

АНАЛИЗ ГЕНЕТИЧЕСКИХ ТЕКСТОВ. I. 1-ГРАММНЫЕ ХАРАКТЕРИСТИКИ

В.Д.Гусев, В.А.Куличков, Т.Н.Титкова

Назовем текстом конечную упорядоченную символьную последовательность, составленную из элементов конечного алфавита. Пусть n - длина текста (в символах), $n = |\Lambda|$ - число элементов (мощность) алфавита Λ . Подпоследовательность текста, состоящую из 1 расположенных подряд символов, будем называть 1-граммой ($1 = 1, 2, \dots, n$), а анализ текста, основанный на подсчете частот встречаемости в нем различных 1-грамм, - 1-граммным анализом.

Под генетическим текстом будем понимать представление молекул нерегулярных полимеров (ДНК, РНК, белков) в виде упорядоченной последовательности мономеров (нуклеотидов или аминокислот). Важность исследования данного класса молекул обусловлена тем, что они определяют протекание наиболее существенных генетических процессов в организмах. Алфавит нуклеотидов состоит из 4 символов (А, Т, Г, Ц для ДНК и А, У, Г, Ц для РНК), алфавит аминокислот - из 20 символов.

В последние годы благодаря совершенствованию техники определения нуклеотидного состава молекул были расшифрованы и опубликованы первые полные нуклеотидные последовательности (геномы) простейших микроорганизмов (4X174, G4, SV40, PBR322, fd, M132 и ряда других). Длины этих геномов составляют порядка ($10^3 - 10^4$) символов. Весьма актуальной задачей становится создание и пополнение банка данных по микробиологическим объектам и пакетов прикладных программ по исследованию различных свойств этих объектов.

Цель данной работы - исследование 1-граммных характеристик текстов, кодирующих наследственную информацию простейших микроорганизмов. Указан ряд содержательных задач, для решения которых представляется перспективным использовать 1-граммное описание текста. Ряд примеров иллюстрирует возможность выявления по 1-грам-

мным характеристикам функционально значимых фрагментов (или особенностей) текста.

Отметим, что наряду с информацией о частотах встречаемости различных 1-грамм в тексте очень важна и информация о местах вхождения каждой конкретной 1-граммы в текст. Эта информация зачастую дает возможность отличать случайные совпадения 1-грамм от неслучайных, функционально значимых.

I. Структура геномов вирусных частиц [1]. Из шести исследованыхся геномов (см. выше) пять первых представляют собой кольцевые молекулы ДНК. Поскольку двухцепочечные молекулы ДНК построены по принципу комплементарности (т.е. любая цепь однозначно получается из другой повсеместной заменой нуклеотидов А, Т, Г, Ц нуклеотидами Т, А, Ц, Г соответственно), анализу подвергалась лишь одна из них. В тех случаях, когда это было возможно, анализировались цепи, соответствующие структуре м-RНК ($\phi X174$, G4, fd, MS2). Шестой геном (MS2) описывает одноцепочечную молекулу РНК.

Структура геномов фагов и вирусов определяется основными генетическими процессами, обеспечивающими воспроизведение этих организмов. К таким процессам относятся редупликация (копирование ДНК или РНК), транскрипция (синтез РНК на ДНК) и трансляция (синтез белка на РНК). Части текста, отвечающие отдельным актам редупликации, транскрипции и трансляции, называются соответственно репликоном, скриптоном и цистроном. Каждая из этих функциональных единиц имеет свои знаки начала и конца - инициаторы и терминаторы, представляющие собой нуклеотидные последовательности длиной от трех до нескольких десятков символов.

Между перечисленными функциональными единицами существует, как правило, отношение иерархического вложения: геном - это совокупность репликонов, репликон - последовательность скриптонов, скриптон, в свою очередь, разбивается на цистроны, каждый из которых уже кодирует определенный белок. Текст, отвечающий цистрону, разбивается при этом на последовательность триплетов (кодонов), кодирующих в соответствии с генетическим кодом определенные аминокислоты. Отдельные цистроны могут располагаться последовательно, не пересекаясь друг с другом, могут частично пересекаться, могут быть даже целиком вложенными один в другой. Различие в синтезируемых ими белках достигается в последнем случае за счет использования отличающихся по фазе триплетных рамок считывания, либо (при совпадающих фазах) за счет несовпадения длин кодирующих последовательностей.

Приведенное выше описание структуры генома весьма схематично, но достаточно для наших целей. Следует отметить лишь, что наряду с указанными функциональными единицами, назначение которых в той или иной степени выяснено, существуют некоторые другие, функциональное назначение которых пока непонятно (нетранскрибуемые участки ДНК, межстронные интервалы, интроны). Исследование различных характеристик этих участков (в том числе и 1-граммных) и сопоставление их с аналогичными характеристиками известных функциональных единиц может дать толчок к выявлению их функционального назначения.

2. Примеры содержательных задач. Укажем несколько задач из области анализа нуклеотидных последовательностей, при решении которых существенную пользу может принести знание 1-граммных характеристик этих последовательностей.

Следует сразу заметить, что 1-граммные характеристики расшифровываемых нуклеотидных текстов (или отдельных участков этих текстов) всегда интересовали исследователей. Однако интерес этот в подавляющем большинстве случаев был односторонним. А именно анализировались лишь участки геномов, соответствующие цистронам (т.е. участки, кодирующие белки), причем анализ последних ограничивался преимущественно подсчетом частот использования различных кодонов (т.е. триграмм) и исследованием распределения нуклеотидов по различным позициям кодонов.

