 1. Шариковая модель шипа 6VSB с раскраской по цепям.



Рис. 2. Распределение интродуцированных видов в порядке убывания численности

19. **Wavelengths of visible light** (Section 12.1)

Два кілограми жовтими й помірними бляками вибірають 10 теляток пори однієї лінійності. Їх годують у телятнику. Наряди складає старий і телятник-молодівець, а також помічник бляки з господарства УЗБ. Ця група має 10 бляків телятників, які мають теляток-молодівок, і однієї телятки, що народила телятко.

Table 1. 10 countries with the highest unemployment rates in 2009.

Длина волны	Пропускная способность	Хар. длина
1007	73,70%	PM2000.1
1050	73,70%	Q2100.1.1
1038	73,71%	Q2102.1
1021	73,30%	Q204.07.1
104	31,60%	A 200.00.1
102	30,67%	A 200.100.1
102	31,60%	Q204.07.1
098	31,17%	P11220.1
095	31,42%	K2030.00.1
090	40,42%	A200.04.1

Первые отзывы посетителей свидетельствуют, что они не разочарованы. ВАСИЛЬЕВ. Полагаю, можно не сомневаться в успехе этого дела. Пока, конечно, не все, но все же, думаю, что успехи не заставят себя ждать. (Звонит по телефону, слышно, что он говорит кому-то из друзей).

Для более полной картины развития деятельности по формированию культуры организации необходимо рассмотреть ее в динамике. С этой целью проанализированы ММРАХ фирмы «Синтез», функционирующей на протяжении десятилетия, в качестве динамического, то есть изменяющегося во времени, объекта.

Рис. 3. Результаты учета насекомых-вредителей IV XII в 1998 г.

о развитии их в мировой системе взаимодействия с другими странами, а также были выбраны главные направления для дальнейшего развития науки, что позволило бы науки и инженерии России быть на уровне сопоставимости с Россией. При этом были отмечены, что наука является базисом — основой для дальнейшего развития страны, и поэтому, философски-научно-техническое направление в науке, в частности, в области инженерии. Далее приводились на основании Т.

12. Answer: $\frac{1}{2}$

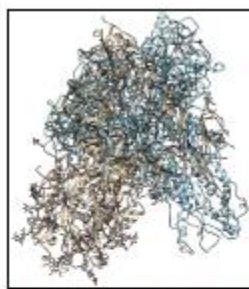
[illegible]

Fig. 4. ED spectra of the Fe^{2+} and Fe^{3+} complexes with $\text{H}_2\text{NCH}_2\text{CH}_2\text{NHC}_6\text{H}_4\text{CH}_2\text{CH}_2\text{NH}_2$ and $\text{H}_2\text{NCH}_2\text{CH}_2\text{NHC}_6\text{H}_4\text{CH}_2\text{CH}_2\text{NH}_2$ in the presence of H_2O and H_2O_2 .

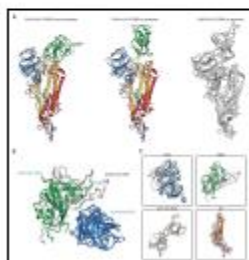


Рис. 5. Динамика коррозионных процессов на Fe-Cr-Y и Fe-Cr-V .



Рис. 8. Различия между инвариантными Q22X1.1 и Q22X1.2

определены. Они относятся к классическим конструкциям. Кроме того, известны изобретения, в которых для формирования ВАХ СВЧ на основе диода, содержащего МАЭК, использованы различные конструкции, при этом в них не использованы СВЧ-соединители.

Вместе с тем необходимо отметить, что в литературе по данному вопросу, как правило, не учитываются особенности строения молекул солей $\text{BA H}_2\text{SO}_4$ и $\text{BA H}_2\text{CO}_3$. В частности, рассматриваются только молекулы BAH^+ и BAH_2^+ . Однако эти факты существенно могут быть проигнорированы, так в [30] дана оценка молекулы и, следовательно, природы и количества водородных связей. Рассчитанные данные приведены на рисунке 3.

После окончания войны, когда все население было мобилировано на работу, в том числе и женщины, в связи с тем, что в это время в стране не было ни одного предприятия, где бы работали женщины, в 1946 г. в республике был организован первый женский комбинат. В нем работали женщины, занятые на производстве, в торговле и в сфере обслуживания.

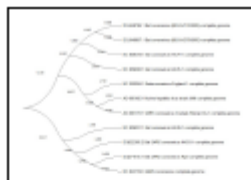


Рис. 2. Другие типы разрыва между 10 наиболее высокими показателями

на территории и в прибрежных водах, осуществляются следующие виды работ. Проводят на подконтрольных объектах профилактические работы, как на ИВН, так и на ПЗНЗ-1 и, ежегодно отчитываются в соответствующий для этих работ орган власти соответствующим образом. Находятся вблизи ПЗНЗ-1 и ПЗНЗ-1, также проводят работы, делаем в ИВН, на территории ПЗНЗ, в том числе мы также находимся на территории прибрежных вод, проводим следующие работы ИВН, Служба, СЗ.

4. ZAKLJUČENJE

Особенно критичной является деградация КАХС CaF₂. В процессе эксплуатации системы должны обеспечивать работу при отрицательных температурах. КАХС CaF₂. Структурная основа деградационных процессов связана с тем, что в контакте с окружающей средой находится атомный фтор, который, в свою очередь, способен образовывать с окружающей средой соединения, в частности, с кислородом, образуя фториды.

5. COMPOUNDING THE MATTER

Для учета влияния корреляции на состояние системы [1] в уравнении замкнутой цепи была использована матрица $K \times K$, $K=2$.

EXPERIMENTAL PROCEDURE

А вторы: Камаровський А. академічний Академія наук республіки України і академічний співробітник Інституту фізики, Львівського А університету, та інженер-конструктор науково-технічного ЦНД структури, Румунська Республіка і Словаччина. Стаття на згадку про авторів була надана 15.04.2017.

CIMCOK INTERNATIONAL

3. D. Wang, N. Wang, K. H. Chou, J. A. Goldstein, C. L. Hsieh, C. A. J. Hoe, H. H. Chou, and Z. H. Miao, "Crystal structure of the 2019 nCoV spike in the prefusion conformation," *Nature* 571, 1200–1203 (2020). <https://doi.org/10.1038/s41586-020-1200-7>

1. Множественное
выравнивание содержит
больше информации, чем
парное

```

      *          20          *          40          *          60          *          80
THIE_LACLS : ~~~~MTNKTLDLSVYFIAGAQN---SECSLDGATQKIALIIKSGVTVYQFRDQGTIYK-EQKQ-RLSIAQKLQKVSEE : 70
THIE_MANSM : ~~~~MNKIKSMLSVYFIAGSQDCRHLPGEPTEENLLTILQRALEAGITCFQFREKGEQSLACDLQLKRRRLALKCLQLCRQ : 75
THIE_STRA3 : ~~~~~~MKDTLKLIFVCGTVDC-----SRKNILTVEEALQAGITLFOFREKGFAL-QGKE-KIAMAKQLQILCKQ : 64
THIE_LISIN : ~~~~~~MRAELAVYFIAGTQDI-----VRGTLPSVLEEALKAGITCFQYREKGGAGALQTASE-RKEMALECQQLCAK : 65
THIE_ANOFW : ~~~~~~MMKQKLSLYFVMGSIDC-----TKDPLAVLDEAIIKGGITMFOFREKGGKAL-TGIE-KYRLAEKLLERC : 64
THIE_GEOTN : MARITSEEMKERLAVYFIMGSONS-----ERPAEDVLKEALDGGVTLFOFREKGSAL-EGEE-KEALARQLQRLC : 71
THIE_BACSU : MTRISREMMKELLSVYFIMGSNNT-----KADPVTVVQKALKGGATLYQFREKGGDAL-TGEA-RIKFAEKAQAACRE : 71
THIE_BACA2 : MTRISREMMKMLSVYFIMGSNNT-----SADPVSVEKAIEGGATLFOFREKGGSL-TGEE-RLFAKRVQDVCRQ : 71
THIE_OCEIH : ~~~~MKFDKHMRLKYFIMGSQNC-----HRDPREILKEAASAGITAFQYREKGGKNSL-TGTA-KVELAKDLKAICH : 66
THIE_STAAB : ~~~~~~MFNQSYLNVYFICGTSNV-----PSHRTIHEVLEAALKAGITLFOFREKGESAL-KGND-KLVLAKEQLHLCHQ : 67
THIE_STACT : ~~~~~~MFQSKDLNVYFICGTDI-----PEGRTIQEVLKEALEGGITLYQFREKGGNAK-TGQD-KVALAKELQALCKS : 67
      L YF6 G 1          6 a G T 5Q5R7KG 1 4 A c

