

Программы множественного выравнивания

На материале последовательностей гомологичных белков.

MSA – Multiple Sequence Alignment

План

1. +Гомология

2. +Чего хотим от MSA

- Реконструкцию эволюция
- Реконструкцию сходства хода полипептидных цепей 3D структур
- Иногда эволюция \neq 3D

3. Алгоритмы и программы

- Прогрессивное выравнивание
- Итеративное
- Регрессивное

4. Сравнение выравниваний тех же последовательностей

5. Верификация MSA.

- 3D совмещение
- Базы данных (MSA benchmarks)

6. Проект блочного MSA

1. Гомология белков?

DNA_methylase (PF00145)

310 architectures

20330 sequences

0 interactions

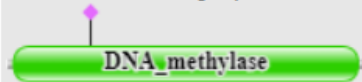
5673 species

Domain organisation

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

There are 9877 sequences with the following architecture: DNA_methylase

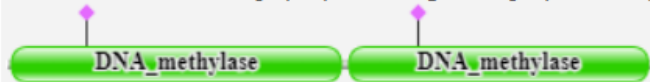
[W8T222_PEPAC](#) [Peptoclostridium acidaminophilum DSM 3953] Cytosine-specific methyltransferase {ECO:0000256|RuleBase:RU000417} (367 residues)



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

There are 2373 sequences with the following architecture: DNA_methylase x 2

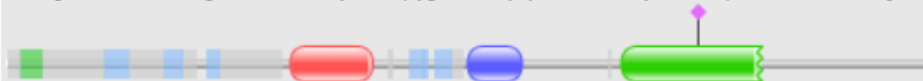
[A0A269TJK7_9MOLU](#) [Mycoplasmag agassizii] Cytosine-specific methyltransferase {ECO:0000256|RuleBase:RU000417} (664 residues)



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

There are 1510 sequences with the following architecture: PWWP, ADD_DNMT3, DNA_methylase

[W5QA98_SHEEP](#) [Ovis aries (Sheep)] DNA (cytosine-5-)-methyltransferase {ECO:0000256|ARBA:ARBA00011975} (948 residues)



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

There are 506 sequences with the following architecture: BAH, Chromo, DNA_methylase

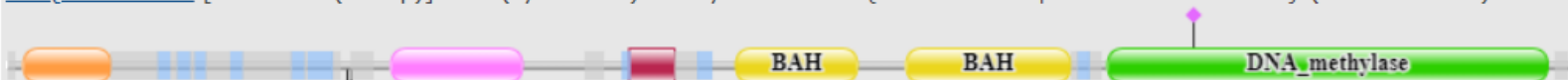
[A0A2G5D446_AQUCA](#) [Aquilegia coerulea (Rocky mountain columbine)] Uncharacterized protein {ECO:0000313|EMBL:PIA38263.1} (854 residues)



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

There are 414 sequences with the following architecture: DMAP_binding, DNMT1-RFD, zf-CXXC, BAH x 2, DNA_methylase

[W5Q1C4_SHEEP](#) [Ovis aries (Sheep)] DNA (cytosine-5-)-methyltransferase {ECO:0000256|PIRNR:PIRNR037404} (1612 residues)




Гомология белков:


- С одинаковой доменной архитектурой
- Доменов из одного семейства из разных белков


DNA_methylase (PF00145) 310 architectures 20330 sequences <= 20330 673 species

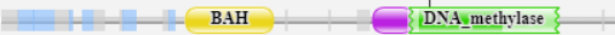
Domain organisation


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

There are 9877 sequences with the following architecture: DNA_methylase
[W8T222_PEPAC](#) [Peptoclostridium acidaminophilum DSM 3953] Cytosine-specific methyltransferase {ECO:0000256|RuleBase:RU000417} (367 residues)
 <= 9877
[Show all sequences with this architecture.](#)

There are 2327 sequences with the following architecture: DNA_methylase x 2
[AQA269TJK7_9MOLU](#) [Mycoplasma agassizii] Cytosine-specific methyltransferase {ECO:0000256|RuleBase:RU000417} (664 residues)
 <= 2327
[Show all sequences with this architecture.](#)

There are 1510 sequences with the following architecture: PWWP, ADD_DNMT3, DNA_methylase
[W5QA98_SHEEP](#) [Ovis aries (Sheep)] DNA (cytosine-5)-methyltransferase {ECO:0000256|ARBA:ARBA00011975} (948 residues)
 <= 1510
[Show all sequences with this architecture.](#)

There are 506 sequences with the following architecture: BAH, Chromo, DNA_methylase
[AQA2G5D446_AQUCA](#) [Aquilegia coerulea (Rocky mountain columbine)] Uncharacterized protein {ECO:0000313|EMBL:PIA38263.1} (854 residues)
 <= 506
[Show all sequences with this architecture.](#)

There are 414 sequences with the following architecture: DMAP_binding, DNMT1-RFD, zf-CXXC, BAH x 2, DNA_methylase
[W5Q1C4_SHEEP](#) [Ovis aries (Sheep)] DNA (cytosine-5)-methyltransferase {ECO:0000256|PIRNR:PIRNR037404} (1612 residues)
 <= 414

2. Чего хотим в результате MSA гомологичных доменов белка?

- Семейство доменов – единица локальной эволюции белков (локальная = точечные замены, делеции и вставки; по большей части короткие)
- В эволюции редко, но бывают также крупные перестройки (слияния и разделения белков, дубликации, рекомбинации, крупные делеции, вставки; много примеров циклических перестановок)

1. Хотим от MSA реконструкцию эволюции доменов семейства от общего предка: в колонке выравнивания потомки одного а.к.о. общего предка

Чего хотим в результате MSA гомологичных доменов белка?

- 3D структура белков консервативнее их последовательности.