Интерес к характеристикам третьего порядка обусловлен, естественно, тривиальной структурой генетического кода. Нет никаких оснований предполагать, что они будут играть основную роль и при исследовании некодирующих частей геномов. Ниже будет показано, что важная информация содержится не только в характеристиках третьего порядка, но и в характеристиках меньшего ($l = 1, 2$) и большего ($l = 4, 5, \dots, l_{\max}$) порядков. Параметр l_{\max} здесь соответствует тому минимальному значению l , начиная с которого в тексте уже отсутствуют повторяющиеся 1-граммы, т.е. для любой l -граммы текста x частота ее встречаемости $F_l(x) \leq 1$ при $l_{\max} \leq l \leq n$.

2.1. Задача разметки текста сводится к построению алгоритмов автоматического выделения функциональных единиц геномов (см. п.1) и соответствующих им знаков пунктуации (инициаторов и терминаторов). Эта задача разбивается на ряд отдельных подзадач, многие из которых в формальной постановке сводятся к задаче распознавания.

Укажем для примера на задачу автоматического выделения промоторов – знаков пунктуации, ответственных за начало процесса транскрипции. Каждый промотор представляет собой последовательность длиной до 50 символов. В настящее время расшифровано уже несколько десятков промоторов, т.е. установлена их первичная структура и доказано (молекулярно-биологическим анализом), что соответствующие последовательности нуклеотидов выполняют функции инициаторов транскрипции. Совокупность этих последовательностей может рассматриваться как обучающая выборка по промоторам в задаче классификации "промотор-не промотор". Для элементов этой выборки существует реальное опознание устройство – РНК-полимераза. Обучающая выборка по "не промоторам" может быть получена из текстов уже расшифрованных геномов с исключенными промоторными зонами.

Целью обучения является построение решающего правила, моделирующего принцип действия РНК-полимеразы, т.е. позволяющего относить произвольную последовательность к одному из двух указанных классов. Результат классификации может быть подтвержден (либо опровергнут) весьма трудоемким молекулярно-биологическим анализом. Удачное решение задачи классификации формальными методами (т.е. прямо по виду символьной последовательности) может существенно упростить и ускорить процедуру выявления генетической структуры генома.

В основу построения решающего правила кладется различие в характеристиках элементов первого и второго образов, выявляемое при анализе обучающих выборок. Нами исследовалась возможность использования для этой цели 1-граммных характеристик. А именно была выдвинута и экспериментально подтверждена гипотеза в том, что различие между промоторами и "не промоторами" наиболее ярко проявляется на уровне длинных 1-грамм ($1 \geq 7$). Было показано, что существуют достаточно длинные 1-граммы, типичные только для промоторов, т.е. 1-граммы, каждая из которых встречается как минимум в двух промоторах, но не встречается у "не промоторов".

Примерами таких 1-грамм являются: ТГТТГАЦА (в промоторах λ-PL, E.COLI-TRP, ΦX174D, S.TURN-TRP), АЦАЦТТ (в промоторах E.COLI-LAC, E.COLI-GAL, TRHК-TTRP, E.COLIK12-ARAC), ГЦГГГГАТА (в промоторах λ-PL, λ-PR, λ-PRM) и ряд других. Анализ позиций, занимаемых этими 1-граммами внутри промоторов, подтверждает, что совпадение 1-грамм у разных промоторов носит неслучайный характер. Каждой из указанных 1-грамм внутри разных промоторов соот-

всегда будут одинаковые или очень близкие позиции. Таким образом, сочетание 1-граммного и позиционного анализа позволяет в данном случае выявить функциональные зоны внутри промоторов и установить возможные варианты нуклеотидных последовательностей для этих зон.

Выделенное множество характерных (в упомянутом выше смысле) 1-грамм естественно назвать "покрытием" обучающей выборки промоторов в случае, если каждый прометер из этой выборки содержит хотя бы одну характерную 1-граммму. В качестве примера укажем, что для обучающей выборки из 24 промоторов удалось получить покрытие из 10 характерных 1-грамм ($6 \leq 1 \leq 10$). Элементы покрытия в сочетании с информацией о наиболее вероятных зонах их расположения внутри промоторов представляют основу для выработки различных стратегий распознавания.

2.2. Задача вычисления эволюционного расстояния между геномами является составным элементом более общей задачи построения филогенетических древ для различных семейств геномов. Можно выделить следующие элементарные преобразования, посредством которых осуществляется эволюция геномов:

- 1) замена одного основания другим;
- 2) вставка цепочки нуклеотидов в исходную последовательность;
- 3) устранение цепочки нуклеотидов из исходной последовательности;
- 4) инверсия (замена цепочки нуклеотидов в исходной последовательности элементами комплементарной цепочки, прочитанными в обратном порядке). Например, цепочка ААТЦГТЦ исходной последовательности будет заменена инвертированной ГАЦЦГАТТ, которая получается из комплементарной к исходной (ТТАГЦЦАГ) прочтением ее справа налево;
- 5) дупликация (повторение цепочки символов);
- 6) транспозиция (перенес цепочки символов из одного места последовательности в другое).

Естественно определить меру различия между двумя геномами как минимальное число осмысленных с биологической точки зрения преобразований типа I-6, переводящих один геном в другой. Вычисление этой меры представляет нетривиальную комбинаторную задачу, пока нерешенную в общем виде. Для частного случая, когда допустимыми являются лишь преобразования I-3, причем два последних при-

менимы только к цепочкам длины 1, расстояние вычисляется с помощью метода динамического программирования. При этом существует однозначная связь между расстоянием и линейкой максимальной длины общей подследовательности двух последовательностей (геномов) [2].