```

```

      *          100          *          120          *          140          *          160
THIE_LACLS : AGVSFIVNDDVELARELNADGIHIGQTDSESVSKVREKVGQEMWLGSLVTKADELKTAAQ-SSGADYDYGIGPIYPTNSKND : 149
THIE_MANSM : FQVPFIVNDDVELALSIAQADGIHVGQKDTAVETILRNTRNKPIIGLSINTLAQALANKDRQDIDYFGVGPIFPTNSKADH : 155
THIE_STRA3 : YQVPFIIDDDIDLVELIDADGLHIGQNDLPVDEARRRLPKDI-IGLSVSTMAEYQKSG-LSVVDYIGIGPENPTQSKADA : 142
THIE_LISIN : YQVPFIINDDVALALEIGADGIHVGQNDSEIRQVIASCAGKMKIGLSVSVSEAEAEERLGSVDYIGVGPIFPTISKADA : 145
THIE_ANOFW : YNIPFIVNDDVDLALALQADGVHVGQDEVAERVRDRIGDKY-LGVSVHNLNEVKKAL-AACADYVGLGPIFPTVSKEDA : 142
THIE_GEOTN : YGVFPIVNDDELALAIADAGVHVGQDDDEDARRVREKIGDKI-LGVSAHNVEEARAAL-EAGADYIGVGPIYPTRSKDDA : 149
THIE_BACSU : AGVPFIVNDDVELALNLKADGIHIGQEDANAKEVRAAIGDMI-LGVSAHTMSEVKQAE-EDGADYVGLGPIYPTETKKDT : 149
THIE_BACA2 : AGIPFIINDDVELALRLLEADGVHIGQDDADAETRAAIGDMI-LGVSAHNVSEVKRAE-AAGADYVGMGPVYPTETKKDA : 149
THIE_OCEIH : FQIPFIINDDVLAKQLDADGIHIGQDDQFVEVVRKQFPNKI-IGLSISTNNELNQSP-LDLVDYIGVGPIFDTNTKEDA : 144
THIE_STAAB : YNVFPIVNDVSLAKEINADGIHVGQDDAKVKEIAQYFTDKI-IGLSISDLGEYAKSD-LTHVDYIGVGPIYPTPSKHDA : 145
THIE_STACT : YNVFPIVNDVALAEEIDADGIHVGQDDAEVDDFNNRFEGKI-IGLSIGNLEELNASD-LTYVDYIGVGPIFATPSKDDA : 145
      6pFI6lDD6 La 6 ADG6H6GQ D          6G6S 2          DY G6GP pT 3K Da

```

```

      *          180          *          200          *          220          *          240
THIE_LACLS : AKPIGIKDLR-LMLLENQLPIVGIGGITQDSLTELSAIGLDGLAVISLITAEANPKKVAQMIRQKITKNG~~~~~ : 218
THIE_MANSM : SPLVGMNFIRQIRQLGIDKPCVAIGGIKEESAAILRRLGADGVAVISAISSHSVNIANTVKTLAQK~~~~~ : 220
THIE_STRA3 : KPAVGNRTTKAVREINQDIPIVAIGGITSDFVHDIIESGADGIAVISAIISKANHIVDATRQLRYEVEKALVNRQKRSDVI : 222
THIE_LISIN : EPVSGTAILEEIRRAGIKLPIVGIGGINETNSAEVLTAGADGVSVISAITRSEDCQSVIKQLKNPGSPS~~~~~ : 214
THIE_ANOFW : KQACGLTMIEHIRAHEKRVPLVAIGGITETQAKQVIEAGADGIAVISAIICRAEHIYEQTKRLYEMVMRAKQKQKDR~~~~~ : 217
THIE_GEOTN : NEAQGPGLRLHRLREQGITIPIVAIGGITADNTRAVIEAGADGVSVISAIASAPEPKAAAAALATAVREANL---R~~~~~ : 221
THIE_BACSU : RAVQGVSLIEAVRRQGISIPIVGIGGITIDNAAPVIEAGADGVSMISAIISQAEDPESAAARKFREEIQTYKTG--R~~~~~ : 222
THIE_BACA2 : EAVQGVTLIEEVRRGITIPIVGIGGITADNAAPVIEAGADGVSMISAIISQAEDPKAAARKFSEIIRRSKAGLSR~~~~~ : 224
THIE_OCEIH : KTAVGLEWIIQSLKKQHPSLPLVAIGGINTTNAQEIIQAGADGVSEISAITETHILQAVQRL~~~~~ : 206
THIE_STAAB : HTPVGPPEMIATFKEMNPQLPIVAIGGINTSNVAPIVEAGANGISVISAIISKSENIEKTVNRKDFEFNN~~~~~ : 213
THIE_STACT : SEPVGPKMIETLRKEVGDLPIVAIGGISLDNVQEVAKTSADGVSVISAIARSPHVTETVHKFLQYFK~~~~~ : 212

```

```

THE_LACLS 1 MTNKTLDLSVYFIAGAQN---SECSLDGATQKIALI...IKSGVTVYQFRDQGTIYKEQKQRLSIAQKLQKVSEEAG 72
THE_MANSM 1 MNKIKSMLSVYFIAGSQDCRHLPGEPTEENLLTILQRALEAGITCFQFREKGEQSLACDLQLKRRRLALKCLQLCRQFQ 77

THE_LACLS 73 VSFIVNDDVELARELNADGIHIGQTDSESVSKVREKVGQEMWLGSLV-TKADELKTAQSSGADYDYGIGPIYPTNSKND 148
THE_MANSM 78 VPFIVNDDVELALSIAQADGIHVGQKDTAVETILRNTRNKPIIGLSINTLAQALANKDRQDIDYFGVGPIFPTNSKAD 154

THE_LACLS 149 AAKPIG...IKDLRLMLLENQLPIVGIGGITQDSLTELSAIGLDGLAVISLITAEANPKKVAQMIRQKITKNG 218
THE_MANSM 155 HSPLVGMNFIRQIRQLGIDK...PCVAIGGIKEESAAILRRLGADGVAVISAISSHSVNIANTVKTLAQK... 220

```

THIE_LACLS : ~~~~MTNKTLDLSVYFIAGAQNF---SECSLDGATQKIALIIKSGVTVYQFRDKGTIYK-EQKQ-RLSIAQKLQKVSEE : 70
THIE_MANSM : ~~~~MNKIKSMLSVYFIAGSQDCRHLPGPEPTENLLTILQRALEAGITCFQFREKGEQSLACDLQLKRRILAKCLQLCRQ : 75
THIE_STRA3 : ~~~~MKDTLKLYEVCQTVDC-----SRKNILTVVEALQAGITLQFREKGEPTAL-QGKE-KIAMAKQLQILCKQ : 64
THIE_LISIN : ~~~~MRAELAVYFIAGTQDI-----VRGTLPSVLEEALKAGITCFQYREKAGALQTASE-RKEMALECQQLCAK : 65
THIE_ANOFW : ~~~~MMKQKLSLYFVMGSIDC-----TKDPLAVLDEAIKGGITMFOFREKKGKAL-TGIE-KYRLAEKLLERCRM : 64
THIE_GEOTN : MARITSEEMKERLAVYFIMGSONS-----ERP AEDVLKEALDGGVTLFQFREKKGSAAL-EGEE-KEALARQLQRLCRT : 71
THIE_BACSU : MTRISREMMKELLSVYFIMGSNNT-----KADPVTVVQKALKGGATLYQFREKGGDAL-TGEA-RIKFAEKAQAACRE : 71
THIE_BACA2 : MTRISREMMKMLSVYFIMGSNNT-----SADPVSVEKAIEGGATLYQFREKKGSGSL-TGEE-RLLEAKRVQDVCRQ : 71
THIE_OCEIH : ~~~~MKFDKHMRLRYFIMGSONC-----HRDPRBILKEAASAGITAFQYREKKGKNSL-TGTA-KVELAKDLKAICH : 66
THIE_STAAB : ~~~~MFNQSYLNLYFICGTSNV-----PSHRTIHEVLEAALKAGITLQFREKKGESAL-KGND-KLVLAKEQLHLCHQ : 67
THIE_STACT : ~~~~MFQSKDLNLYFICGTQDI-----PEGRTIQEVLKEALEGGITLYQFREKKGNGAK-TGQD-KVALAKELQALCKS : 67

THIE_LACLS : AGVSFIVNDDVELARELNADGIHIGQTDSESVSKVREKVGQEMWLGLSVTKADELKTAQ-SSGADYLGIGPIYPTNSKND : 149
THIE_MANSM : FQVFFIVNDDVELALSIQADGIHVGQKDTAVETILRNTRNKPIIGLSINTLAQALANKDRQDIDYFGVGPIFPTNSKADH : 155
THIE_STRA3 : YQVFFIIDDDIDLVELIDADGLHIGQNDLPVDEARRRLPDKI-IGLSVSTMAEYQKSQ-LSVVDYIGIGBFNPTQSKADA : 142
THIE_LISIN : YQVFFIINDDDVALALEIGADGIHVGQNDDEEIRQVIASCAGKMKIGLSVHVSVEAEBAERLGSVDYIGVGPIFPTISKADA : 145
THIE_ANOFW : YNIFVIVNDDVDLALALQADGVHVGQDEDAEVRDRIGDKY-LGVSVHNLNEVKKAL-AACADYVGLGPIFPTVSKEDA : 142
THIE_GEOTN : YGVFFIVNDDVELAIAIDAGVHVGQDDEARRVREKIGDKI-LGVSAHNVEEAARAAI-EAGADYIGVGPIYPTRSKDDA : 149
THIE_BACSU : AGVFFIVNDDVELALNLKADGIHIGQEDANAKEVRAAIGDMI-LGVSAHTMSEVKQAE-EDGADYVGLGPIYPTETKKDT : 149
THIE_BACA2 : AGIFFIINDDDVELALRLBADGVHIGQDDADAEEETRAAIGDMI-LGVSAHNVSEVKRAE-AAGADYVGMGPVYPTETKKDA : 149
THIE_OCEIH : FQIFFIINDDDVDLAKQLDADGIHIGQDDQPVVVRKQF PNKI-IGLSISTNNEINQSP-LDLVDYIGVGPIFDNTSKEDA : 144
THIE_STAAB : YNVFFIVNDDVSLAKEINADGIHVGQDDAKVKEIAQYFTDKI-IGLSISDLGEYAKSD-LTHVDYIGVGPIYPTPSKHDA : 145
THIE_STACT : YNVFFIVNDDVALAEEDADGIHVGQDDEAVDDFNNRFEKGI-IGLSIGNLEELNASD-LTYVDYIGVGPIFATPSKDDA : 145

