

*А.В. Бутвиловский, Е.В. Барковский, В.Э. Бутвиловский*  
**Выравнивание аминокислотных и нуклеотидных последовательностей**

*Белорусский государственный медицинский университет*

Статья посвящена теоретическим аспектам выравнивания аминокислотных и нуклеотидных последовательностей. Рассмотрены основные принципы работы программ Clustal.

Ключевые слова: выравнивание, серия программ Clustal, матрицы выравнивания. Выравнивание аминокислотных или нуклеотидных последовательностей – это процесс сопоставления сравниваемых последовательностей для такого их взаиморасположения, при котором наблюдается максимальное количество совпадений аминокислотных остатков или нуклеотидов. Различают 2 вида выравнивания: парное (выравнивание двух последовательностей ДНК, РНК или белков) и множественное (выравнивание трех и более последовательностей). Наиболее популярной серией программ для множественного выравнивания последовательностей является Clustal. Первая программа серии Clustal была создана Д.Хиггинсом в 1988 году [8]. Затем она была усовершенствована Д. Фенгом, Р. Дулиттл и В. Тейлором путем добавления прогрессивного выравнивания, то есть созданием множественного выравнивания в результате серий попарных выравниваний, следуя ветвлению направляющего дерева, построенного методом UPGMA [3, 10].

В 1992 году появилась второе поколение программ Clustal. Программа, названная Clustal V, отличалась способностью проводить сопоставления существующих выравниваний и построением направляющего дерева методом NJ [6, 7, 9].

Третье поколение программ, появившееся в 1994 году и названное Clustal W, стало значительно проще в работе благодаря усовершенствованному алгоритму [12]. Кроме этого появилась возможность выбирать матрицы сравнения аминокислот и нуклеотидов, а также устанавливать штрафы за внесение пробелов. Следует отметить, что высокая совместимость программ этого поколения с другими пакетами программ обусловлена за счет предоставления результатов выравнивания в виде формата FASTA. Последним представителем серии является программа Clustal X, для которой характерен более удобный интерфейс и более легкая оценка результатов выравниваний [11]. В настоящее время именно последние программы серии Clustal этого поколения (версия 1.83) позволяют создавать наиболее биологически корректные множественные выравнивания дивергировавших последовательностей [1].

Программы третьего поколения серии Clustal доступны на многих серверах (<http://npsa-pbil.ibcp.fr>, <http://www.ebi.ac.uk>) в двух вариантах – интерактивном и почтовом. Интерактивный вариант предполагает ожидание пользователем получения результатов выравнивания (целесообразно применять при небольшом (<100) количестве последовательностей), а почтовый – по электронной почте (применяется при большом числе последовательностей).

Принципы работы CLUSTAL

Первоначально необходимо ввести на одном из серверов изучаемые аминокислотные или нуклеотидные последовательности в одном из 7 возможных форматов (NBRF/PIR, EMBL/SWISSPROT, Pearson (Fasta), Clustal (\*.aln), GCG/MSF (Pileup), GCG9/RSF, GDE).

Наиболее часто используется формат FASTA, сущность которого заключается во введении знака > перед названием каждой последовательности, а затем (с новой строки) однобуквенном обозначении аминокислот и нуклеотидов. Суммарная длина вводимых последовательностей не должна превышать 40000 для WWW и 60000 для e-mail серверов.

При использовании данной программы выравнивание состоит из трех этапов: парных выравниваний, построения направляющего дерева и множественного выравнивания.

1. В ходе парных выравниваний предварительно сравниваются все возможные пары изучаемых последовательностей. На основании проведенных сравнений вычисляются показатели сходства в соответствии с выбранными матрицами. Существуют 2 разновидности парного выравнивания: медленное (slow) и быстрое (fast). Медленное выравнивание является более точным, но его не рекомендуется применять в случае большого количества (более 20) последовательностей значительной длины (более 1000 остатков). Медленное выравнивание характеризуется 4 параметрами:

- штрафом на внесение делеции (gap open penalty). Уменьшение этого параметра способствует внесению разрывов в выравнивание, что ухудшает качество.

Увеличение – приводит к тому, что выравнивание будет представлять собой длинные участки последовательностей почти без вставок или делеций.

· штраф на продолжение делеции (gap extension penalty). Этот параметр контролирует возможность внесения длинных вставок или делеций.