Отметим, что преобразования 2-6 несёт "групповой" характер, т.е. применимы к цепочкам, состоящим из многих (иногда до нескольких десятков) символов. Для вычисления расстояния необходимо уметь отыскивать в текстах участки, соответствующие вставкам (устранениям), инверсиям, дупликациям и транспозициям. Основу для этого дает вычисление и сопоставление полных частотных спектров обоих текстов.

Так, например, различие в частотах встречаемости какой-либо достаточно длинной^{*)} 1-грамм в двух родственных текстах, как правило, свидетельствует о том, что в процессе эволюции имела место дупликация данной 1-граммы. Сопоставление результатов частотного и позиционного анализа зачастую позволяет выявить дупликацию даже при наличии одного текста. Об этом сигнализирует близость (а при многократной дупликации – периодичность) расположения повторяющихся 1-грамм.

Для выявления инверсий необходимо сравнить частотные спектры двух текстов: любого из исходных и инвертированного оставшегося. Наличие достаточно длинных (например, по отношению к схеме независимого порождения) совпадающих 1-грамм в обоих текстах говорит о возможности инверсий. Если в результате позиционного анализа выявляется, что номера позиций, занимаемых интересующей нас 1-граммой, в обоих исходных текстах совпадают или близки друг к другу, то предположение о наличии инверсии подтверждается, в противном случае возможна транспозиция с инверсией, либо случайное совпадение.

Соображения, аналогичные вышеприведенным, могут быть положены и в основу отыскания транспозиций. О наличии вставок (устранений) сигнализирует различие в длинах анализируемых текстов. Для локализации вставок целесообразно осуществлять предварительную синхронизацию обоих текстов по выявленным характерным участкам, в качестве которых могут выступать знаки пунктуации, совпадение 1-грамм

*) При малых значениях λ частоты одинаковых 1-грамм в родственных геномах будут отличаться почти всегда, поскольку степень изменения частотных характеристик под влиянием любого из преобразований 1-6 обычно тем больше, чем меньше λ .

мы, инверсии, дупликации. Различие длин последовательностей, за ключенных между идентичными знаками синхронизации в обоих геномах, свидетельствует о наличии вставки (утраления).

2.3. Восстановление текста по фрагментам. Определение первичной структуры длинных последовательностей (например, целых геномов) в настоящее время технически трудно осуществимо. С помощью существующих методик удается восстанавливать первичную структуру лишь относительно коротких участков (~ 20–300 нуклеотидов). Поэтому для определения первичной структуры всего генома его предварительно разрезают на более короткие участки (фрагменты) специальными ферментами-рестриктазами, определяют первичную структуру каждого участка, а затем, соединяя нужным образом расшифрованные куски, восстанавливают последовательность, соответствующую геному в целом.

Определение порядка следования фрагментов является нетривиальной задачей. Для восстановления текста используют наборы фрагментов, соответствующие разным рестриктазам. Поскольку каждая рестриктаза разрезает геном лишь в характерных для нее участках (например, рестриктаза НaeII реагирует только на сочетание ГГЦЦ, разрывая связи между Г и Ц), фрагменты, получаемые от разных рестриктаз, оказываются перекрывающимися. Это свойство и используется для упорядочения фрагментов.

С формальной точки зрения задача восстановления текста по набору фрагментов, полностью его покрывающих, сводится к исследованию вопроса о существовании единственного (с учетом перекрываемости) варианта упорядочения фрагментов, полученных при фиксации определенного набора стратегий расщепления текста (набора рестриктаз). Предложено несколько способов сведения данной задачи к задаче отыскания эйлеровых цепей на графе. Структура соответствующего графа определяется типом текста и используемым набором рестриктаз.

Поскольку с каждой рестриктазой ассоциируется определенная 1-грамма, то количество таких 1-грамм в тексте и их расположение будут определять соответствующий набор фрагментов. Зная 1-граммные и позиционные характеристики родственных текстов, можно целенаправленно формировать набор рестриктаз, так чтобы добавление каждой новой рестриктазы в максимальной степени уменьшало неоднозначность, возникающую при восстановлении текста по уже имеющемуся набору фрагментов.

3. Частотные характеристики генетических текстов. В табл.1 приведены частотные характеристики текстов ФХ174, Г4, РД, СВ40, РВР322 и МС2 для значений $l = 1, 2, 3$. Все 1-граммы упорядочены по убыванию частоты встречаемости их в тексте. Порядковый номер g каждой 1-граммы будем называть ее рангом. Параметр g изменяется в диапазоне от 1 до M_1 , где M_1 - количество различных 1-грамм в тексте.

Поскольку с увеличением l объем таблиц возрастает очень быстро (при малых l , как правило, $M_1 = n^l$; при $l \rightarrow l_{\max}$ значение M_1 ограничено величиной $N-l+1$), для значений $l > 3$ выборочно приведены лишь интегральные частотные характеристики E_1^k и S_n (табл.2). Здесь E_1^k - количество различных 1-грамм текста, каждая из которых встречается ровно k раз; $S_n = \sum_{i=1}^n k(i) \cdot E_1^{(i)}$ - текущая сумма числа 1-грамм, соответствующих первым частотным градациям, упорядоченным по убыванию.

Наиболее длинные повторяющиеся 1-граммы, встретившиеся в каждом из текстов, приведены в табл. 3.