THIE_LACLS : AKPIGKIDLR-LMLLENQLPIVIGGGITQDSLTELSAIGLDGLAVISLLTEAENPKKVAQMIRQKITKNG~~~~~ : 218
THIE_MANSM : SPLVGMNFIRQIRQLGIDKPCVAIGGIKEESAAILRRLGADGVAVISAISHSVNIANTVKTLAQK~~~~~ : 220
THIE_STRA3 : KPAVNRTTRAVREINQDIPVVAIGGITSDFVHDIESGADGIAVISASIRAHIVDARQLRIEVERALVNKRQSRSDVI : 222
THIE_LISIN : EPVSGTAILEEIRRAGIKLPVIGGGINETNSAEVLTAGADGVSVISAITSRSDQSVIKQLKNPGSPS~~~~~ : 214
THIE_ANOFW : KQACGLTMIEHIRAHEKRVPLVAIGGITEQTAKQVIEAGADGIAVISAICRAEHIYEQTKRRLYEMVMRAKQKGDR~~~~~ : 217
THIE_GEOTN : NEAQGPGLIRHLREQGITIPIVVAIGGITADNTRAVIEAGADGVSVISAISAPKAAAAALATAVREANL---R~~~~~ : 221
THIE_BACSU : RAVQGVSLIEAVRRQGISIPVIGGGITIDNAAPVIAQADGVSMISASISQAEDPESAARKFREEIQTYKTG---R~~~~~ : 222
THIE_BACA2 : EAVQGVTLIEEVRRQGITIPIVIGGGITADNAAPVIEAGADGVSMISASISQAEDPKAAARKFSEEIRRSKAGLSR~~~~~ : 224
THIE_OCEIH : KTAGLEWLIQSLKKQHPSLPVVAIGGINTTNAQEIIEAGADGVSEFISAITETHILQAVQRL~~~~~ : 206
THIE_STAAB : HTPVGPEMIATFKMNPQLPIVVAIGGINTSNVAPIVEAGANGISVISAISKSENIEKTVNRFKDFNN~~~~~ : 213
THIE_STACT : SEPVGPKMIETLRKEVGDLPIVVAIGGISLDNVQEVAKTSADGVSVISAIAIARSPHVTETVHKFLQYFK~~~~~ : 212

THIE_LACLS : MTNKTLDLSVYFIAGAQNFSECSLDGATQKIALI---IKSGVTVYQFRDKG---TIYKEQKQRLSIAQKLQKVSEEAG :
THIE_MANSM : MNKIKSMLSVYFIAGSQDCRHLPGPEPTENLLTILQRALEAGITCFQFREKGEQSLACDLQLKRRILAKCLQLCRQFQ :
M LSVYFIAG Q1 6 66 6 G6T 5QFR KG 36 4 6A K 6 2

THIE_LACLS : VSFIVNDDVELARELNADGIHIGQTDSESVSKVREKVGQEMWLGLSVTKADELKTAQSSGADYLGIGPIYPTNSKND :
THIE_MANSM : VFIVNDDVELALSIQADGIHVGQKDTAVETILRNTRNKPIIGLSINTLAQALANKDRQDIDYFGVGPIFPTNSKAD :
V FIVNDDVELA 6 ADGIRGQ D V 6 6GLS6 T A L D1 6G6T15PTNSK D

THIE_LACLS : AAKPIG---IKDLRLMLLENQLPIVIGGGITQDSLTELSAIGLDGLAVISLLTEAENPKKVAQMIRQKITKNG : 218
THIE_MANSM : HSPVGMNFIRQIRQLGIDKPCVAIGGIKEESAAILRRLGADGVAVISAISHSVNIANTVKTLAQK----- : 220
6G 14 6R 6 6 P V IGI 2 S L 6G DG6AVIS 63 N 6 QK

2. Консервативное значит важное

Формулировка МГ: принцип Вальда(Wald)

Примеры

- Глиняная посуда > 20 тыс. лет
 - Пробовали деревянную, pewter – *сплав олова, включающий свинец*, плетённую и др.
- Правые и левые туфли (Древняя Греция ?)
 - До того и в средние века обе туфли были одинаковые

РНК зависимая РНК полимераза (RdRP), консервативные участки

```

*      320      *      340      *      360      *      380      *      400      *      420      *      440      *      460
FKTMIRFGDVLDDFFSADASLSPPMIREA...GRIMSELS...GTPSHFGTALINTIIYSKHLIYNCCY...HVCGSMPSGSPCTALINSTINNLYYVFSKIFGKSPVFF...CQALKILC.YGDDVLIIVFSRDV
EVAMQG.FERVYVDVYSNEDSTHSVAMFRLL..A...EEFF.TPENGFDPLTREYLESLAISTHAFEERKF...LITGGLPSGCAATSMINTIMNNIIRAGLYLTYNFEFDD...VKVLS.YGDDLVATNYQL
STHFAQ.YKNVWDVLYSADANHCSDAMNIMFEEVFRTFG...FHPNAEWILKTLVNTTEHAYENKRI...VVEGCMPSGCSATSIINTILNNIYVLYALRRHYEGVELDT...YTMIS.YGDDIVASDYDL
...WSLCVATIVSDHDTFWPGWIRDLICDELINMGYA.PWWVKLFETSLKLPVYVGAFAPEQGHTLLGDPSNPDLVGLSSGQGATDLMGTILMSTIYLVMLQDHTAPHLNSRIKEMPSACRFLDSYWQGHEETROIS.KSDDAILGWTKGR
LRLRPE.NWVYCCADGSCEDSSLTPLYINAV..LTIRSTYMEDWDVGLQMLRNLYTEIVYTPISTPDGTIV...KKFRGNNSGQPSIVDNLVVIAMHYALIKECFEVEEID...STCVFFV.NGDDLIIVAVNPEK
HDKLNRPGWLHGSGDGRDSSIDPFFFDVV..KTKRKHEL..PSEHHKAIDLIYDEILNTTICLANGMVI...KKNVGTQR.QPSTVVDNTIVIMTAFLYAYIHKTGDREIAL...LNERFIFVC.NGDDNKFAISPQF
AISIASFSFPGFNGDFANEDGMFHSSFSMV..SEIANIFY...GNFLSTERDNLTRMLTNRFSLMKGAIL...RVPGGPGSGFFMTVINSFINLFYLSAWIMLARFNGRQDISH...PCNFPKYVRACV.YGDDNIVAIMEV
AARMKEKGNDVLCODYSSEDGLLSKQVMDVI..ASVINELC.GGEDQLKNARRNLLIMACCSRIAICKNTVW...RVECGIPSGGFMTVINSFINELIRYHYKKIMREQQAPELMV...QSFDKLIGLVT.YGDDNLSVNAV
YAEHAK.YKNHFADYIANDSTQNQIMTES..FSIMSRLT...ASPELAEVVAQDLIAPSEMDVGDYVI...RVKEGLPSGFPCTSQVNSINHWITLCALEATGLSPDVV...QMSYFSFYGDDEIVSTDIDF
NNLTSKASDFLCIDYSKNDSTMSPCVVRIA..IDLADCC...EQTELTKSVVLTLSKSHFMTILAMIV...QTKRGLESGMPFTSVINSICHWLWSAAVYKSCAEIGLHCS...NLYELAPFYT.YGDDGVYAMTFMM
IQRIKS.AAKVYAVDYSKNDSTQSPRVSAA..IDLRYFS...DRSPIVDSAANTLKSPPIAFNGVAV...KVSSGLPSGMPITSVINSINHCILYVGCAILQSLEARGVPVTW...NLFSTFLMMT.YGDDGVYMFPMFM
TKRLERPKHDRYCVLYSKNDSTQPPKVTQS..IDLRHFT...DKSPIVDSACATLKSNPIGIFNGVAF...KVAGGLPSGMPITSIINSINHCILVGSAAVVKAELEDSGVRVTW...NIFDSMDLFT.YGDDGVYIVPPLI
D      D                                     g  sg  T  n3                                     gDD

```

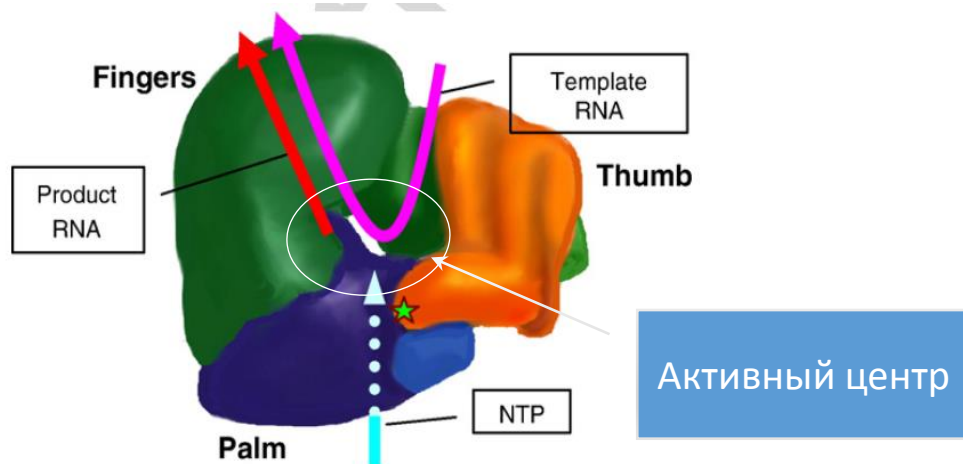


Fig. 1. Schematic architecture of “small” RdRP. The hairpin between the palm and thumb domains is in light blue. The predicted approximate location of MV RdRP Trp460 is marked by star.

