

72

Алгоритм определения вторичной структуры РНК

An algorithm that searches for a secondary structures in a given RNA sequences is suggested. The complexity of the algorithm is close to quadratic. The algorithm was tested on simulated and natural sequences

Приводится алгоритм решения одной из фундаментальных задач алгоритмической (или, как иногда говорят, вычислительной) генетики – алгоритм определения вторичной структуры РНК. Сообщение состоит из трех пунктов: математической постановки задачи, описания алгоритма (включая оценку времени его работы), нескольких слов, относящихся к биологической мотивировке. Алгоритм был компьютерно реализован, был проведен счет на модельных и реальных исходных данных (который показал его удовлетворительность); результаты счета, как и их интерпретации, представлены в журнал «Молекулярная биология». Строго математическое исследование этой задачи важно, но весьма трудно

1. Постановка задачи

Имеется последовательность нуклеотидов длины ДЛИНА (где ДЛИНА порядка 1000). Под отрезком $[a, b]$, где $a \leq b$, понимается множество натуральных чисел $\{a, a+1, \dots, b\}$. Два отрезка $[a, b]$ и $[c, d]$ комплементарны, если в исходной последовательности a -ый нуклеотид комплементарен d -ому, $(a+1)$ -ый нуклеотид комплементарен $(d-1)$ -ому и т.д. (в частности длины отрезков одинаковы). Шпилькой назовём последовательность комплементарных пар отрезков $[a_1, b_1], [c_1, d_1], \dots, [a_n, b_n], [c_n, d_n]$, все отрезки которой попарно не перекрываются, расположены на числовой оси в порядке возрастания следующим образом: $[a_1, b_1], [a_2, b_2], \dots, [a_n, b_n], [c_n, d_n], \dots, [c_1, d_1]$, длина каждого её отрезка не меньше МИНПОДРЯД (где МИНПОДРЯД – целое число порядка 3), все разности $a_2 - b_1, a_3 - b_2, \dots, d_3 - c_2, d_2 - c_1$ меньше МАКСВЫПЯЧИВАНИЕ (где МАКСВЫПЯЧИВАНИЕ – целое

2037

число порядка 20), а разность $(c_n - b_n)$ не больше МАКСПЕТЛЯ и не меньше МИНПЕТЛЯ (где МАКСПЕТЛЯ – целое число порядка 20, МИНПЕТЛЯ – целое число порядка 3) Кроме того требуется, чтобы *длина шпильки* (определяемая как $d_1 - a_1 + 1$) не превосходила ДЛИНА ШПИЛЬКИ (где ДЛИНА ШПИЛЬКИ – целое число порядка 40)

Мы говорим, что отрезок $[a_1, b_1]$ склеен с отрезком $[c_1, d_1]$, (a_1 с d_1 , $a_1 + 1$ с $d_1 - 1$ и т.д.), отрезок $[a_2, b_2]$ – с отрезком $[c_2, d_2]$ и т.д. Отрезок $[a_1, b_n]$ называется *левым основанием*, отрезок $[c_n, d_1]$ – *правым основанием*. Число a_1 называется *началом шпильки*, d_1 – *концом шпильки*, b_n – *началом петли*, c_n – *концом петли*. Также $[a_1, b_1]$ называется *внешним краем*, а $[c_n, b_n]$ – *внутренним краем* шпильки

Перекрытием называется пара шпилек, у которых, во-первых, правое основание первой шпильки имеет не менее ДЛИНА ПЕРЕСЕЧЕНИЯ общих элементов с левым основанием второй шпильки (где ДЛИНА ПЕРЕСЕЧЕНИЯ – целое число порядка 5), во-вторых, конец первой шпильки находится левее конца петли второй шпильки и, в-третьих, начало второй шпильки находится правее начала петли первой шпильки. Левое и правое основания первой шпильки перекрытия будут обозначаться соответственно $[A', B']$ и $[C', D']$, а второй шпильки – $[A, B]$ и $[C, D]$. Таким образом, $D' < C$ и $A > B'$. Пересечение отрезков $[C', D']$ и $[A, B]$ называется *базой перекрытия* (как говорилось его длина не менее ДЛИНА ПЕРЕСЕЧЕНИЯ)

Мощностью шпильки называется количество склеенных пар нуклеотидов. *Мощностью перекрытия* называется сумма мощностей его шпилек. Два перекрытия назовем *подобными*, если у них одинаковые мощности, а также одинаковые левые и правые основания обеих шпилек