Табл. 4 содержит энтропийные характеристики текстов [4] (значения \hat{H}_1 и \hat{R}_1 для $l = 2, \dots, 6$). Ввиду ограниченности длины текстов энтропийные оценки с увеличением l становятся статистически недостоверными. Формально это выражается в том, что с ростом l величины $\hat{H}_1 \rightarrow 0$, а $\hat{R}_1 \rightarrow 1$. Ввиду отсутствия эффективных подходов к определению порогового значения $l(N)$, определяющего статистически достоверные результаты от недостоверных, мы можем воспользоваться для этой цели лишь некоторыми косвенными соображениями. Не обсуждая их детально, укажем, что в данном случае можно, по-видимому, в качестве оценки для $l(N)$ взять значение $\hat{l}(N) = 4$.

В табл.6 приведены величины коэффициентов ранговой корреляции Спиримана между каждой парой текстов для значений $l = 2, 3, 4$. При больших значениях l уже не все из 4^l возможных комбинаций нуклеотидов присутствуют в тексте, поэтому использовать данную меру сходства без соответствующей модификации нецелесообразно. Нуклеотиды Т (в ДНК-содержащих геномах) и У (в РНК-содержащих) при вычислении ранговых мер сходства отсутствуют.

4. Сравнительный анализ 1-граммных характеристик генетических текстов. Ниже обсуждаются наиболее интересные, на наш взгляд, закономерности, выявленные при 1-граммном анализе текстов. Далеко

не для всех из них удается предложить соответствующую интерпретацию. И даже там, где это делается, интерпретация, как правило, носит характер гипотезы. Укажем вначале на два интересных свойства, вытекающих из анализа характеристик первого и второго порядка.

4.1. Классификация по Т, А- или Ц, Г-преблажанию. Сопоставление характеристик первого порядка показывает, что все тексты можно разделить на две группы в зависимости от ранжировки нуклеотидов по частоте встречаемости. В первую группу входят тексты ФХ174, G4, FD и SV40, в которых преобладают нуклеотиды Т и А (Т, А - богатые тексты), во вторую - PBR322 и MS2 (Ц, Г - богатые). Промежуточные варианты, когда Т и А перемежаются с Ц и Г при частотном упорядочении, отсутствуют.

Данная классификация, по-видимому, характеризует вторичную структуру молекул РНК, кодируемых этими геномами, т.е. их пространственную конфигурацию. Известно, что бактериофаг MS2 в отличие от объектов первой группы имеет сложную вторичную структуру. Ее отличительным признаком является наличие большого количества длинных "шпилек" - участков самокомplementарности одноцепочечной молекулы MS2. В [3] показано, что шпильки формируются преимущественно из кодонов, содержащих нуклеотиды Ц и Г в третьей позиции. Поскольку комплементарная пара Ц = Г обладает тремя водородными связями, а пара А = У лишь двумя, преимущественное использование Ц и Г в третьей позиции (а только здесь и допустимы вариации) должно обеспечить большую устойчивость шпилек.

Возможность подтверждения либо опровержения предложенной интерпретации содержится в анализе пока еще не установленной вторичной структуры РНК, транскрибируемых с генома плазмида PBR322. Он существенно отличается от генома MS2 тем, что является ДНК-содержащим и двухцепочечным, аналогично объектам из первой группы. Если его РНК обладают слабо выраженной вторичной структурой, то гипотеза о связи ранжировки нуклеотидов в частотной характеристике первого порядка со сложностью и стабильностью вторичной структуры соответствующих молекул отпадает. В противном случае эта гипотеза получает сильное подтверждение.

4.2. ЦГ-эффект у SV40. Анализ частотных характеристик второго порядка (табл. I) обнаруживает очень сильный аномальный эффект в частоте встречаемости биграммы ЦГ у вируса SV40. Эта биграмма встречается в тексте всего лишь 27 раз ($\tau = 16$). Для

1 - б о с л и ц

Использование характеристики гомотетических генетических генотипов для определения
1 = 1,2,3

1	2	Ф274		G4		FD		ST40		PBR322		MS2	
		1-граам	F ₁ (x)										
1	1	I	1684	A	1519	I	2210	I	1582	I	1208	I	932
2	2	A	1291	I	1510	A	1577	A	1518	I	1134	I	928
3	3	I	1254	I	1446	I	1321	I	1094	I	1036	I	875
4	4	I	1157	I	1102	I	1299	I	1033	A	984	A	834
5	1	II	572	AA	554	II	815	II	575	II	578	II	266
6	2	II	480	II	503	AA	544	AA	533	II	329	II	263
7	3	II	404	II	443	AI	515	II	422	II	306	II	250
8	4	AA	395	AT	392	II	483	II	398	IA	300	IP	242
9	5	AT	383	II	391	II	481	II	389	III	292	II	231
10	6	II	347	AI	386	II	469	AT	352	II	287	II	230
11	7	II	327	II	352	II	445	AT	345	II	285	III	223
12	8	II	325	II	342	II	397	II	326	II	281	AA	222
13	9	II	320	II	334	II	322	II	292	II	259	II	217
14	10	II	312	II	313	II	316	AI	288	AA	258	AI	215
15	11	II	267	III	304	II	286	II	272	AT	254	II	211
16	12	AI	260	II	287	II	278	II	266	AT	239	II	207
17	13	II	257	II	279	II	274	II	257	II	236	AT	204
18	14	II	255	II	272	AI	274	II	247	II	235	II	201