2. Хотим от MSA чтобы C-alpha атомы а.к.о. из одной колонки при совмещении 3D структур этих белков оказались в одном положении

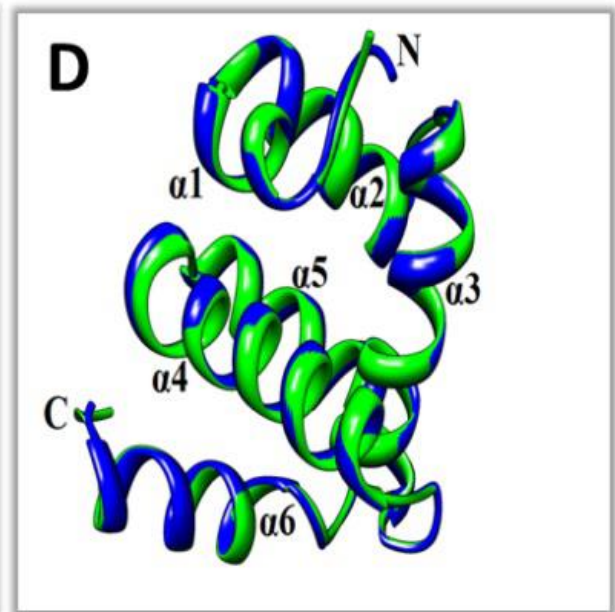
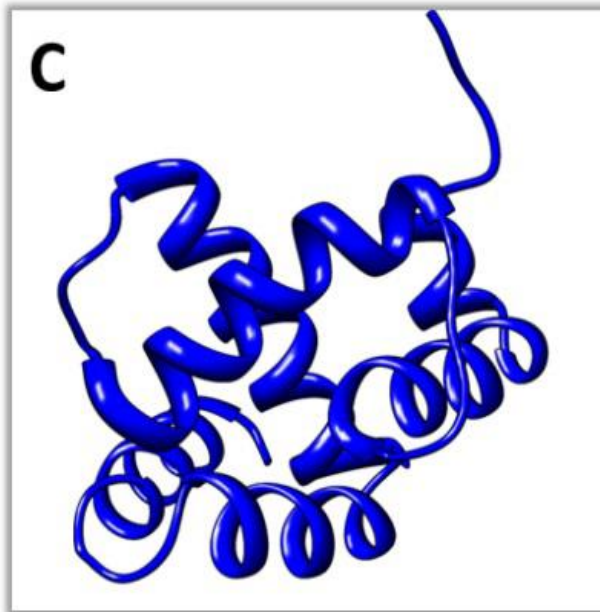
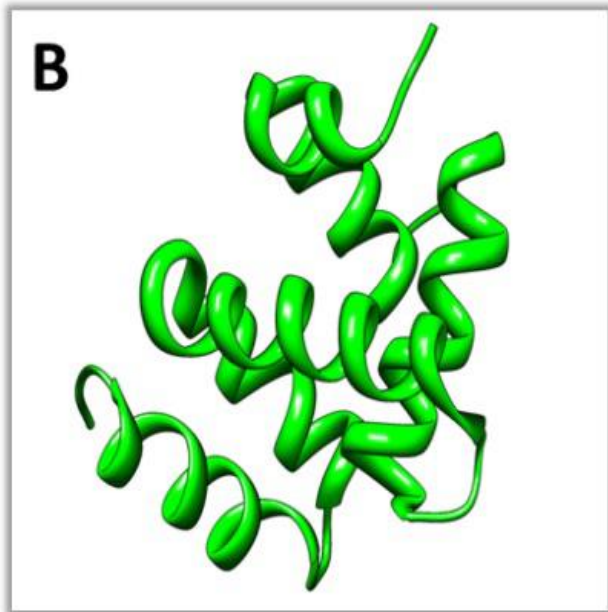
A

```

      *           20           *           40           *           60           *
MsepCSP5_ : MNSFTVLCLFALVALAVARPDGKYTD RYDSV NLDQILSNRRLIVPYIKMLDQ GKCTPDGKELKTHIREA : 70
MbraCSFA6 : -----EDKYTDKYDNINLDEILANKRLIVAYVNCVMERGKCSPEGKELKEHLQDA : 50
              KYTD YD  NLD IL N RLIV Y  C      GKC P GKELK H  A

      80           *           100          *           120           *
MsepCSP5_ : LEQDCAKCTKAQRDGT RQVMGHLINHEVDYWNELKAKYDFKNLYSTHHEQELRKLKQ----- : 127
MbraCSFA6 : IENGCKKCTENQEKGAYRVIEHLIKNEIEIWRELTAKYDPTGNWRKKYEDRAKAAGIVIP EE : 112
              E  C KCT  Q  G  V  HLI  E  W  EL  AKYDF           K  E

```



2.3 Хотим в результате MSA ГОМОЛОГИЧНЫХ ДОМЕНОВ белка:

1. - реконструкцию эволюции доменов семейства от общего предка: в колонке выравнивания потомки одного а.к.о. общего предка

НЕ ПРОВЕРЯЕМО!

2. - чтобы C-alpha атомы а.к.о. из одной колонки при совмещении 3D структур этих белков оказались в одном положении

ПРОВЕРЯЕМО:

а. Подкопить денюжат на RCSA всех белков выравнивания

б. Проверить на белках с известной 3D

с. Предсказать 3D структуры всех белков и сравнить их.

БЕДА: основой предсказания служит сходство фрагментов последовательности с фрагментами белков с известной 3D

Хотим в результате MSA ГОМОЛОГИЧНЫХ доменов белка

Хотеть не вредно!

Верно ли, что эволюция = сохранение хода
полипептидной цепи в 3D?

Отчасти, да. Убеждают примеры совмещаемых
структур со слабым сходством последовательностей.

Отчасти нет. См. ниже.

Chatzou et al., Multiple sequence alignment modeling: methods and applications, 2016

If the alignment is meant to be a structural model, aligned residues should have comparable positions in their respective 2D or 3D structures. If the alignment is functional, as may happen when analyzing genomic data, aligned positions are expected to support similar functions. Even though it is reasonable to expect a significant overlap between these criteria,

It must be stressed that the complexity of evolutionary forces is such that their full agreement cannot be taken for granted.

For instance, two structures may be similar as a consequence of convergent evolution but non-homologous from an evolutionary point of view.

Например, сравнительно длинные вставки в одном и том же месте структуры могут произойти независимо и не быть гомологичными. Возникают потому, что вставки в этом месте не нарушают структуру и функцию белка

ПРИМЕР



Younas A et al., A chemosensory protein MsepCSP5 involved in chemoreception of oriental armyworm *Mythimna separata*. *Int J Biol Sci.* 2018

Mythimna separata is one of the most serious pests affecting the quality and quantity of crop yield [25, 26]

The binding sites of MsepCSP5 to candidate volatiles were well predicted by three-dimensional structure modeling and molecular docking experiments.