- матрица сравнений нуклеотидов (DNA weight matrix, Clustal W 1.6). В наиболее широко используемой матрице DNA identity (рис. 1) совпадение нуклеотидов оценивается в 1 балл, а несовпадение – -10000 баллов. Такой высокий штраф за несоответствие облегчает внесение пробелов.

	<b>A</b>	<b>T</b>	<b>G</b>	<b>C</b>
<b>A</b>	<b>1</b>			
<b>T</b>	-10000	<b>1</b>		
<b>G</b>	-10000	-10000	<b>1</b>	
<b>C</b>	-10000	-10000	-10000	<b>1</b>

Рис. 1. Матрица DNA identity.

- матрица сравнения аминокислот (protein weight matrix) – PAM, Blosum и Gonnet.

Выбор матрицы оказывает большое влияние на получаемые результаты, так как каждая матрица представляет отражение отдельных эволюционных гипотез. Известно, что все замены аминокислот не являются равновероятными и в ходе эволюции чаще происходят замены на сходные по физико-химическим свойствам аминокислоты. Так в ходе эволюции гидрофобный изолейцин достаточно часто заменяется на гидрофобный валин и редко на гидрофильный цистеин. Исследования эволюционных изменений различных белковых семейств позволили установить частоты фиксированных мутаций аминокислот и





построения выравнивания выбираются только сегменты, превышающие этот порог. Для увеличения скорости можно уменьшить этот параметр, а для увеличения точности – увеличить.

- длина сегмента, включающего “наилучший выровненный сегмент” (window size). Для увеличения скорости надо уменьшать этот параметр, для увеличения точности – увеличивать.

2. Построение на основании попарных сравнений направляющего дерева (guide-tree). Первоначально методом NJ (neighbor-joining, связывания ближайших соседей) строится бескорневое дерево. Затем устанавливается корень по методу Томсона-Хиггинса-Гибсона таким образом, чтобы значения длин ветвей по отношению к корню остались неизменными.

3. Множественное выравнивание является основой программ Clustal, однако детали его очень сложны. Каждый этап множественного выравнивания состоит из сопоставления двух последовательностей или выравниваний, выполняемого в соответствии с ветвлением дендрограммы. Основными параметрами множественного выравнивания являются:

- штрафы за внесение делеции (gap penalties) устанавливаются как в попарном выравнивании.

- отсрочка различающихся последовательностей (delay divergent sequences) обеспечивает первоочередное выравнивание более сходных последовательностей.

- вес транзиций (transition weight) (А-Г или Ц-Т ) имеет значения между 0 и 1.

Если вес равен 0, то транзигия рассматривается как несовпадение. Если вес равен 1, то транзигия рассматривается как совпадение (алфавит из 4-буквенного вырождается в двухбуквенный пурин-пиримидин). Для слабо сходных последовательностей вес транзиций должен быть близок к 0, для близкородственных – к 1.

- матрицы сравнения нуклеотидов или аминокислот.

Полученное выравнивание может быть отображено в черно-белой или цветной гамме. Идентичные аминокислотные остатки или нуклеотиды отмечаются звездочкой (\*), консервативные замены – двоеточием (:), а полуконсервативные – точкой (.). Консервативность и полуконсервативность аминокислотных замен определяются в соответствии с таблицей 1. Если заменяемые аминокислоты расположены в одной группе, то замена считается консервативной. Результаты выравнивания можно загрузить с помощью приложения Jal View для последующего анализа последовательностей.

Таблица 1

Группы аминокислот, используемые для определения консервативности и полуконсервативности замены при выравнивании последовательностей

Аминокислоты	Цветовое обозначение
AVFPMILW	красный
DE	синий
RHK	сиреневый
STYHCNGQ	зеленый
Другие	серый

Используемые по умолчанию параметры выравниваний

В таблице 2 приведены параметры, используемые по умолчанию на двух серверах (<http://npsa-pbil.ibcp.fr> и <http://www.ebi.ac.uk>).