16	AP	15	AP	20	TTT	17	ATT	12	TTT	11	TTT	10	TTT	9	TTT	8	TTT	7	TTT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT		
15	AP	16	AP	21	TTT	18	ATT	13	TTT	12	ATT	11	TTT	10	TTT	9	TTT	8	TTT	7	TTT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT
14	AP	15	AP	20	TTT	16	ATT	11	TTT	10	ATT	9	TTT	8	TTT	7	TTT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT				
13	AP	14	AP	19	TTT	15	ATT	10	TTT	9	ATT	8	TTT	7	TTT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT						
12	AP	13	AP	18	TTT	14	ATT	9	TTT	8	ATT	7	TTT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT								
11	AP	12	AP	17	TTT	13	ATT	8	TTT	7	ATT	6	TTT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT										
10	AP	11	AP	16	TTT	12	ATT	7	TTT	6	ATT	5	TTT	4	TTT	3	TTT	2	TTT	1	TTT												
9	AP	10	AP	15	TTT	11	ATT	6	TTT	5	ATT	4	TTT	3	TTT	2	TTT	1	TTT														
8	AP	9	AP	14	TTT	10	ATT	5	TTT	4	ATT	3	TTT	2	TTT	1	TTT																
7	AP	8	AP	13	TTT	9	ATT	4	TTT	3	ATT	2	TTT	1	TTT																		
6	AP	7	AP	12	TTT	8	ATT	3	TTT	2	ATT	1	TTT																				
5	AP	6	AP	11	TTT	7	ATT	2	TTT	1	ATT		TTT																				
4	AP	5	AP	10	TTT	6	ATT		TTT		ATT		TTT																				
3	AP	4	AP	9	TTT	5	ATT		TTT		ATT		TTT																				
2	AP	3	AP	8	TTT	4	ATT		TTT		ATT		TTT																				
1	AP	2	AP	7	TTT	3	ATT		TTT		ATT		TTT																				
				196	TTT	197	ATT		TTT		ATT		TTT																				
				244	TTT	223	ATT		TTT		ATT		TTT																				
				264	TTT	273	ATT		TTT		ATT		TTT																				
				277	TTT	258	ATT		TTT		ATT		TTT																				
				279	TTT	247	ATT		TTT		ATT		TTT																				
				282	TTT	232	ATT		TTT		ATT		TTT																				
				190	TTT	214	ATT		TTT		ATT		TTT																				
				192	TTT	191	ATT		TTT		ATT		TTT																				

I Podobno tak

T 86 π π 2

ОБЩЕСТВЕННО-ПРАВОВЫЙ АСПЕКТ ПРОЦЕССА РЕГИСТРАЦИИ

сравнения укажем, что частоты встречаемости биграмм с рангом I6 во всех остальных текстах в 7-9 раз выше.

Отметим, что исследователи, занимавшиеся изучением кодирующих частей генома SV40 и родственных вирусов, уже обращали внимание на чрезвычайно малое количество кодонов, содержащих биграмму ЦГ в позициях 1-2 или 2-3. Наш анализ показывает, что данный эффект возникает уже на уровне биграмм, т.е. аномалии в использовании кодонов – его следствие, а не причина. Эффект наблюдается для генома в целом, как для кодирующих частей (при любых фазах считывания), так и для некодирующих. Выявление данного и ему подобных эффектов не требует предварительной разметки текстов и может быть автоматизировано с помощью достаточно развитой в настоящее время техники обнаружения аномальных наблюдений.

Мы не можем пока указать, какую функциональную нагрузку несет данный эффект. В [6] предпринята попытка объяснения данного эффекта тем, что в геномах вирусов эукариотов существует механизм перехода ЦГ в ТГ (5-метил-цитозин часто мутирует в Т, причем Ц метилируется именно в составе биграммы ЦГ). Этот процесс должен был бы приводить к соответствующему избытку ТГ. Проверка данной гипотезы на схеме независимого порождения (см.п.4.5) показала, однако, что некоторый избыток ТГ имеет место, но он далеко не компенсирует недостатка ЦГ.

Отметим, что из шести анализировавшихся микроорганизмов SV40 является единственным, который паразитирует на эукариотах, остальные – на прокариотах. Если данный эффект подтвердится и на других вирусах, паразитирующих на эукариотах, он может быть положен в основу алгоритма классификации текстов по признаку: эукариотический-прокариотический.

4.3. Наиболее длинные повторения. Значения параметра l_{\max} составляют: для текста SV40 – 59, РД – 33, G4 и РВВ322 – 14, ФХ194 – 13, М52 – 11. При $l = l_{\max} - 1$ в каждом из текстов уже появляются повторяющиеся 1-граммы (см. табл.3). Анализ их расположения в текстах проливает некоторый свет на их происхождение.

Самый длинный повтор ($l = 58$) наблюдается в тексте SV40 (1-грамма №14 из табл.3). Эта 1-грамма начинается и заканчивается триграммой ТГГ. В тексте обе 58-граммы расположены подряд так, что последние три символа первой 1-граммы накладываются на первые три символа второй. Таким образом, наличие данного повтора можно