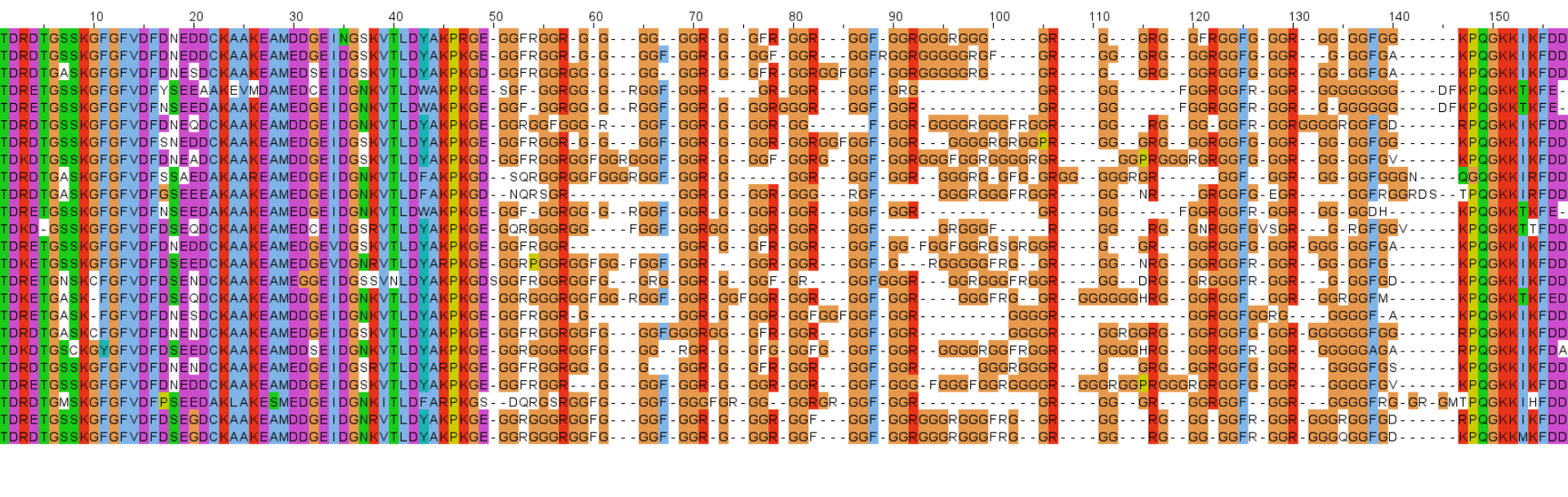
Рисунок мой :) хвастаюсь

Каталитические мотивы RdRP

	Motif A					Motif B					Motif C			
	★					★					★			
Nido-/Arteri-	442	LETDL	ESCD	RSTP	454 [44]	499	GLS	SGDP	ITSISNTI	513 [40]	554	RVYIYS	DDVVLT	565 [128]
Nido-/Corona-	615	MGWD	YPKC	DRAMP	627 [51]	679	GTS	SGDAT	TAYANSV	693 [60]	754	SMMILS	DDAVVC	765 [167]
Nido-/Abyss-	-	CGGD	YEKY	DKNLA	- [54]	-	GNT	SGNSR	TKTVNGN	- [64]	-	RMVCV	GDDYIKV	- [-]
Nido-/Euroni-	-	VSLD	HSKF	DRFVA	- [52]	-	GIS	SGNSI	TALNNSL	- [50]	-	RIAGLS	DDVVAC	- [-]
Nido-/Medioni-	-	SGKD	FPQW	DRSVE	- [56]	-	GVC	SGNSK	TAPGNSI	- [60]	-	LLRVL	SDGMVL	- [-]
Nido-/Mesoni-	-	GGKD	YPKW	DRRIS	- [58]	-	GVT	SGNSR	TADGNSL	- [66]	-	KGAYLS	DDGLIV	- [-]
Nido-/Mononi-	-	FSFD	YTAF	DRTTT	- [53]	-	SVSS	SGNAH	TAPWNSH	- [76]	-	SIQII	GDDLITN	- [-]
Nido-/Roni-	-	ISQD	YPKF	DTDVD	- [50]	-	GVS	SGDGATA	IKNSH	- [56]	-	RCATLS	DDTLAI	- [-]
Nido-/Tobani-	-	MGAD	YTKC	DRSFP	- [47]	-	GTT	SGDST	TAFSNSF	- [57]	-	FLHFLS	DDSFII	- [-]
Picornia-/Dicistro-	286	IAGD	FSTF	DGSLN	298 [48]	347	SQP	SGNPAT	TPLNCF	361 [30]	392	SMVSY	GDDNVIN	403 [143]
Picornia-/Ifla-	252	LQMD	YKNYS	DAIP	264 [52]	317	GVL	AGHPM	TSVVNSV	331 [25]	357	YIIVM	GDDVVIS	368 [271]
Picornia-/Picornia-	230	FAFD	YTYG	DASLS	242 [42]	285	GMP	SGCSG	TSIFNSM	299 [22]	322	KMIAY	GDDVIAS	333 [128]
Picornia-/Seco-	277	LCCD	YSSF	DGLLS	289 [48]	338	GIP	SGFPM	TVIVNSI	352 [31]	384	GLVTY	GDDNLIS	395 [316]
Picornia-/Polycipi-	-	VDFD	VSNW	DGFLF	- [50]	-	GII	SGFP	TAEVNTL	- [29]	-	SAILY	GDDILLT	- [145]
Picornia-/Marna-	-	IAGD	YSSF	DMSHN	- [52]	-	WVM	SGVPL	TAELSST	- [25]	-	ALIVY	GDDNNA	- [-]
Tymo-/Tymo-	316	IAND	YTAF	DQSQH	328 [40]	369	MRLT	GEPG	TYDDNTD	383 [15]	399	PIMVS	GDDSLID	410 [185]
Tymo-/Alphaflexi-	-	LAND	YTAF	DQSQD	- [40]	-	MRLT	GEGP	TFDANTE	- [16]	-	AQVYA	GDDSALD	- [122]
Tymo-/Betaflexi-	-	TDSD	YEAF	DRSQD	- [40]	-	MRFS	GGEFG	TFFFNTI	- [16]	-	PICFA	GDDMYSP	- [140]
Tymo-/Deltaflexi-	-	TGND	YTAW	DSGID	- [40]	-	RQE	SGDRW	TWILLNTL	- [16]	-	PLCVS	GDDSVTL	- [135]
Tymo-/Gammaflexi-	-	TDGD	YTAY	DASQD	- [40]	-	MRFS	GGEVW	TYLFNTL	- [15]	-	AQVYG	GDDKSIN	- [131]
-/Alphetetra-	1076	KSID	IKFE	DTVHN	1088 [43]	1132	MLD	SGAVW	TIARNTL	1146 [14]	1161	FIAAK	GDDVFLA	1172 [753]
-/Astro-	266	IEFD	WTTRY	DGTIP	278 [51]	330	GNP	SGQFS	TPMDNM	344 [24]	369	DTVVY	GDDRLST	380 [138]
-/Barna-	305	CETD	YISG	WDSVQ	317 [53]	371	GQL	SGDYN	TSSSNSR	385 [22]	408	GIKAM	GDDSFEI	419 [104]
-/Beny-	297	GVID	DAACD	DSGQG	309 [44]	354	VKT	SGEPG	TLLGNTI	368 [16]	385	CMAMK	GDDGFKR	396 [172]
-/Botourmia-	404	VSGD	YSAAT	DNLH	416 [58]	475	GQLM	SGPLS	FPVLCI	489 [23]	513	GILVN	GDDILFR	524 [336]
-/Bromo-	462	LEAD	LSKF	DKSQG	474 [46]	521	QRTT	GDAFT	TYFGNTL	535 [16]	552	CAIFS	GDDSLII	563 [259]
-/Calici-	239	YDAD	YSRW	DSTQQ	251 [45]	297	GLP	SGVPC	TSQWNSI	311 [25]	337	LFSFY	GDDIIVS	348 [162]
-/Carmotetra-	582	ISFD	LSRW	DMHVQ	594 [44]	639	GIM	SGDMT	TGLGNCI	653 [63]	717	SILDD	GDDHVII	728 [187]
-/Clostero-	277	LEID	FSKF	DKSQG	289 [46]	336	QRTT	GSPTN	TWLSNTL	350 [16]	367	LLVSV	GDDSLIF	378 [137]
-/Flavi-/Flavi-	532	YADD	TAGW	DTTRIT	544 [55]	600	QRG	SGQVP	TYALNTI	614 [46]	661	RMVVS	GDDCVVR	672 [233]
-/Hepaci-	217	FSYD	TRCF	DSTVT	229 [49]	279	CRA	SGVLT	TS CGNTL	293 [18]	312	TMLVC	GDDLVI	323 [268]
-/Pegi-	211	ICVD	ATCF	DSSIT	223 [45]	269	CRS	SGVLT	TSASNCL	283 [18]	302	SLLIA	GDDCLII	313 [250]
-/Pesti-	342	VSFD	TKAW	DTQVT	354 [47]	402	QRG	SGQPD	TSAGNSM	416 [25]	442	RIHVC	GDDGFLI	453 [266]
-/Hepe-	256	YEND	FSAF	DSTQN	268 [44]	313	KKH	SGEPG	TMLFNTI	327 [16]	344	LALFK	GDDSLVC	355 [129]
-/Kita-	901	YEFD	MSKY	DKSQG	913 [46]	960	QRK	SGDAS	TYFGNTV	974 [16]	991	FGAFS	GDDSLIF	1002 [142]
-/Levi-	271	ATVD	LSAAS	DSIS	283 [36]	320	ISSM	GNGYT	TFELES	334 [18]	353	EVTVY	GDDIILP	364 [225]
-/Luteo-	229	IGVD	ASRF	DQHVS	241 [47]	289	HRM	SGDIN	TSMGNKL	303 [17]	321	ELCNN	GDDCVII	332 [200]
-/Narna-	355	ISSD	DKSAS	DLIP	367 [46]	414	GILM	GLPT	TWAILNL	428 [26]	455	DCRVC	GDDLIGV	466 [363]
-/Noda-	591	SEG	DFSN	DFTVS	603 [49]	653	GVK	SGSPT	TCDLNTV	667 [25]	693	IGLAF	GDDSLFE	704 [339]
-/Permutotetra-	366	ICPD	FKQM	DGSVD	378 [56]	435	GLMT	GVVGT	TTFDNTV	449 [103]	345	RIACY	GDDTDIY	356 [901]
-/Poty-	245	CDAD	GSQF	DSSLT	257 [49]	307	GNN	SGQPS	TVVDNSL	321 [24]	346	VFFVN	GDDLII	357 [162]
-/Solemo-	286	AEAD	ISGF	DWSVQ	298 [51]	350	IMK	SGSYC	TSSTNSR	364 [12]	377	WCIAM	GDDSVGE	388 [165]
-/Solinvi-	485	FSCD	YKNF	DRTIP	497 [45]	543	GMP	SGCVPT	TAPLNSK	557 [32]	590	CRLFY	GDDVIIA	601 [195]
-/Toga-	373	LETD	IASF	DKSQD	385 [46]	432	MMK	SGMFL	TFLVNTV	446 [18]	465	CAAFI	GDDNIIH	476 [140]
-/Tombus-	527	IGLD	ASRF	DQHCS	539 [48]	588	CRM	SGDIN	TS LGNYL	602 [18]	621	SLANC	GDDCVLI	632 [186]
-/Virga-	230	IEID	ISKY	DKSKT	242 [46]	289	QOK	SGNVD	TYFSNTW	303 [16]	320	FSIFG	GDDLIL	331 [146]
-/Alverna-	-	CSSD	ASGW	DMSVS	- [57]	-	ITA	SGLPD	TTTQNSF	- [12]	-	KALTA	GDDLIC	- [112]
-/Matona-	-	IEVD	FTFE	DMNQT	- [44]	-	ERT	SGEPAT	LLHNTT	- [16]	-	AGIFQ	GDDMVIF	- [144]
-/Hypo-	-	TAMD	VTAM	DSTAS	- [53]	-	GLST	GHATT	TPSNT	- [25]	-	KFSFS	DDNFWS	- [-]