Pursuing further, biological activities of *M. separata* to highly bound compounds elicited strong behavioral responses, such as alcoholic compounds displayed strong attractiveness whereas terpenes showed repellency to *M. separata*

Multiple Sequence Alignment of CsCP against Plant CPs

3UBE/1-222 1 - - APASIDWRKK GAVTSVKDQ GAC G M C W A F G A T G A I E G I D A I T T G R L I S V S E Q Q I V D C 56
1S4V/1-223 1 - TVPASVDWRKK GAVTSVKDQ G Q C G S C W A F S T I V A V E G I N Q I K T N K L V S L S E Q E L V D C 57
1CQD/1-221 1 D D L P D S I D W R E N G A V V P V K N Q G G C G S C W A F S T V A A V E G I N Q I V T G D L I S L S E Q Q L V D C 58
1IWD/1-215 1 - - L P S F V D W R S K G A V N S I K N Q K Q C G S C W A F S A V A A V E S I N K I R T G Q L I S L S E Q E L V D C 56
1POP/1-212 1 - - I P E Y V D W R Q K G A V T P V K N Q G S C G S C W A F S A V V T I E G I I K I R T G N L N Q Y S E Q E L L D C 56

conservation



3UBE/1-222 57 D T - A A A A A A G C D A D D A F R W V I T N G G I A S D A N Y P Y T G V D G T C D L N - - K P I A A R I D G Y T N 111
1S4V/1-223 58 D T D Q N Q G C N G C L M D Y A F E F I K Q R G G I T T E A N Y P Y E A Y D G T C D V S K E N A P A V S I D G H E N 115
1CQD/1-221 59 T T - A N H G C R G G W M N P A F Q F I V N N G G I N S E E T Y P Y R G Q D G I C N S T - V N A P V V S I D S Y E N 114
1IWD/1-215 57 D T - A S H G C N G G W M N N A F Q Y I I T N G G I D T Q Q N Y P Y S A V Q S C K P - - Y R L R V V S I N G F Q R 111
1POP/1-212 57 D R - R S Y G C N G G Y P W S A L Q L V A Q Y G - I H Y R N T Y P Y E G V Q R Y C R S R E K G P Y A A K T D G V R Q 112

conservation



3UBE/1-222 112 V P - N S S S A L L D A V A K Q P V S V N I Y T S S T S F Q L Y T G P G I F A G S S C S D D P A T V D H T V L I V G 168
1S4V/1-223 116 V P E N D E N A L L K A V A N Q P V S V A I D A G G S D F Q F Y S E - G V F T G - S C G - - - T E L D H C V A I V G 168
1CQD/1-221 115 V P S H N E Q S L Q K A V A N Q P V S V T M D A A G R D F Q L Y R S - G I F T G - S C N - - - I S A N H A L T V V G 167
1IWD/1-215 112 V T R N N E S A L Q S A V A S Q P V S V T V E A A G A P F Q H Y S S - G I F T G - P C G - - - T A Q N H C V V I V G 164
1POP/1-212 113 V Q P Y N Q G A L L Y S I A N Q P V S V V L Q A A G K D F Q L Y R G - G I F V G - P C G - - - N K V D H A V A A V G 165

conservation



3UBE/1-222 169 Y G S N G T N A D Y W I V K N S W G T E W G I D G Y I L I R R N T N R P D G V C A I D A W G S Y P T K S T S - 222
1S4V/1-223 169 Y G T T I D G T K Y W T V K N S W G P E W G E K G Y I R M E R G I S D K E G L C G I A M P A S Y P I K K S S N 223
1CQD/1-221 168 Y G T E N D - K D F W I V K N S W G K N W G E S G Y I R A E R N I E N P D G K C G I T R F A S Y P V K K G T N 221
1IWD/1-215 165 Y G T Q S G - K N Y W I V R N S W G Q N W G N Q G Y I W M E R N V A S S A G L C G I A Q L P S Y P T K A - - - 215
1POP/1-212 166 Y G - - - - P N Y I L I K N S W G T G W G E N G Y I R I K R G T G N S Y G V C G L Y T S S F Y P V K N - - - 212

conservation



-- Conserved

-- Catalytic

-- Active site

-- 3UBE (*Crocus sativus*)

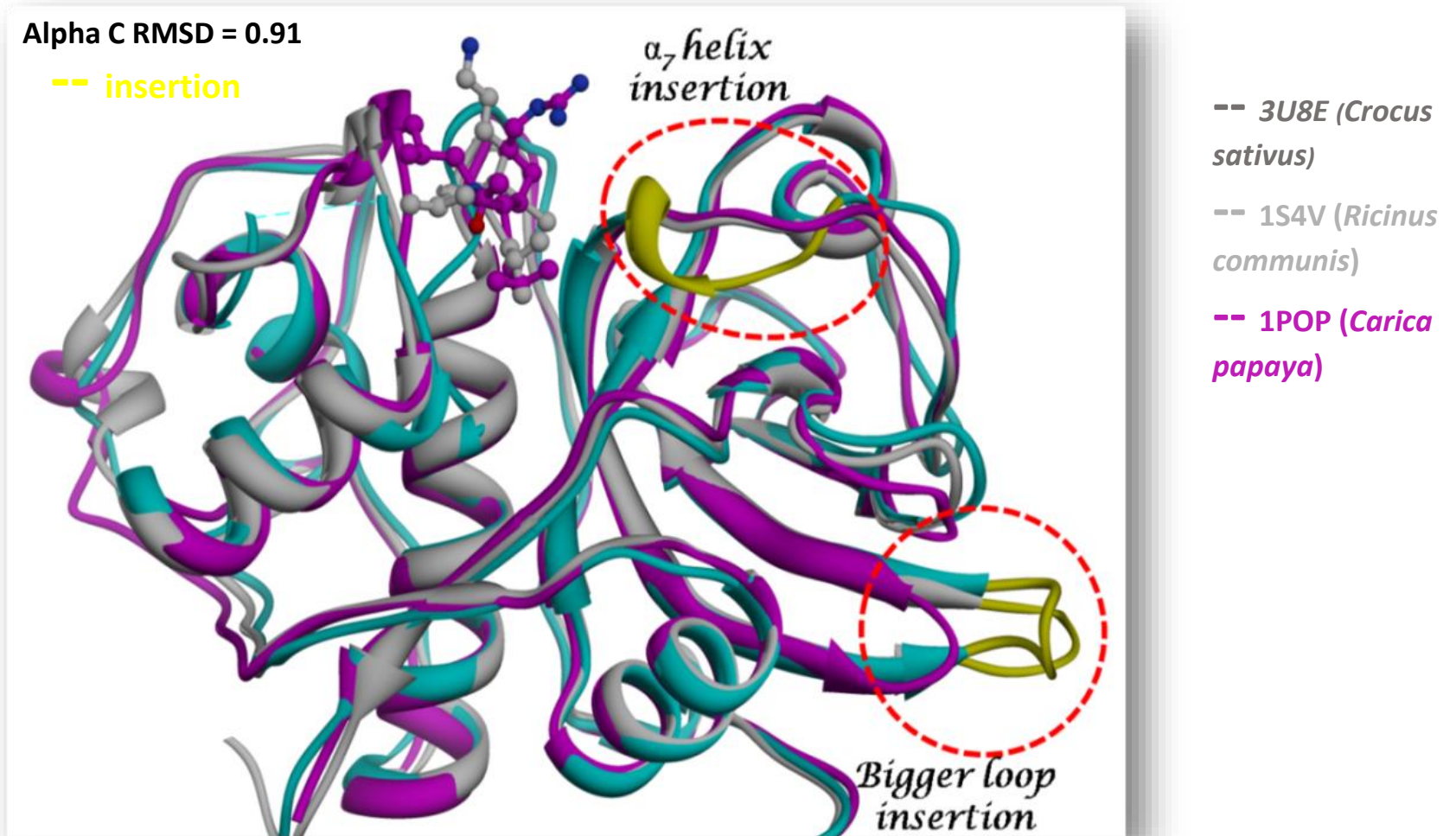
-- 1S4V (*Ricinus communis*)