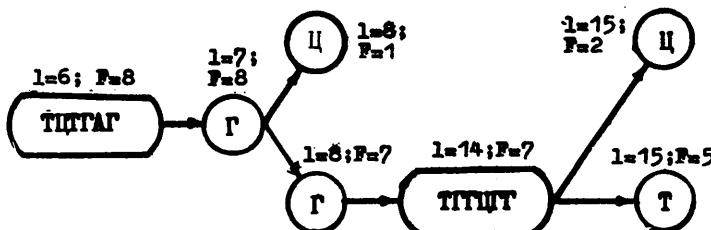
Т а б л и ц а 3

Наиболее длинные повторяющиеся 1-граммы
генетических текстов

Текст	l	F_1	1-граммы	№ 1-граммы
ФХ174	I2	2	ЦГЦЦААГТАЦТТ	1
	I2	2	ЦЦЦДГГЦЦГГТТ	2
	II	2	ТТЦЦГГГГАТТ	3
	II	2	ЦЦЦЦГГГААГ	4
G4	I3	2	ГГГГГГЦГГАТТЦ	5
	II	2	ГГГГГГААТТЦ	6
	II	2	ЦГАЦГГГГГГ	7
	II	2	АЦЦЦЦТГААГ	8
FD	32	2	ГГГГЦГГГГГГГГГГГГГГГГГГГГГГГГГГГГГГ	9
	32	2	ТГГ	10
	20	2	ГГ	11
	I3	2	ААГГГГААТТЦААА	12
	I2	2	ЦЦЦААТТЦГГ	13
SV40	58	2	ТГГГГГГГГАААТТГГГГААГГГГААГГГГГГГГГГГГГ	14
	21	2	ТГГГГГГГГААГГГГГГГГГГГГГГ	15
	I3	2	ГАГГАЦАЦАГАГГ	16
	I3	2	ААААЦЦАГААГ	17
	II	2	ЦЦЦАААААА	18
РВР322	I3	2	АГЦААААГГЦАГ	19
	II	2	ТГГГГГГГГГ	20
	II	2	ЦГГГГГГГГ	21
	II	2	АААААГГААГ	22
MS2	I0	2	АГАЦГГГГГ	23
	I0	2	АУУЦЦЦУЦАГ	24
	8	4	АГГУГГЦУ	25
	8	3	ГУУУУАЦА	26

было бы объяснить дупликацией соответствующей 55-грамм без наложения (но при этом маловероятно появление еще одного ТТГ после такой композиции) либо дупликацией всей 58-граммы с наложением (одним из возможных объяснений механизма наложения может быть "неравный" кроссинговер). Наличие данного повтора уже отмечалось. Функциональное назначение его пока не выяснено, известно лишь, что в этом районе инициируется репликация генома SV40. С интервалом в три символа от данного участка находится еще одна дупликация (дважды повторяется 21-гамма №15 из табл. 3).

Интересный пример дупликаций наблюдается в тексте № 1. Ядро этой дупликации составляет 14-гамма X = ТЦТГАГГГГГЦТГ. Она легко выделяется при анализе дерева всевозможных продолжений 6-грамм ТЦТГАГ:



Из приведенного рисунка видно, что при изменении l от 8 до 13 дерево не имеет разветвлений, т.е. продолжается формирование какого-то сообщения (с неясной пока семантикой). При $l = 14$ формирование заканчивается, и далее дерево расщепляется. Три 15-граммы, заканчивающиеся на Т, следуют в тексте друг за другом, образуя тройную дупликацию, которой предшествует триграмма ГГЦ(ГГЦХТХТХ). На расстоянии порядка 400 символов от этого участка расположен другой, образованный четырехкратным (с интервалом в 1 символ) повторением базовой 1-граммы и заканчивающийся триграммой ТГЦ(ХЦХТХЦХТЦ). Обладая внутренней периодичностью, два этих участка содержат в себе более длинные (по сравнению с базовой) повторяющиеся последовательности (1-граммы № 9, 10 и II из табл. 3).

Интересно отметить, что расположение базовых 1-грамм синхронизовано относительно рамки считывания при трансляции так, что каждая из них кодирует идентичные фрагменты белка. Это указывает на то, что в данном случае, в отличие от предыдущего, семантическую нагрузку дупликаций следует искать на уровне функционирования кодируемого ими полипептида.

Примером дупликаций являются такие 1-граммы №6 (самый длинный повтор в G4) и №13 (PBR322). Интересно отметить, что последняя 1-грамма, аналогично 58-граммме из SV40, имеет конец, совпадающий с началом (биграмма АГ). Механизм дупликации вновь улавливает это и осуществляет дупликацию с наложением. Более короткие 1-граммы, например №1 (из ФХ174), №23 (из MS2), встречающиеся каждая по два раза, уже разнесены внутри своих текстов, т.е. наличие таких повторов может трактоваться и как случайный фактор.

Т а б л и ц а 4
Оценки энтропии и избыточности
генетических текстов

Текст	под- стр. 1-гра. ммы	1				
		2	3	4	5	6
ФХ174	\hat{H}_1	1,968	1,932	1,900	1,798	1,453
	\hat{R}_1	0,015	0,033	0,049	0,100	0,273
G4	\hat{H}_1	1,963	1,937	1,903	1,815	1,479
	\hat{R}_1	0,018	0,031	0,048	0,092	0,260
FD	\hat{H}_1	1,963	1,928	1,890	1,798	1,493
	\hat{R}_1	0,018	0,036	0,055	0,101	0,254
SV40	\hat{H}_1	1,906	1,980	1,862	1,810	1,451
	\hat{R}_1	0,047	0,055	0,069	0,095	0,274
PBR322	\hat{H}_1	1,989	1,974	1,944	1,855	1,389
	\hat{R}_1	0,006	0,013	0,028	0,073	0,305
MS2	\hat{H}_1	1,997	1,985	1,956	1,823	1,284
	\hat{R}_1	0,001	0,007	0,021	0,088	0,357

4.4. Низкая избыточность. С учетом замечания о достоверности энтропийных оценок, сделанного в п.3, избыточность генетических текстов, как это следует из табл.4, очень невелика (не превышает 0,1). Для сравнения укажем, что оценки избыточности современных естественных языков (русского, английского, немецкого) дают значение $R \approx 0,7$.