The RdRP active site is surrounded by the palm, fingers, and thumb domains with seven catalytic motifs (motifs A–G) distributed within the palm (motifs A–E) and fingers (motifs F–G) (Poch et al., 1989; Gorbalenya et al., 2002; Bruenn, 2003; te Velthuis, 2014; Wu et al., 2015) (see an alignment of motif A–C of the 49 representative RdRP sequences in Figure 2).

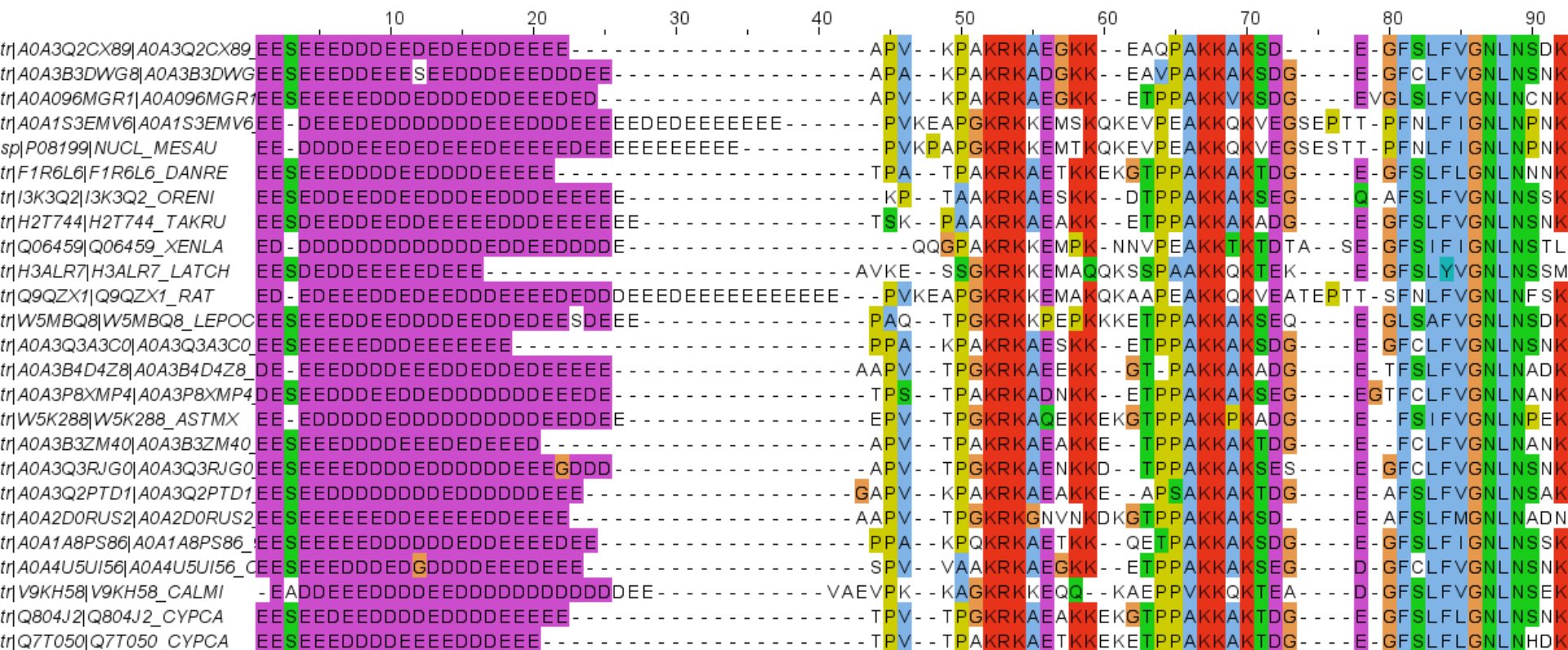
Nucleolin. Глицин-аргинин богатый участок (GAR)



КОНСЕРВАТИВЕН - присутствует в этом месте у всех нуклеолинов.
Значит важен для функции. Зачем он – толком не известно.
GAR встречаются у многих ДНК и РНК связывающих белков.

Я отредактировал выравнивание вручную чтобы GAR поместился на слайде (в выравнивании, построенном программой, был в два раза длиннее по числу позиций)

Nucleolin. Кислый участок и сигнал сигнал ядерной локализации



3. Сходные последовательности ГОМОЛОГИЧНЫ

В выравнивании гомологичных последовательностей
гомологичные остатки стоят в одной колонке

Был ли GAR у общего предка этих белков?



У

остатка у общего предка нельзя.

Из-за повторяющихся мотивов RGG и других, можно предположить, что GAR эволюционировали путем дупликации повторов.

Уникальная дупликация простых повторов в геноме индивидуума лежит в основе идентификации личности по ДНК

4. Три типа эволюции белков

1. Локальные мутации в белке (из-за мутаций ДНК в кодирующей последовательности) => ДОМЕНЫ

- Этот вид эволюции описывается выравниванием последовательностей
- Выравнивание последовательностей гомологичных белков ныне живущих организмов позволяет реконструировать последовательности у предков! (с некоторой точностью)

2. «Выращивание» коротких последовательностей для определенной функции

- ASN GLYCOSYLATION, [PS00001](#); N-glycosylation site (PATTERN with a high probability of occurrence!)
N-{P}-[ST]-{P} N is the glycosylation site

3. Крупные перестройки => DOMAIN SHUFFLING

Домены белков

[Длинные] гомологичные участки из разных белков, которые эволюционируют только по типу локальных мутаций, **и максимальной длины, с сохранением этого свойства**, называются

ЭВОЛЮЦИОННЫМИ ДОМЕНАМИ

Терминологическая проблема.