-- 1CQD (*Zingiber officinale*)

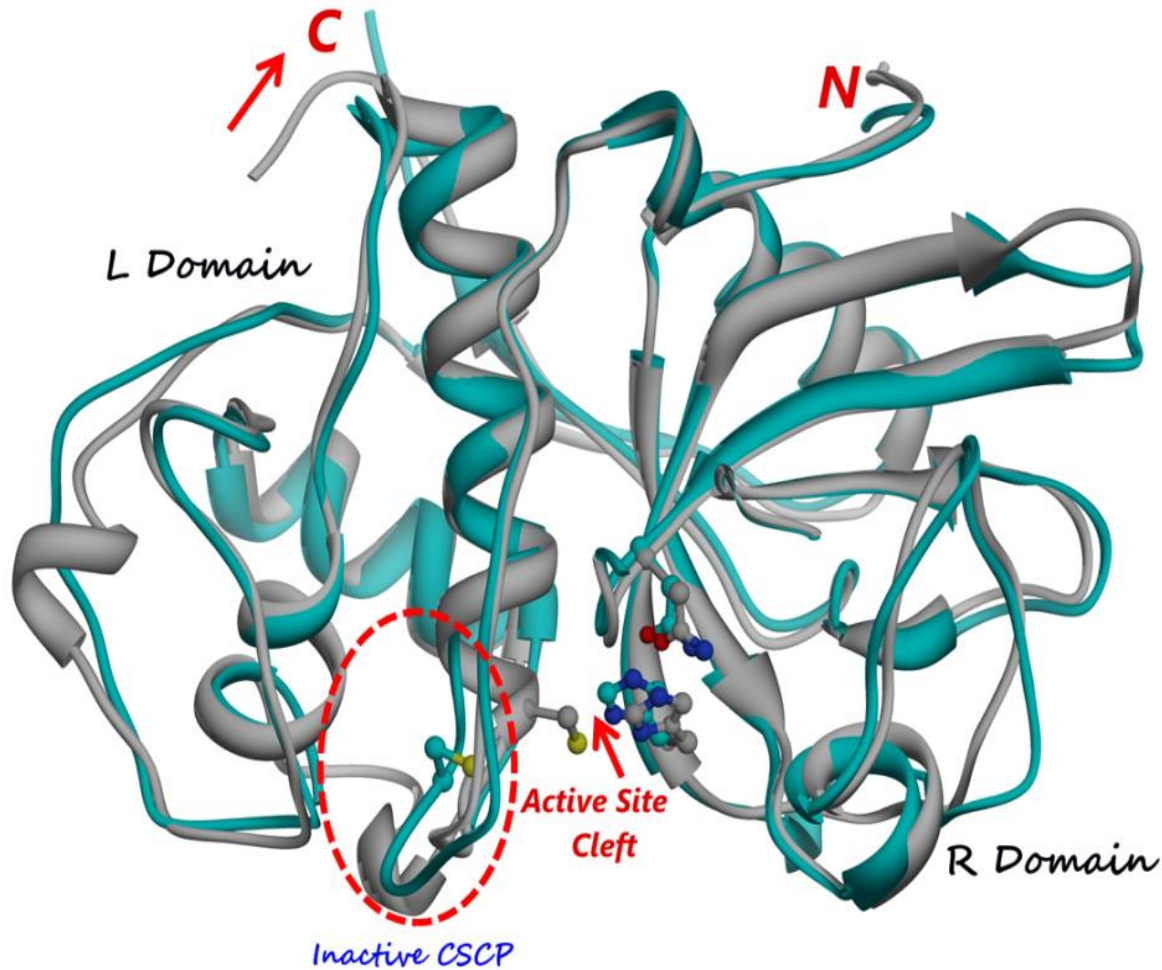
-- 1IWD (*Ervatamin B*)

-- 1POP (*Carica papaya*)

Conserved Structure of CsCP Superimposition with Plant CPs



Overall Topology of CsCP



-- Inactive Form

-- Active Form

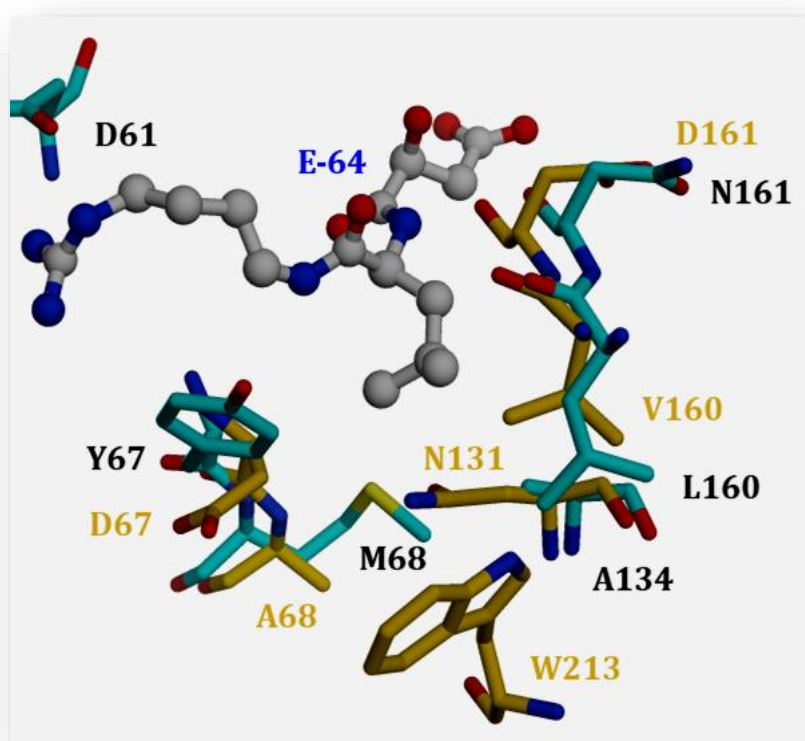
Catalytic Triad:
C25, H162, and N183
(Active Site Cleft)