Отсутствие избыточности в языке означает, что любая комбинация символов в нем является осмысленной. Именно по такому принципу построен генетический код, использующий все 64 возможные комбинации троек нуклеотидов. Анализ некодирующих частей показывает, что и там на уровне $l = 3$ нет запрещенных комбинаций.

В порядке нарастания избыточности тексты ранжируются следующим образом: MS2, PBR322, G4 и ФХ174, РД, SV40 . Приведенная ранжировка, возможно, отражает многообразие и степень перекрывания различных функциональных единиц в геномах (в некотором смысле сложность организации геномов). У генома MS2 , избыточность которого минимальна, гены практически не перекрываются друг с другом. В геномах G4, ФХ174, SV40 имеются зоны перекрытия, иногда весьма обширные. Наложение различных функциональных единиц (гена на ген, промотора на ген и т.д.) увеличивает количество ограничений, которым должна удовлетворять символическая последовательность, т.е. повышает избыточность языка. Аналогии подобного рода проводил К.Шеннон [4], обсуждая связь между избыточностью языка и возможностью построения кроссвордов.

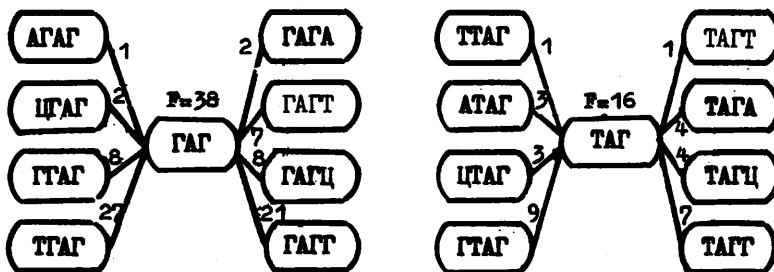
Низкая избыточность находит свое отражение в слабом наклоне частотных характеристик соответствующего порядка. Наклон характеристики 1-го порядка может быть определен как МИК – оценка параметра γ_1 в выражении $F_1(r) = C_1 r^{\gamma_1}$, используемом для аппроксимации частотной кривой. Значения оценок этого параметра довольно сильно отличаются для текстов двух выделенных групп, но не превышают 0,5 ($l = 2,3$). Это говорит о том, что закон Ципфа (близость параметра γ к единице для частотных словарей естественного языка) не выполняется для 1-граммных характеристик генетических текстов. Более того, различия в значениях этого параметра могут быть положены в основу классификации генетических текстов.

4.5. Схема независимого порождения как грубая модель текста. Низкие значения избыточности R для $l = 2,3,4$ наводят на мысль о том, что, по крайней мере, в данном диапазоне значений l можно воспользоваться в качестве грубой модели текста схемой независимых испытаний с вероятностями порождения каждого символа « равными $P(\alpha)/N$, где $\alpha \in \{A,T,G,C\}$, а $P(\alpha)$ – частота встречаемости символа α в соответствующем тексте. Основанием для этого служит тот факт, что для последовательности, полученной в соответствии со схемой независимых испытаний, $H_1 = H_2 = \dots = H_l$ для любого l . Анализ табл.4 показывает, что для $l = 2,3,4$ мы делаем не слишком большую натяж-

ку, допуская эту гипотезу. В наилучшем соответствии с ней находятся тексты с наименьшей избыточностью (MS2, PBR322). Это подтверждается и прямым сопоставлением частот 1-грамм, наблюдаемых фактически и вычисляемых по схеме независимых испытаний.

Схема независимого порождения представляет интерес с точки зрения оценки многих интересующих нас параметров, например, таких как $\max M_1, E_1^1$, определяющих трудоемкость алгоритмов обработки символьных последовательностей. С биологической же точки зрения несомненный интерес представляет анализ очень сильных отклонений фактических частот от модельных. Зачастую такие отклонения являются функционально значимыми.

Рассмотренный выше Ц-эффект у SV40 можно трактовать, с одной стороны, как аномалию в частотной характеристике, с другой стороны, – как сильное отклонение от схемы независимых испытаний ($F_{\text{мод}}(\text{ЦГ}) = 217$, $F_{\text{эксп}}(\text{ЦГ}) = 27$). Пример дупликации в тексте FD, описанный в п. 4.3, можно также трактовать как аномальное отклонение от схемы независимых испытаний, поскольку при увеличении 1 от 8 до 14 частота соответствующей 1-граммы не должна была оставаться постоянной, а должна была уменьшаться на каждом шаге за счет разветвления дерева на 4 ветви с вероятностями $P(\alpha)$ ($\alpha \in \{A, T, G, C\}$) для каждой ветви. Другой пример аномального ветвления – лево- и правосторонние расширения 1-грамм ГАГ(из FD) и ТАГ(из G4):



В табл. 5 приведен еще ряд примеров аномальных отклонений экспериментально наблюдаемых частот ($F_{\text{эксп}}$) от вычисляемых в соответствии с моделью независимого порождения ($F_{\text{мод}}$). Отметим, что во всех текстах фактически наблюдаемые частоты встречаемости 1-грамм АА, ААА выше, а триграммы ТАГ (терминальный кодон) существенно ниже, чем это следует из схемы независимых испытаний.