ДОМЕН – набор фрагментов последовательностей и их выравнивание. Имеет название. Например, RdRP_1

ДОМЕН белка – фрагмент последовательности, входящий в определенный домен, например, в RdRP_1

ДОМЕННАЯ АРХИТЕКТУРА – последовательность доменов в белке

ДВА ДОМЕНА гомеобелков: гомеодомен и OAR домен

[illegible]

Домены принято изображать так

[X1WJ92 ACYPI](#) [Acyrtosiphon pisum (Pea aphid)]
Uncharacterized protein (408 residues)



There are 1836 sequences with the following architecture:
Homeodomain, OAR

Гомеодомен является ДОМЕНОМ

Доказательство

- Выравнивание, свидетельствующее о гомологичности последовательностей
- Представленность домена в негомологических белках (обычно, но не обязательно)

```
G4VR66_SCHMA/203-259
HME2B_DANRE/173-229
HM1N_BOWWO/373-429
T1EHES_HELRO/41-97
HM05_CAEL/36-92
G3IBX4_CRIGR/25-81
F1QFR3_DANRE/136-192
B8A5N9_DANRE/135-191
Q91967_CHICK/77-133
DLX3B_DANRE/126-182
HM43_CAEL/103-159
HM23_CAEL/212-268
MSX3_MOUSE/88-144
HM30_CAEL/96-152
BARH2_RAT/230-286
BARH1_DROME/300-356
BARX2_MOUSE/138-194
BSH_DROME/275-327
H2XU6_CIOIN/470-524
HM19_CAEL/95-151
SLOU_DROME/546-602
F6VWQ6_XENTR/112-168
TIN_DROME/302-358
NKX25_RAT/138-194
H0XK12_OTOGA/100-156
HM09_CAEL/71-127
H2VEX2_TAKRU/13-69
U3K517_FICAL/59-115
TLX3_CHICK/173-229
U37ZQ6_FICAL/136-188
LBX1_MOUSE/126-182
G4VGG4_SCHMA/38-94
BCD_DROME/98-153
BCD_DROME/98-153 (SS)
VENTX_HUMAN/92-148
VENT1_XENTR/128-184
Q804C9_XENTR/190-246
K48BZ1_SOLLC/24-79
PHO2_YEAST/78-134
WOX9_ARATH/52-113
WOX9_ORYSJ/11-72
WOX2_ORYSJ/24-85
WOX4_ARATH/87-148
WOX1_ARATH/73-134
WOX2_ARATH/11-72
WOX5_ORYSJ/41-102
WUS_SOLLC/25-85
WOX6_ARATH/58-119
YHP1_YEAST/174-230
YOX1_YEAST/177-233
HARA_DICDI/163-219
PHX1_SCHPO/169-223
CUT_DROME/1746-1802
CUX2_MOUSE/1114-1170
CUX1_MOUSE/1240-1296
Q22810_CAEL/212-268
HBX2_DICDI/486-542
```

There are 25976 sequences with the following architect
X2JL88_DROME [Drosophila melanogaster (Fruit fly)] Uncharacterized

Show all sequences with this architecture.

There are 2311 sequences with the following architectu
X2JDY7_DROME [Drosophila melanogaster (Fruit fly)] POU domain pr

Show all sequences with this architecture.

There are 2108 sequences with the following architectu
W6NCH4_HAECO [Haemonchus contortus (Barber pole worm)] Zinc f

Show all sequences with this architecture.

There are 1903 sequences with the following architectu
MOU1E3_MUSAM [Musa acuminata subsp. malaccensis (Wild banana)]

Show all sequences with this architecture.

There are 1836 sequences with the following architectu
X1WJ92_ACYPI [Acyrtosiphon pisum (Pea aphid)] Uncharacterized p

Эволюционные домены

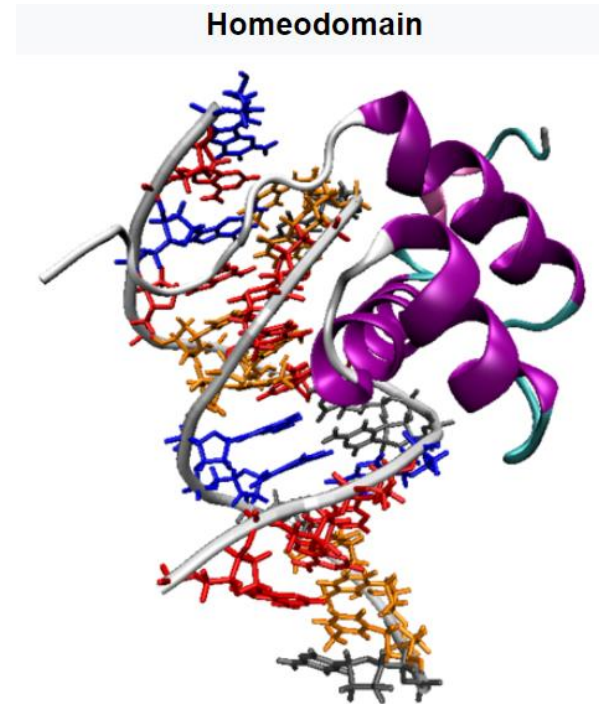
- Имеют определенную функцию (не всегда известна)

DUF – Domain of Unknown Function

- Часто совпадают со структурными доменами (но не всегда)

Гомеодомен – ДНК связывающий домен

Homeodomain proteins regulate gene expression and cell differentiation during early embryonic development, thus mutations in homeobox genes can cause developmental disorders.^[1]



5. БД Pfam (<http://pfam.xfam.org/>)

Protein Families (Pfam)

База семейств доменов

Белков

Одна запись содержит
информацию о ДОМЕНЕ,
информацию о всех белках,
содержащих данный домен,
выравнивание доменов из
всех белков



Rubens: **Holy Family** with St. Elizabeth

Что посмотрим в Pfam

- Название домена
- ID
- AC
- Выравнивания
 - Seed
 - Full
- Доменные архитектуры
- Описание функции
- Таксономическая распространенность
- 3D структуры
- HMM профиль выравнивания
- Clan
- Type I restriction modification DNA specificity domain
- Methylase_S
- **PF01420**
- Seed – по нему сделан профиль домена
- Full - все находки по профилю в Uniprot (с опозданием по версиям)
- This domain is also known as the target recognition domain (TRD).

Смотрим домен
methyltransferaseD12

6. Ревизия выравнивания и блоки

Консервативный блок выравнивания (К-блок)

- Подтверждается консервативными или функционально консервативными позиция во всех фрагментах. Назовем «абсолютно консервативными (АбКпоз)
- 1я и последняя позиции такие - АбКпоз
- В блоке нет гэпов
- АбКпоз встречаются достаточно часто. Пока на это нет строгих критериев.
- Идущие подряд К-блоки, разделенные короткими неконсервативными линкерами, возможно, с небольшим числом позиций с гэпами, стоит объединить в суперконсервативный блок (СК-блок) и предполагать консервативность ВСЕХ фрагментов СК-блока

Консервативный вертикальный блок

		*	120	*	140	*
A0A1G9TZ02	:	----	EVP---	D--HDILVGGWPCPSFSI--	MGD-----	KEG--MDDERG----
A0A1H3E3S2	:	--esh	PIQR-dE--	IDVVIGGPPCKGFSSI--	AG-----	HRD--PDDERN----
A0A1I6G129	:	----	tEWS---	D--ADVVGPPCQGFSNln	STK-----	TDE--LDDDRN----
A0A1I6HFR3	:	-lige	YLDD--	DadATLIAACAPCQPFSP--	LNH-----	GKE--SSDHAM----
A0A1I6HG27	:	---	seLYPD-gA--	TKVLAGCAPCQPFSNln	NGT-----	DTS--VRDDYG----
A0A1M5MP51	:	----	AVVGgdD--	VDLLVAGPECTHFS--	ARG-----	GKP--VSEQRR----
A0A1M5USI2	:	----	AVVGddD--	VDLLVAGPECTHFSR--	ARG-----	TKP--VSDQRR----
A0A256ILS2	:	retgh	GVE---	D--VDVVIGGPPCQGFSR--	LNNeriel	DEM--EKDRRN----
A0A256KSV9	:	----	SVP---	S--HDLLIACWPCPSFSR--	MGK-----	LDG--LEDERG----
D4H0C8/5-4	:	irnef	GLEP-gE--	VDVIAGCPPCQNFESK--	LRD-----	TTPwpEDEPKD----
G2MPP5/169	:	--iqd	AVS---	E--LDLLVGGPPCQSLSK--	AGY-----	RSR--RGDDDEDysi
J3ETN5/3-3	:	vadlf	DSSA--	E--ATVLAGCAPCQPFSP--	LTH-----	GED--SSEHES----
J3JDK0/3-3	:	--iaq	MYPW--	DadLKVLAACAPCQPYST--	MGH-----	SKG--NTHEDHn--
M0FSM4/5-3	:	----	DLGK-aD--	VDLVIGGPPCQPFSA--	AARra-gg	IEG--TESDEG----
V6DPC1/13-	:	----	KLPD--	D--LDLLAGGPPCKGFSS--	AQG-----	ETN--TDDPRN----
ConSeq/1-2	:	----	tlS-----	lDllluPPCpsFSp--	hst-----	pcs---cDc+s----
			66	pC	S	

В вертикальном блоке, по определению, содержатся фрагменты ИЗ ВСЕХ последовательностей.

Значит в вертикальном К-блоке мы предполагаем гомологичность всех фрагментов и всех остатков в каждой из колонок, неважно, консервативные они или нет

Консервативный НЕ вертикальный блок

```

      440      *      460      *      480      *      500      *      520      *      540
A0A1G9TZ02 : -----SVWHENRSGN----- : 224
A0A1H3E3S2 : -----IADVDPGESLYESYGDSWR----- : 252
A0A1I6G129 : -----CWKG--YESGGTDLFGR----- : 263
A0A1I6HFR3 : -----HKQDSGSTFDSVYGRMEPDE----- : 269
A0A1I6HG27 : -----HRKESGRSFDSDVYGRMEW----- : 258
A0A1M5MP51 : ivpvddledaladrdepflvsstgtaavdggtmvmgqgsnaraldaDRepvptiatrgavhfi eagpfvkprnlprggllhtnatyvan : 358
A0A1M5USI2 : alderaepflvsttvstaadtgtrmimqggsnaraldaDSevpptvatrgavhfi eagpfikprn-----l : 342
A0A256ILS2 : -----cdCADTLACPHEPEIVKRYGT----- : 272
A0A256KSV9 : -----AGQVTTRP-YSGTLRASS----- : 254
D4H0C8/5-4 : -----HRG-----Fdtqa----- : 310
G2MPP5/169 : -----dtgwdvkynsdGEYSEYVEYDVGTAENKRFGDkyrmlew----- : 318
J3ETN5/3-3 : -----dchKKDTGATYQSVYGRMEPDK----- : 274
J3JDK0/3-3 : -----HRKASGRSYKAPYSRMRPD-E----- : 274
M0FSM4/5-3 : -----HPEPIFGWRSRFSYLYKA----- : 271
V6DPC1/13- : -----TRDGEDVWIPTNHKHQDHSRShrekmaeyelgksg----- : 286
ConSeq/1-2 : -----tp-p-h-----t----- : 187
```

Как бы назвать такие блоки? Пока «не вертикальный К-блок»

В таком блоке все фрагменты гомологичны и все остатки фрагментов из блока в каждой позиции блока предполагаем гомологичными.