Inactive Form:
C25 bridged with C22

Active Form:
C22 bridged with C64
C25 is free

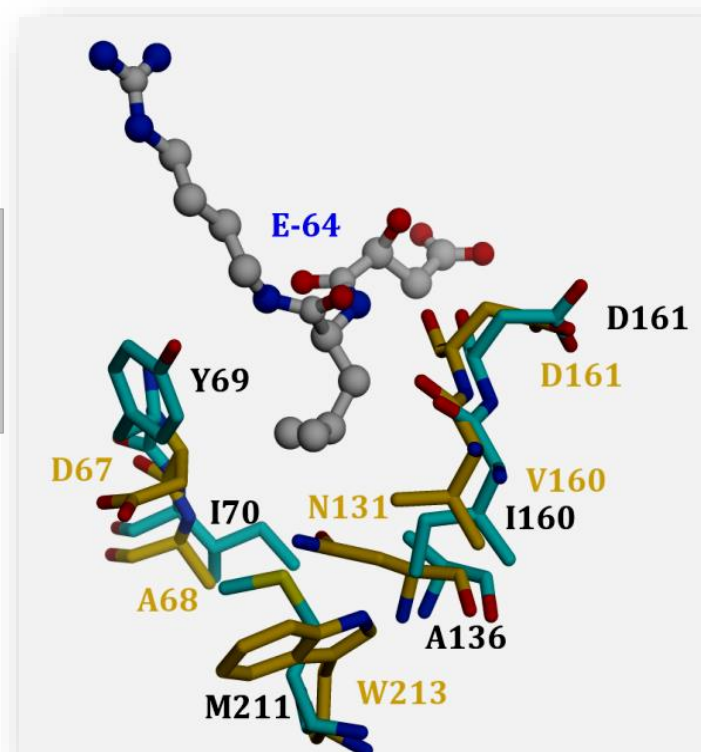
KEY Findings Identification of S₂ Pocket Residues of CSCP

Cathepsin K-E64



-- Others
-- CsCP

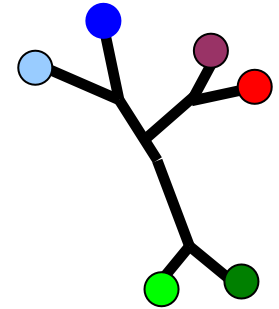
Actinidin-E64



3. Алгоритмы и программы MSA

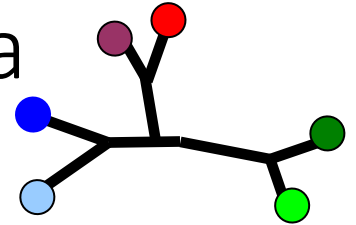
- Прогрессивное выравнивание (ClustalW)
- Итеративное выравнивание (MUSCLE, PRALINE, IterAlign)
- Методы, основанные на согласованности (MAFFT, ProbCons)
- Методы, использующие структурные данные

Иерархический алгоритм выравнивания 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

Книга

Multiple Sequence Alignment Methods and Protocols

Edited by Kazutaka Katoh

Research Institute for Microbial Disease, Osaka
University, Osaka, Japan

https://link.springer.com/protocol/10.1007/978-1-0716-1036-7_17

Free

Алгоритмы и программы MSA

Progressive Method

A reasonable and widely used heuristic is the progressive method [2–4].

In this strategy, a tentative tree, called a “guide tree,” is built based on an all-to-all approximate comparison. Then, the sequences are aligned from the leaves to the root on the tree, in a group-to-group manner.

When the calculation reaches the root, the full MSA is obtained.

Many MSA programs use the progressive method as a part of the calculation. Among them, PRANK (Chapter 2) has a notable point that it rigorously considers insertions and deletions on the guide tree.

Алгоритмы и программы MSA

Iterative Refinement

The progressive method has a well-known problem that errors can occur in early steps (i.e., close to a leaf) of the guide tree, and those errors remain in the final step (the root of the guide tree).

One effective solution is to correct this type of mistake is iterative refinement [5–7]. The procedure is:

- (i) construct an initial MSA;
 - (ii) divide the MSA into two groups; (iii) re-align the two groups;
- repeat (ii) and (iii). This technique is used in Prrn5 (Chapter 5), Clustal Omega (Chapter 1), MAFFT (Chapter 11), and MSAProbs (Chapter 3).

Алгоритмы и программы MSA

Consistency

Another important idea to overcome the limitations of the progressive method is consistency [8, 9]. In the tree-based consistency transformation technique, proposed first in Notredame et al. [10], when aligning two sequences (A and B), other sequences (e.g., C) in the dataset are also used.

That is, in addition to direct alignment between A and B, an alternative alignment between A and B is synthesized by using alignments AC and BC.

Alignment AB is recalculated by considering such alternative alignments and used in the progressive alignment step. As a result, alignment errors in early steps are efficiently suppressed.

This method was further elaborated to use probabilistic pairwise alignments by a pair hidden Markov model in ProbCons [11]. MSAProbs (Chapter 3) represents this type of MSA method and gives highly accurate MSAs.

Популярные программы

1. MAFFT
2. Clustal Omega
3. ClustalW
4. Muscle
5. Probcons
6.

4. Сравнение двух выравниваний тех же последовательностей

ььь

Какие выравнивания тех же последовательностей совпадают?