Таблица 5

АНОМАЛЬНЫЕ ОТКЛОНЕНИЯ ОТ СХЕМЫ НЕЗАВИСИМОГО ПОРОЖДЕНИЯ

Текст	1-граммы	F _{мод}	F _{эксп}	1-граммы	F _{мод}	F _{эксп}	Коммен- тарий
ФХ174	АА ТА ТГ	310 403 392	395 312 480	ААА АТГ ЦТГ ИТГ ТАГ	74 94 63 68 94	133 140 31 20 24	иниц. терм.
G4	АГ АА ЦТ	300 413 391	186 554 503	ААА АГТ АТА ТАГ ТТГ	112 81 112 81 111	197 36 66 16 66	терм.
FD	АА АГ	388 325	544 244	ААА ЦАТ ТАГ АГТ ГАГ	95 110 112 112 67	201 69 58 53 38	терм.
SV40	ЦГ ТА ЦА	217 459 318	27 325 424	Все ЦГ-содержащие триплеты			
PBR322				ААА ЦАЦ АТЦ ТАГ	128 133 63 65 40 144	212 74 113 132 73 238	терм.
MS2	УЦ	228	263	ЦАЦ УЦГ	57 59	37 89	

Поскольку схема независимых испытаний является весьма грубым приближением к описанию генетических текстов, предпринимаются попытки использовать для описания отдельных частей геномов более сложные модели, например, в виде марковских цепей невысокого (первого и второго) порядка [5]. Оценки переходных вероятностей получаются при этом на основе 1-граммных характеристик соответствующих частей геномов. Такого рода кусочные представления могут быть использованы для решения задач разметки и генерации текстов.

Таблица 6

Коэффициенты ранговой корреляции

Пара текст/текст	1		
	2	3	4
ФХ174/БV40	0,55	0,43	0,37
ФХ174/G4	0,57	0,51	0,58
ФХ174/MS2	-0,08	0,01	
ФХ174/FD	0,86	0,77	0,67
ФХ174/PBR	0,10	0,13	0,21
BV40/G4	0,49	0,41	0,34
BV40/MS2	-0,30	-0,32	-0,15
BV40/FD	0,67	0,55	0,48
BV40/PBR	0,22	-0,04	-0,03
G4/MS2	-0,11	0,14	0,11
G4/FD	0,69	0,56	0,50
G4/PBR	0,24	0,15	0,18
MS2/FD	0,02	0,01	0,10
MS2/PBR	0,62	0,32	0,19
FD/PBR	0,04	-0,02	0,01

Весьма близки тексты ФХ174 и G4 ($r_1 > 0.5$ для всех 1), что неудивительно, так как эти тексты кодируют родственные вирусы. Неожиданным оказывается факт близости этих текстов с текстом FD, поскольку внутренняя организация вирусов ФХ174 и G4 сильно отличается от структуры FD. Объяснение может заключаться либо в том, что при различной внутренней организации похожими могут оказаться белки, кодируемые этими текстами, либо в том, что ранговые меры

4.6. Количественные оценки близости текстов. Простейшими мерами близости текстов являются коэффициенты ранговой корреляции r_1 , вычисленные с использованием частотных характеристик 1-го порядка (см.табл. 6, $l = 2,3,4$). Значения r_1 в целом согласуются с тем разбиением текстов на 2 класса (по ТА- и ЦГ-преобладанию), которое было приведено в п. 4.1, поскольку коэффициенты ранговой корреляции между представителями разных классов близки к нулю или отрицательны. Представляет интерес проследить степень близости между текстами одного класса.

близости не отражают адекватным образом различия во внутренней организации геномов. Последнее может объясняться тем, что ранговые меры близости не учитывают информации о расположении 1-грамм в тексте. Те же меры, которые это учитывают, весьма сложны для вычисления.

Тексты PBR322 и MS2, будучи весьма близкими по биграммным характеристикам ($\rho_2 = 0.62$), с увеличением 1 становятся менее похожими ($\rho_4 = 0.18$). Так как характеристики четвертого порядка учитывают большую информацию, чем характеристики второго порядка, эти тексты, видимо, не следует считать слишком близкими.

5. Заключительные замечания. Проведенный анализ затронул всего лишь 6 текстов, поэтому многие выводы, сделанные выше, носят предварительный характер и подлежат дальнейшему уточнению по мере накопления материала. Однако уже сейчас можно сказать, что 1-граммный способ описания генетических текстов, сводящийся фактически к поиску всех периодичностей в символьной последовательности, является в сочетании с позиционным анализом адекватным средством для решения многих содержательных генетических задач.

Обобщением 1-граммного способа представления символьных последовательностей является представление, допускающее объединение в один класс не только совпадающих 1-грамм, но и 1-грамм, отличающихся от них в пределах заданной меры различия. Такое описание представляет интерес и с биологической точки зрения, например, в задаче поиска гомологичных участков в родственных геномах.

Л и т е р а т у р а

1. PATHEP B.A. Молекулярно-генетические системы управления.- Новосибирск: Наука, 1975. - 286 с.
2. WAGNER B.A., FISHER M.I. The String-to-String Correction Problem.-J.Assoc.Comput.Machinery, 1974, v.21, N 1, p. 168-173.
3. MASAMI Hasegawa, TERUO Yasunaga, TAKASHI Miyata. Secondary structure of MS2 phage RNA and Bias in code word usage. - Nucl. Acid.Res., 1979, v.7, N 7, p.2073-2079.
4. ШЕННОН К. Математическая теория связи. - В кн.: Работы по теории информации и кибернетике. М., 1963, с. 243-332.
5. ERICKSON J.W., ALTMAN J.J. A Search for Patterns in the Nucleotide Sequence of the MS2 Genome.- J.Math.Biology, 1979, v.7, p.219-230.
6. BIRD A.P. DNA methylation and the frequency of CpG in animal DNA. - Nucl.Acid.Res., 1980, v.8, N 7, p.1499-1505.

Поступила в ред.-изд.отд.
5 ноября 1980 года