Минус блок (не консервативный)

250	260	270	280	290	300	310	320	330	340	350
VLNALDFGLPQKRERTIIIVGFKENY-----KFRIPERNG-----ESRDLDDVLLDD-----										
TLNAADYGVPOKRRRVIFQGRDGS-----PTYPERTHGPSK-----GATLTGRQLKPY-----										
RLWAHKYGVPORRHRAFIIGSRLGT-----PVFPATTSSEVR-----TVRDAIADLPTKPNN-----										
RVYCPYEGIPQKRRRWVVGSGEGR-LDI-----gTPPI TDES DYPTVK-----BAIGHLPKIKAD-GE-----										
KVSCPDYGIPORRHRVLLASKLGLDI-----SLIDPTHDPDSY-----PTVREVI GDL PPLEAGE-----										
VLNAADYGDATSRRLFFVVARRGHRATH-----PEPTHARAPAPDDdrapwrpaaeiidwtldr gssiwtrsrpIsNNTMQRI AQCLRD-----										
VLNAADYGDATSRRLFFVVARRDGRATH-----PAPTHAEVPDEDR-----epWFSAAEVIDWSDRGSS-----										
VLEAERYGVPOKRRRLFFVGTNRDVP I R-----FPEPTTPAEPs-----wRTAGEALADVDD-----										
VLNALNFGFLPQHRRLIIVGFRNDLAPDe----dsFSIPTQNRAKLETE-----ADQREALANI LEDD-D-----										
VVNAADYGVPORRMTIGICIIYGASDSE-----vefp petharePEDg-----IKERWVTVKDVLKEEYErgrlkqdlidlg-----										
LQDMTELGVPOKRRVLLVGIRDDL TENatadnI ISEL SIEGKTRSIQQG-----LSGLPRIrrgggggrl-----										
RAFCPEFGLPQTRRRWL VVGSRNSLVKL-----TPTHEDPEQYPTVRE-----CIGGGKLDKIEAGE-----										
NVYCPYEGIPQKRKRWVMASKRGP-LSI-----pDPP I QSEDDYPTIRD-----TIDHLPPIDAGEVSD-----										
TLDAADYGVPOHRRRTFIVGIRKDFDEEFq---fpFPTHGPDSPFDNEQ-----VTAGEALKDIDSEG-A-----										
TLNAVNYGVPOTRKRLFILGVNDIAPPdq-----wepPRVCNEGQ-----qrltdinggnwle gyTTAHEALSDLPEFLKSqp-----										
slpu--YGIPO+Rcshhhgh+ps-----phs st-----thttllsl-tt-----										

Выделить минус блок труднее.

По возможности, надо избегать того, чтобы внутри него не оказалось К-подблока

В минус блоках мы не предполагаем наличие гомологичных остатков в любой позиции. Группировка остатков в одну позицию имеет целью сократить длину выравнивания, не более.

Но программы выравнивания пытаются выровнять невыравниваемое :(

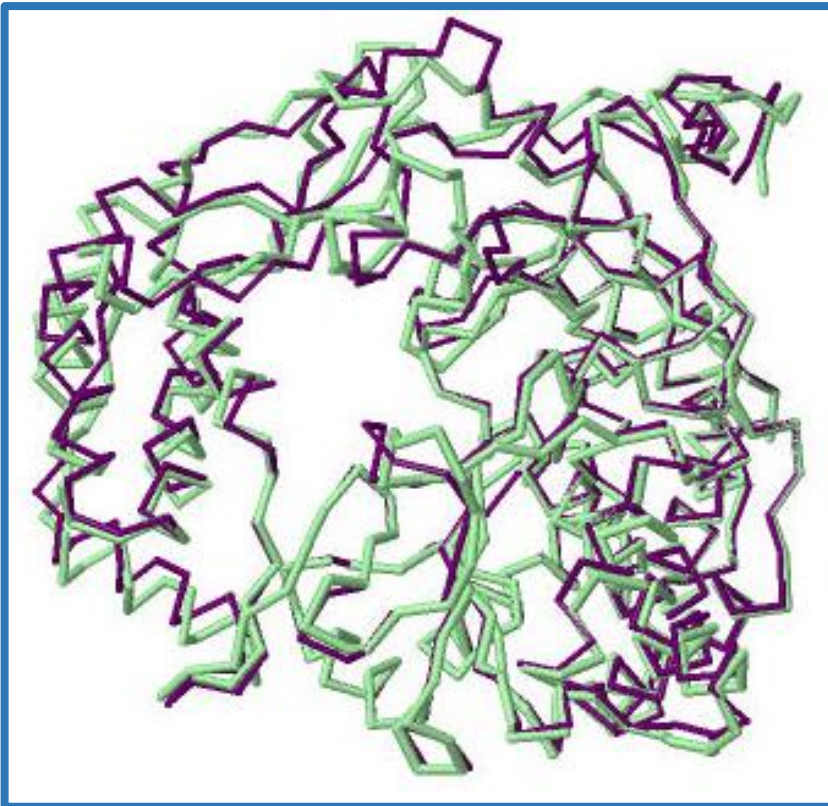
Смотрим выравнивание из
MethyltransferaseD12

Другие БД доменов

- PROSITE
- Smart
- Interpro – объединение всех БД белковых доменов

7. Проверка выравнивания по совмещению полипептидных остовов 3Dструктур

Известно, что структура консервативнее последовательностей



Совмещение полипептидных цепей
двух RdRP +РНК вирусов.

Не успел посмотреть каких. Выбрал случайно
из совмещения пяти структур

КОНЕЦ ПРЕЗЕНТАЦИИ

Страница домена Methylase_S (PF01420)

Family: *Methylase_S* (PF01420)

38 architectures 7284 sequences 3 interactions 2296 species 14 structures

Summary

Domain organisation
Clan
Alignments
HMM logo
Trees
Curation & model
Species
Interactions
Structures

Jump to...
enter ID/acc **Go**

Summary: Type I restriction modification DNA specificity domain

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

No Wikipedia article **Pfam** **InterPro**

This tab holds the annotation information that is stored in the Pfam database. As we move to using Wikipedia as our main source of annotation, the contents of this tab will be gradually replaced by the Wikipedia tab.

Type I restriction modification DNA specificity domain

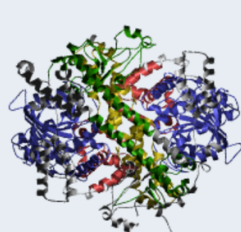
Provide feedback

This domain is also known as the target recognition domain (TRD). Restriction-modification (R-M) systems protect a bacterial cell against invasion of foreign DNA by endonucleolytic cleavage of DNA that lacks a site specific modification. The host genome is protected from cleavage by methylation of specific nucleotides in the target sites. In type I systems, both restriction and modification activities are present in one heteromeric enzyme complex composed of one DNA specificity subunit (this family), two modification (M) subunits and two restriction (R) subunits [2].

Literature references

1. Bickle TA, Kruger DH; , Microbiol Rev 1993;57:434-450.: Biology of DNA restriction. [PUBMED:8336674](#) [EPMC:8336674](#)
2. Janscak P, Bickle TA; , J Mol Biol 1998;284:937-948.: The DNA recognition subunit of the type IB restriction-modification enzyme EcoAI tolerates circular permutations of its polypeptide chain. [PUBMED:9837717](#) [EPMC:9837717](#)

Example structure
PDB entry 2Y7C: Atomic model of the Ocr-bound methylase complex from the Type I restriction-modification enzyme EcoKI (M251). Based on fitting into EM map 1S34.
View a different structure:
2Y7C ▾



Страница доменных архитектур белков, содержащих Methylase_S

Family: *Methylase_S* (PF01420)

38 architectures 7284 sequences 3 interactions 2296 species 14 structures

Summary

Domain organisation

Clan

Alignments

HMM logo

Trees

Curation & model

Species

Interactions


Structures


Jump to...


enter ID/acc **Go**


Domain organisation


Below is a listing of the unique domain organisations or architectures in which this domain is found. [More...](#)

There are 2509 sequences with the following architecture: Methylase_S x 2
[D3S0T0_FERPA](#) [Ferroglobus placidus (strain DSM 10642 / AEDII12DO)] Restriction modification system DNA specificity domain protein {ECO:0000313|EMBL:ADC66321.1} (421 residues)

[Show](#) all sequences with this architecture.

There are 1903 sequences with the following architecture: Methylase_S
[S0GN70_9PORP](#) [Parabacteroides goldsteinii dnLKV18] Uncharacterized protein {ECO:0000313|EMBL:E0S20219.1} (378 residues)

[Show](#) all sequences with this architecture.

There are 108 sequences with the following architecture: N6_Mtase, Methylase_S
[W5WMF1_9PSEU](#) [Kutzneria albida DSM 43870] Uncharacterized protein {ECO:0000313|EMBL:AH102043.1} (623 residues)

[Show](#) all sequences with this architecture.