Таблица 2

Используемые по умолчанию параметры Clustal W

Параметр / сервер		<a href="http://npsa-pbil.ibcp.fr">http://npsa-pbil.ibcp.fr</a>	<a href="http://www.ebi.ac.uk">http://www.ebi.ac.uk</a>
Медленное парное выравнивание	штраф за внесение делеции	15 / 10	15 / 10
	штраф за продолжение делеции	6.66 / 0.1	6.66 / 0.2
	матрица сравнений нуклеотидов / аминокислот	IUB / Gonnet	DNA identity / Gonnet
Быстрое парное выравнивание	размер участка максимального совпадения	2 / 1	.
	штраф на внесение делеции	5 / 3	15 / 10
	число непрерывно совпадающих k-плетов	4 / 5	.
	длина сегмента, включающего "наилучший выровненный сегмент"	4 / 5	.
Множественное выравнивание	штрафы за внесение делеции	15 / 10	15 / 10
	матрицы сравнения нуклеотидов и аминокислот	IUB / Gonnet	DNA identity / Gonnet

Примечание. Первое значение в каждой ячейке относится к последовательностям нуклеиновых кислот, второе – к последовательностям белков.

При выравнивании последовательностей нет необходимости указывать все используемые параметры. Как правило, достаточно указать сервер, на котором проводилось выравнивание, и отметить стандартность его условий.

Применение программ Clustal

Основным предназначением выравниваний, проведенных с помощью программ Clustal, является вычисление на их основании эволюционных дистанций между аминокислотными и нуклеотидными последовательностями, синонимичной и несинонимичной дистанций, определение характера и типа аминокислотных замен и т. д.

В ходе выравнивания также выявляются консервативные участки последовательностей, которые могут являться элементами вторичной структуры, сайтами связывания лигандов и другими функциональными мотивами. Это используется для предсказания вторичной и третичной структуры и функции белков, а также для идентификации новых представителей белковых семейств. Кроме этого, программы Clustal используются для построения дендрограмм, показывающих филогенетические отношения сравниваемых последовательностей без учета (кладограммы) или с учетом длин ветвей (филограммы).

Дот-матрицы

При высоком сходстве последовательностей и небольшом числе делеций выравнивание последовательностей может быть получено с помощью точечных матриц гомологии (дот-матриц). При этом одна из последовательностей располагается горизонтально, а другая – вертикально. Если символы в строке и колонке совпадают, то на их пересечении ставится точка. При сопоставлении двух идентичных последовательностей поставленные точки образуют сплошную линию, при наличии делеций – линия разрывается и ее участки смещены влево или вправо. В случае большой длины последовательностей они разбиваются на подслова длины  $n$ , точка в позиции  $i, j$  ставится только, когда в двух подсловах, начинающихся в позициях  $i$  и  $j$ , совпадает не менее  $k$  символов.

#### Литература

1. Chenna, R., Sugawara, H., Koike, T., Lopez, R., Gibson, T.J., Higgins, D.G., Thompson, J.D. //Nucl. Acids Res. – 2003. – Vol. 31 913). – P. 3497-3500.
2. Dayhoff, M.O., Schwartz, R.M., Orcutt, B.C. // In atlas of protein sequence and structure. – 1978, NBRF, Washington.-Vol. 5, suppl. 3 (Dayhoff, M.O., ed.). – P. 345-352.
3. Feng, D.F., Doolittle, R.F. //J. Mol. Evol. – 1987. – Vol. 25. – P. 351-360.
4. Gonnet, G.H., Cohen, M.A., Benner, S.A. // Science. – 1992. – Vol. 256. – P. 1443-1445.
5. Henikoff, S., Henikoff, J.G. //Proc. Natl. Acad. Sci. – 1992. – P. 10915-10919.
6. Higgins, D.G. //Methods Mol. Biol. – 1994. – Vol. 25. – P. 307-318.
7. Higgins, D.G., Bleasby, A.J., Fucks, R. //Comput. Appl. Boisci. – 1992. – Vol. 8. – P. 189-191.
8. Higgins, D.G., Sharp, P.M. //Gene. – 1988. – Vol. 73. – P. 237-244.
9. Saitou, N., Nei, M. //Mol. Biol. Evol. – 1987. – Vol. 4. – P. 406-425.
10. Taylor, W.R. //J. Mol. Evol. – 1988. – Vol. 28. – P. 161-169.
11. Tompson, J.D., Gibson, T.J., Plewniak, F., Jeanmougin, F., Higgins, D.G. //Nucl. Acids Res. – 1997. – Vol. 25. – P. 4876-48882.
12. Tompson, J.D., Higgins, D.G., Gibson, T.J. //Nucl. Acids Res. – 1994. – Vol.22. – P.4673 – 4680.