	1	2	3	4	5	6	7	8	9	10	11	12	13
Seq1	M	K	F	R	-	S	S	H	Y	A	-	S	
Seq2	M	K	Y	R	R	R	-	H	Y	A	-	S	
Seq3	M	E	F	R	R	R	R	S	H	Y	A	-	R

	1	2	3	4	5	6	7	8	9	10	11	12
Seq1	M	K	F	-	R	S	S	H	Y	A	S	
Seq2	M	K	Y	R	R	R	-	H	Y	A	S	
Seq3	M	E	F	R	R	R	R	S	H	Y	A	R

	1	2	3	4	5	6	7	8	9	10	11	12	13
Seq1	M	K	F	R	-	S	S	H	Y	A	-	S	
Seq2	M	K	Y	R	-	R	S	H	Y	A	-	S	
Seq3	M	E	F	R	R	R	R	S	H	Y	A	-	R

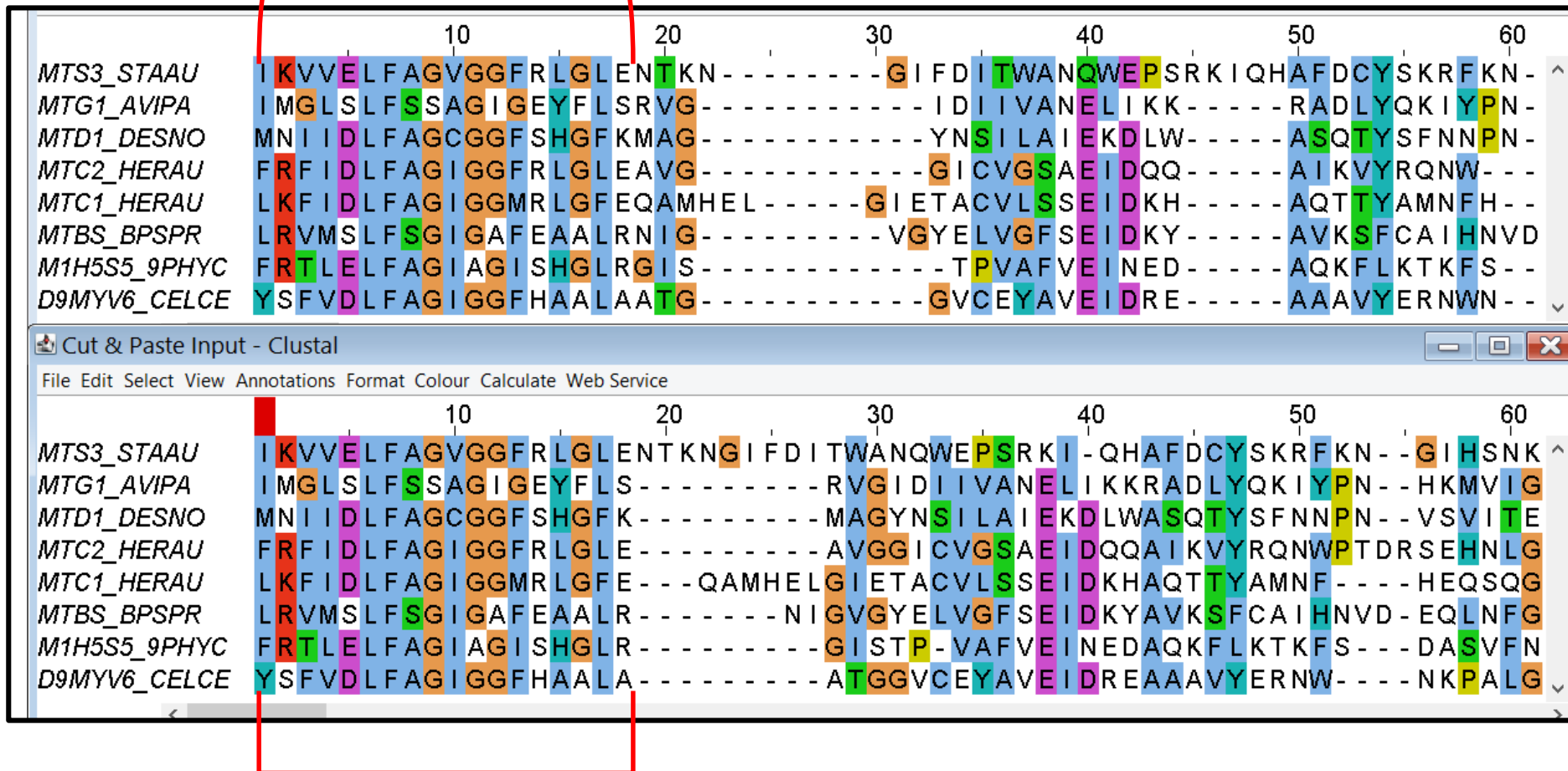
	1	2	3	4	5	6	7	8	9	10	11	12	
Seq1	-	M	K	F	R	-	S	S	H	Y	A	S	
Seq2	-	M	K	Y	R	-	R	S	H	Y	A	S	
Seq3	-	M	E	F	R	R	R	R	S	H	Y	A	R

Колонка i выравнивания X совпадает с колонкой j выравнивания Y если в них – те же самые остатки; те же самые значит – с теми же номерами, а не с теми же буквами!

ПРИМЕР

Выравнивание PF00145-seed из PFAM

Прорежённое до 8 посл-й с помощью EDIT => remove redundancy



Выравнивание с помощью MAFFT тех же 8 посл-й

Совпадение выравниваний в позициях 1 – 18. Больше нет (в показанной части)

Продолжение

		70		80		90		100		110																																							
<i>MTS3_STAAU</i>	---	G	I	H	S	N	K	D	I	A	Q	V	S	D	E	---	E	M	A	N	---	T	E	A	D	M	I	V	G	G	F	P	C	Q	D	Y	S	V	A	R	S	---	---						
<i>MTG1_AVIPA</i>	---	H	K	M	V	I	G	D	I	R	D	Q	R	I	F	N	K	V	L	N	I	A	L	T	N	-	Q	V	D	F	L	I	A	S	P	P	C	Q	G	M	S	V	A	G	K	N	R	D	V
<i>MTD1_DESNO</i>	---	V	S	V	I	T	E	D	I	T	T	L	D	P	G	---	D	L	K	I	S	V	S	D	---	V	D	G	I	I	G	G	P	P	C	Q	G	F	S	L	S	G	---	---	---				
<i>MTC2_HERAU</i>	---	-	P	T	D	R	S	E	H	N	L	G	D	I	T	-	T	L	Q	Q	L	P	A	---	-	H	D	L	V	V	G	G	V	P	C	Q	P	W	S	I	A	G	K	---	---				
<i>MTC1_HERAU</i>	---	-	-	-	E	Q	S	Q	G	D	I	T	Q	I	Q	---	---	-	D	F	P	---	S	---	F	D	F	L	L	A	G	F	P	C	Q	P	F	S	Y	A	G	K	---	---					
<i>MTBS_BPSPR</i>	D	-	-	E	Q	L	N	F	G	D	V	S	K	I	D	K	---	-	K	L	P	---	E	---	F	D	L	L	V	G	G	S	P	C	Q	S	F	S	V	A	G	H	---	---					
<i>M1H5S5_9PHYC</i>	---	D	A	S	V	F	N	D	V	T	K	F	T	K	S	---	---	D	F	P	E	D	---	I	D	M	I	T	A	G	F	P	C	T	G	F	S	I	A	G	S	---	---						
<i>D9MYV6_CELCE</i>	---	-	-	K	P	A	L	G	D	I	T	D	D	A	N	D	---	E	G	V	T	L	R	G	Y	D	G	P	I	D	V	L	T	G	G	F	P	C	Q	P	F	S	K	S	G	A	---	---	