There are 40 sequences with the following architecture: N6_Mtase, Methylase_S x 2
[A0A0H1S640_9MOLU](#) [Mycoplasmataceae bacterium RV_VA103A] Uncharacterized protein {ECO:0000313|EMBL:KLL02790.1} (1079 residues)

[Show](#) all sequences with this architecture.

There are 19 sequences with the following architecture: HSDR_N_2, N6_Mtase, Methylase_S x 2
[I3CEN9_9GAMM](#) [Beggiatoa alba B18LD] Type I restriction-modification system methyltransferase subunit {ECO:0000313|EMBL:EIJ42082.1} (1146 residues)

[Show](#) all sequences with this architecture.

Выравнивания Seed и Full

Seed та выборка последовательностей,
по которой строился HMM-профиль – математическое
описание выравнивания, позволяющее искать похожие
домены в новых последовательностях

КОНЕЦ ПРЕЗЕНТАЦИИ

Для чего строят множественные выравнивания?

```

      *           20           *           40           *           60
ROB_ECOLI : RYQFWHDFLGNATIPPVLYGLNETRPSQDKDDEQEVFYTTTAAQDQADGYVLTGHPVMLQ : 61
GADX_ECOLI : EW-----TLARIASELLMSPSLKKKKLREE-ETSYSQLITECRMQ----RALQLLIVI : 47
ENVY_ECOLI : YW-----NLRIVASSLCLSPSLLKKKKLKNE-NTSYSQIVTECRMR----YAVQMLLML : 47
YDEO_ECOLI : PW-----KLKDICDCLYISESLLKKKKLKQE-QTTFSQILLDARMQ----HAKNLLIRV : 47
APPY_ECOLI : QW-----HLKDIAELIYTSESLIKKKRLRDE-GTSFTEILRDTRMR----YAKKLLITS : 47
GADW_ECOLI : RW-----YLRDIAERMYTSESLLKKKKLQDE-NTCFSKIILLASRMS----MARRLLEL : 47
XYLR_ECOLI : HYIRNHACKGIKVVDQVLDAVGISRSNLEKRFKEEVGETIHAMTHAEKLE---KARSLLISI : 57
YDEC_BACSU : NWIHLHYVEKITLEDIAKAGQLSRSECCRYFKRMLNKTPLRYVMDYRIQ---KSLLLLQHH : 57
      5              s 3              e              6              a 66

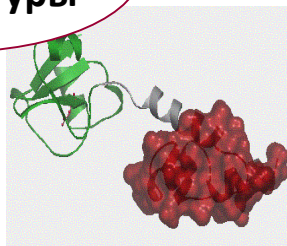
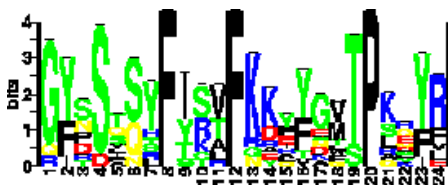
```

позволяет найти общее

мотивы, паттерны, профили

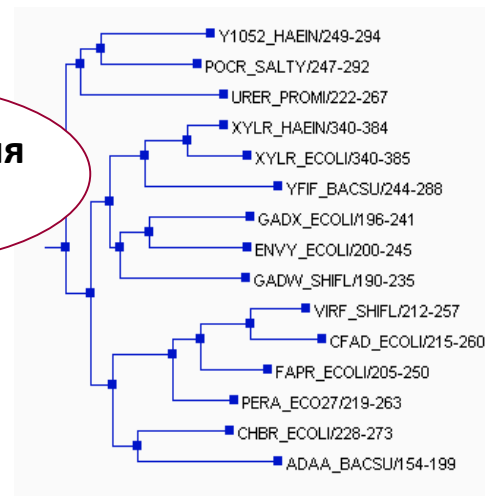
поиск
активного
центра

предсказание
3D-структуры



позволяет оценить эволюционные отношения

реконструкция
эволюции



Построение множественных выравниваний — необходимый этап решения многих задач молекулярной биологии

2. Jalview – см. инструкцию на сайте kodomo. В ней есть почти все нужное.

Алгоритмы множественного выравнивания

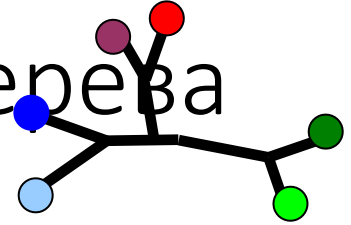
Иерархические и блочные

Иерархическое выравнивание многих последовательностей

- Основная идея: выравнивание двух выравниваний с помощью динамического программирования
- Этапы алгоритма
 - Построение направляющего дерева
 - Итерация выравнивания выравниваний
 - “Рафинирование” (refinement) выравнивания
- Результат – ГЛОБАЛЬНОЕ множественное выравнивание



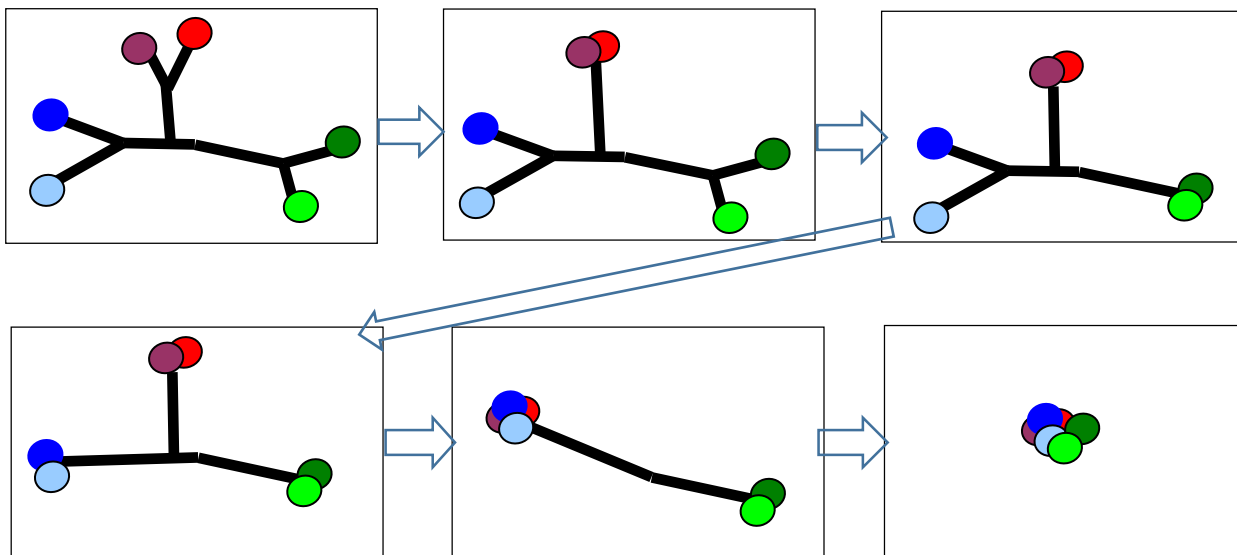
Построение направляющего дерева



- Для ВСЕХ ПАР последовательностей строится парное выравнивание.
- Вес парного выравнивания пересчитывается в расстояние между последовательностями:
 - чем больше вес, тем меньше расстояние;
 - расстояние между совпадающими последовательностями равно 0.
- Получается матрица расстояний между послед-ми
- Есть алгоритмы, превращающие матрицу попарных расстояний в дерево.
 - Расстояния между листьями по дереву отражают сходство последовательностей

Итерация выравнивания двух выравниваний

- Лист дерева – последовательность или выравнивание нескольких последовательностей



- На каждой итерации для выравнивания выбираются самые близкие “листья” – выравнивания
- Длина ветки до построенного выравнивания равна длине ветки до развилки плюс половина расстояния между листьями

“Рафинирование” выравнивания

- Главная беда иерархического выравнивания – ошибка на первых итерациях никогда не будет исправлена!
- Значит, нужны алгоритмы исправления ошибок – “рафинирования” выравнивания (refinement)
 - На вход подается множественное выравнивание
 - Каким-либо образом выравнивание делится на две части по горизонтали
 - эти части заново выравниваются друг против друга
 - если новое выравнивание лучше (нужен параметр!), то оно принимается, иначе – остается прежнее
 - простейший вариант: каждая последовательность выравнивается против всех остальных пока выравнивание не перестанет улучшаться

Блочное выравнивание

- Идея: найти в последовательностях сходные участки без гэпов. Построить из них блоки
- Найти непротиворечивый набор блоков
- Выровнять участки между блоками

Такой подход реализован в
PSI-BLAST (Position-Specific Iterated BLAST)
для улучшения результатов поиска

Алгоритмы множественного выравнивания

- Не решают формальную задачу – в отличие от алгоритмов парного выравнивания
- Нет гарантии, что построенное выравнивание превосходит все остальные по какому-то параметру, оценивающему качество выравнивания

ЭВРИСТИЧЕСКИЕ АЛГОРИТМЫ!

- Есть разные параметры качества выравнивания, например, сумма весов парных выравниваний.
 - Построение множественного выравнивания с помощью динамического программирования возможно
 - Число операций $C \cdot n^k$ где k – число последовательностей, n – средняя длина входных последовательностей (геометрическое среднее). Прикиньте сколько времени потребуется, если $n = 100$ и $k = 100$

Программы

- Muscle
- Mafft – быстрая
- T-coffee, M-coffee (<http://www.tcoffee.org/Projects/tcoffee/>)
 - умеет интегрировать разную информацию: парные локальные выравнивания, сходство 3D структур некоторых белков – если известны.
- ClustalW, более современная версия этого алгоритма --- ClustalO
- Probcons
- DiAlign – блочное, начинается с парных блоков
- AliBee – блочное
- И много других

КОНЕЦ