		62		72		82		92		102																																								
<i>MTS3_STAAU</i>	N	-	-	G	I	H	S	N	K	D	I	A	Q	V	S	---	D	E	E	M	A	N	T	E	A	---	-	D	M	I	V	G	G	F	P	C	Q	D	Y	S	V	A	R	S	L	N	G	E		
<i>MTG1_AVIPA</i>	N	-	-	H	K	M	V	I	G	D	I	R	D	Q	R	I	F	N	K	V	L	N	I	A	L	T	N	-	Q	V	D	F	L	I	A	S	P	P	C	Q	G	M	S	V	A	---	-	G	K	N
<i>MTD1_DESNO</i>	N	-	-	V	S	V	I	T	E	D	I	T	T	L	D	P	G	D	L	K	I	S	V	S	D	V	---	-	D	G	I	I	G	G	P	P	C	Q	G	F	S	L	S	---	-	G	N	R		
<i>MTC2_HERAU</i>	T	D	R	S	E	H	N	L	G	D	I	T	T	L	Q	---	---	-	Q	L	P	A	H	---	-	-	-	D	L	V	V	G	G	V	P	C	Q	P	W	S	I	A	---	-	G	K	N			
<i>MTC1_HERAU</i>	-	-	-	H	E	Q	S	Q	G	D	I	T	Q	I	Q	---	---	-	D	F	P	S	F	---	-	-	-	D	F	L	L	A	G	F	P	C	Q	P	F	S	Y	A	---	-	G	K	Q			
<i>MTBS_BPSPR</i>	V	D	-	E	Q	L	N	F	G	D	V	S	K	I	D	---	-	-	K	K	K	L	P	E	F	---	-	-	D	L	L	V	G	G	S	P	C	Q	S	F	S	V	A	---	-	G	H	R		
<i>M1H5S5_9PHYC</i>	---	-	-	D	A	S	V	F	N	D	V	T	K	F	T	---	---	-	K	S	D	F	P	E	D	---	-	I	D	M	I	T	A	G	F	P	C	T	G	F	S	I	A	---	-	G	S	R		
<i>D9MYV6_CELCE</i>	---	-	-	N	K	P	A	L	G	D	I	T	D	D	A	-	N	D	E	G	V	T	L	R	G	Y	D	G	P	I	D	V	L	T	G	G	F	P	C	Q	P	F	S	K	S	---	-	G	A	Q

Take home message

To compare alignments A and B

1. Sort sequences by ID in both alignments.
2. Если в двух блоках без гэпов из разных выравниваний есть хотя бы одна одинаково выровненная позиция, то блоки выровнены одинаково
3. Если колонка букв в выравнивании A отличается хотя бы одной буквой от колонки букв в выравнивании B, то они выровнены не одинаково.

Следствие

- Два выравнивания A и B тех же последовательностей совпадают если в каждой колонке i ($i = 1, 2, \dots, N$) выравниваний A и B стоят те же самые буквы. Буквы не в смысле а.к.о., а в смысле номера буквы в последовательности (рис. 1.)
- Различие выравниваний A и B определяется числом колонок, для которых это не так и их расположением

Алгоритм сравнения двух выравниваний А и В тех же n последовательностей

1. Сортировка последовательностей по ID в обоих выравниваниях

2. Для выравнивания А шифруем каждую колонку i вектором $A(i)$

$$(i; A(i)), A(i) = (s_1, s_2, \dots, s_n)$$

n = число последовательностей

s_k зависит от того, стоит ли в колонке i в последовательности s_k буква или гэп (-);

если буква, то s_k = номер этой буквы в последовательности s_k ;

если гэп, то минус номер последней буквы прочтенной в последовательности s_k

3. Аналогично определяем $(j, B(j))$

Позиция i выравнивания А одинаково выровнена с позицией j выравнивания В если $A(i) = B(j)$

i может быть равно j , может быть не равно

4. Вычисляем список одинаково выровненных позиций (i,j) .

Процент длины этого списка от длины выравнивания – мера совпадения выравниваний.

Порядок позиций (i, j) в одном из выравниваний описывает расположение совпадающих частей выравнивания.

Готовые программы для сравнения выравниваний тех же последовательностей

- Нашел одну
VerAlign: a multiple sequence alignment assessment
tool

5. Верификация MSA

- По структуре
 - BaliBase
 - Программа FATCAT множественного совмещения структур белков. **Не нашел доступного сервиса.**
 - Реализация FATCAT в PDB для парных сравнений работает.
- Базы эталонных выравниваний для сравнения программ MSA
 - OxBanch **не удалось использовать**

Le et al., 2016

Table 1. The prediction accuracy for alignments of 200 sequences for 238 Pfam families. Aligner settings Prediction Accuracy (in %)

MAFFT L-INS-i	78.94 *
MAFFT—Default	78.19
MAFFT—Fast Mode	77.53 *
Clustal Omega—2 iter	78.36 *
Clustal Omega—1 iter	78.56 *
Clustal Omega—Default	78.63
MUSCLE—2 iter	78.17
MUSCLE—Default	78.13
MUSCLE—1 iter	77.29 *
PASTA—Default	78.70
T-Coffee—Default	78.45
Kalign 2—Default	77.93
Clustal W2—Default	77.13
HMMER—Default	77.86

For aligner settings from the same aligner, the sign (*) signifies that the score is significantly different (higher or lower) from the default score with $P < 0.01$ using the Wilcoxon signed rank test.

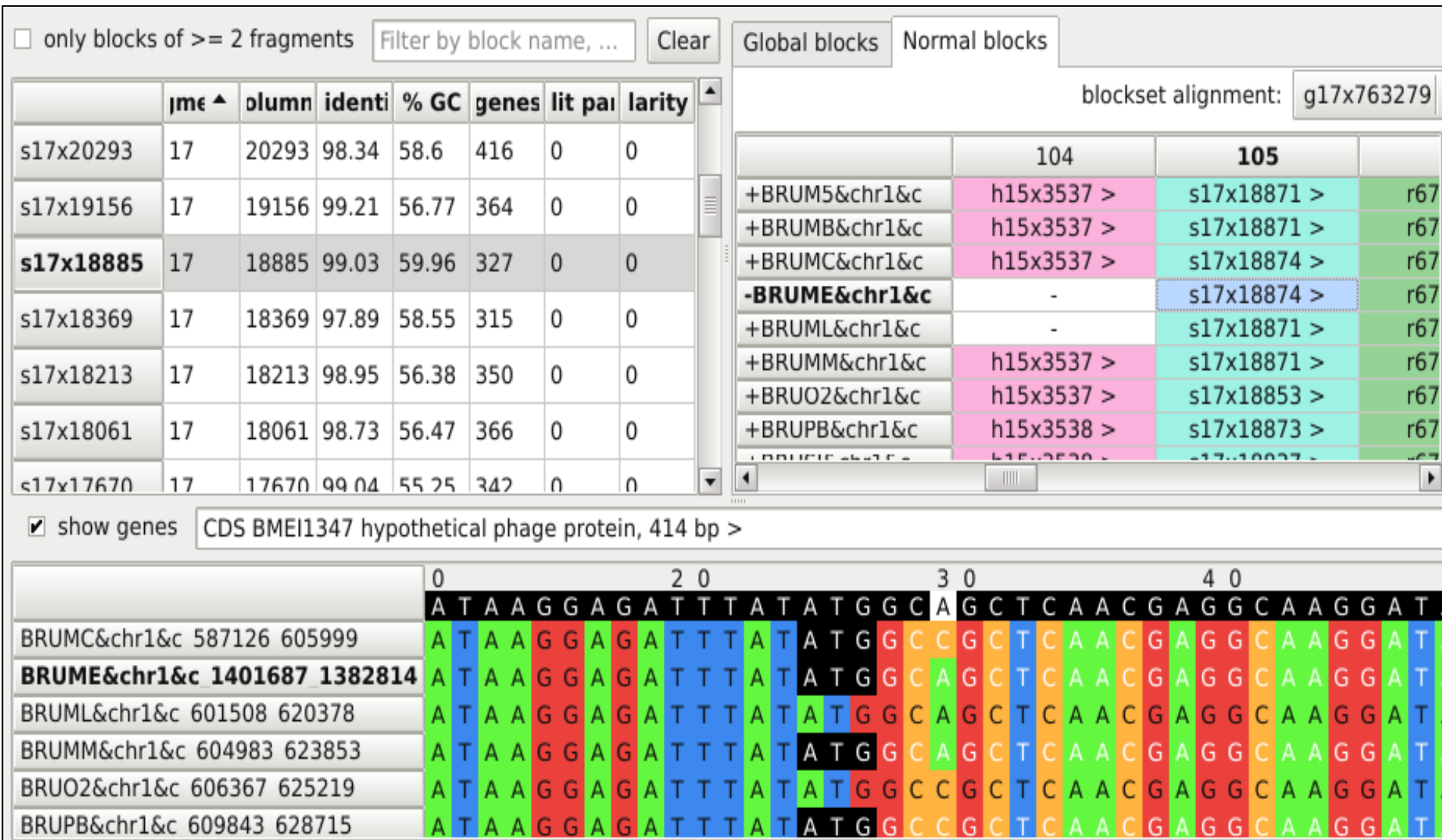
5. Проект блочного MSA

- Определение. MSA = набор блоков достоверного выравнивания во всех или части последовательностей. AAЛ
- Этап I: во входном выравнивании найти все блоки достоверного выравнивания.
- Этап II: найти все блоки достоверного выравнивания во входном множестве последовательностей
- Этап III: создать сервис для блочного MSA

Пример: блочное выравнивание геномов

- Борис Нагаев NPGexplorer

Визуализация выравнивания 17 геномов бруцелл (рис. из работы К.Худяковой)



Каждая ДНК представляется последовательностью блоков

Global blocks Normal blocks

blockset alignment:

	35	36	37	38	39
+BRUA1&chr1&c	s17x21286 >	-	s17x10111 >	-	s17x7319 >
+BRUA2&chr1&c	s17x21286 >	-	s17x10111 >	-	s17x7319 >
+BRUA8&chr1&c	s17x21286 >	-	s17x10111 >	-	s17x7319 >
-BRUA0&chr1&c	s17x21286 >	-	s17x10111 >	-	s17x7319 >
+BRUC2&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >
-BRUCA&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >
+BRUM5&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUMB&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUMC&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >
-BRUME&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUML&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUMM&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUO2&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUPB&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	-	s17x7319 >
+BRUSI&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >
+BRUSS&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >
+BRUSU&chr1&c	s17x21286 >	h13x107 >	s17x10111 >	h6x2653 >	s17x7319 >

Про типы блоков следующий слайд

КОНЕЦ ПРЕЗЕНТАЦИИ

ВЫРАВНИВАНИЕ

DNAK_THEAC 82 KFKVFDKEFTPQQISAFILQKIKKDA-EAFLGEPVNEAVITVPAYFNDNQR 131
DNAK_PICTO 82 KYKIFGKEYTPQQISAFILQKIKRDA-EAFLGEPVTDAVITVPAYFNDNQR 131
HSCA_ACIF2 116 RLRTVAGEKSPVEVSAEILRVLKERAVETLGGEPEGAVITVPAYFDEAQR 165
HSCA_ACIF5 116 RLRTVAGEKSPVEVSAEILRVLKERAVETLGGEPEGAVITVPAYFDEAQR 165

DNAK_THEAC 132 QATKDAGT IAGFDVKRIINEPTAAALAYGVDKSGKSEKILVFDLGGGTLDV 182
DNAK_PICTO 132 QATKDAGAIAGLNVRRINEPTAACLAYGIDKLNQTLKIVIYDLGGGTLDV 182
HSCA_ACIF2 166 QATKDAARLAGLNVLLAEPTAAAVAYGLDKGSEGI-FAIYDLGGGTFDI 215
HSCA_ACIF5 166 QATKDAARLAGLNVLLAEPTAAAVAYGLDKGSEGI-FAIYDLGGGTFDI 215

DNAK_THEAC 183 TIMDFGDGVFQVLSSTSGDTRLGGTDMDEAIVNYIADDFQKKEGIDLKDRS 233
DNAK_PICTO 183 TIMDFGQGVFQVLSSTSGDTHLGGTDMDEAIVNFLADNFQRENGIDLKDHHS 233
HSCA_ACIF2 216 SILRLQAGVFEVLATAGDSALGGDDMDHALAEWLMQE-EGGDASDPLW 262
HSCA_ACIF5 216 SILRLQAGVFEVLATAGDSALGGDDMDHALAEWLMQE-EGGDASDPLW 262

DNAK_THEAC 234 AYIRLRDAAEKAKIELSTTLSTDIDL PYITVTNSGPKHKIKMTLTRAKLEEL 284
DNAK_PICTO 234 AYIRLRDAAEKAKIELSTVLETEINL PYITATQDGPKHLQYTLTRAKFEEL 284
HSCA_ACIF2 263 RRQVLQQ-ARTAKEALSVAEET-MIVLTPSGRAAREIKLSRGRLES 307
HSCA_ACIF5 263 RRQVLQQ-ARTAKEALSVAEET-MIVLTPSGRAAREIKLSRGRLES 307

DNAK_THEAC 285 ISPIVERVKGPIDKALEGAKLKKTEITKLLFVGGPTRIPYVRKYVEDYLG I 335
DNAK_PICTO 285 IAPIVDRSKVPLDTALEGAKLKKGDIDKIILIGGPTRIPYVRKYVEDYFGR 335
HSCA_ACIF2 308 IQPVIQRS LPACRRALRDAGLKLDEIEGVVLVGGATRVPVAVRAMVEEFFRQ 358
HSCA_ACIF5 308 IQPVIQRS LPACRRALRDAGLKLDEIEGVVLVGGATRVPVAVRAMVEEFFRQ 358