Итак, постановка задачи. Приблизительно говоря, алгоритм получив на вход исходную последовательность и число КОЛИЧЕСТВО ПЕРЕКРЫТИЙ (где КОЛИЧЕСТВО ПЕРЕКРЫТИЙ – целое число порядка 5) должен напечатать КОЛИЧЕСТВО ПЕРЕКРЫТИЙ перекрытий наибольшей мощности, при этом из подобных перекрытий печатается не более одного. Точнее говоря, алгоритм будет находить следующие перекрытия. Рассмотрим на множестве перекрытий линейный порядок чтобы сравнить два перекрытия, сначала сравниваются их мощности, при равных мощностях сравниваются A , затем D' , затем C, D, B, B', A', C' в указанном порядке. При этом чем больше мощность, больше D', C, A', C' и меньше A, D, B, B' , тем перекрытие «больше». Если

отождествлять подобные перекрытия, то указанное отношение действительно является линейным порядком. Для каждого X ($1 \leq X \leq \text{ДЛИНА} - \text{ДЛИНА ПЕРЕСЕЧЕНИЯ} + 1$) определим $Y = Y(X) = X + \text{ДЛИНА ПЕРЕСЕЧЕНИЯ} - 1$. Определим $K(X)$ как *наибольшее перекрытие*, база которого включает отрезок $[X, Y]$. Рассмотрим последовательность перекрытий P_1, P_2, \dots, P_n , определяемую по индукции так. Пусть P_1, P_2, \dots, P_i уже определены. Пусть M наибольшая из мощностей перекрытий, неподобных ни одному из перекрытий P_1, P_2, \dots, P_i . Пусть X наименьшее такое, что $K(X)$ имеет мощность M и не подобно ни одному из перекрытий P_1, P_2, \dots, P_i . Тогда $P_{i+1} = K(X)$. Тогда алгоритм находит КОЛИЧЕСТВО ПЕРЕКРЫТИЙ первых элементов этой последовательности, если ее длина не менее КОЛИЧЕСТВО ПЕРЕКРЫТИЙ, а иначе - находит всю эту последовательность.

Верхняя оценка времени работы алгоритма первый этап $O(\text{ДЛИНА} \times \text{ДЛИНА ШПИЛЬКИ}^2 \times \text{МАКСВЫПЯЧИВАНИЕ})$ арифметических операций, второй этап $O(\text{КОЛИЧЕСТВО ПЕРЕКРЫТИЙ} \times \text{ДЛИНА ШПИЛЬКИ} \times \text{МАКСВЫПЯЧИВАНИЕ} \times (\text{ДЛИНА ШПИЛЬКИ} + \text{МАКСВЫПЯЧИВАНИЕ}))$ арифметических операций.

2. Алгоритм решения задачи

Этап 1. На этом этапе мы для всех X находим две четверки (m, A, C, D) и (m', A', B', D') такие, что первая шпилька перекрытия $K(X)$ имеет мощность m' , основания $[A', B']$ и $[C', D']$, а вторая шпилька перекрытия $K(X)$ - соответственно m и $[A, B]$ и $[C, D]$ (для некоторых B, C'). К сожалению, значения B, C' мы не можем вычислять быстро. Мощность $K(X)$ храним в массиве «мощность» ($\text{мощность}[X] = \text{мощность } K(X)$). Любое перекрытие, база которого включает отрезок $[X, Y]$, целиком лежит на отрезке $[\text{НАЧАЛО}, \text{КОНЕЦ}]$, где $\text{НАЧАЛО} = \max\{1, Y - \text{ДЛИНА ШПИЛЬКИ}\}$, $\text{КОНЕЦ} = \min\{\text{ДЛИНА}, X + \text{ДЛИНА ШПИЛЬКИ}\}$.

Для нахождения значения $\text{мощность}[X]$ мы сначала рассматриваем шпильки, начало и конец которых принадлежат отрезку $[\text{НАЧАЛО}, \text{КОНЕЦ}]$, начало петли больше или равно Y ($B \geq Y$). Назовем такие шпильки *допустимыми*. Определяем качество такой шпильки как тройку (m, C, D) , где m - ее мощность, C - конец петли, а D - конец. Качество тем лучше, чем больше m и C и меньше D . Точнее, на множестве таких троек мы рассматри-

вам отношение частичного порядка $(m, C) \leq (m', C')$, если $m \leq m'$ и $C \leq C'$, $(m, C, D) \leq (m', C', D')$, если $(m, C) < (m', C')$ или $(m, C) = (m', C')$ и $D \leq D'$. То есть на парах (m, C) порядок локомпонентный, а на тройках лексикографический, если тройку (m, C, D) понимать как пару $((m, C), D)$. При каждом A из отрезка [НАЧАЛО, КОНЕЦ] мы находим множество, называемое максшпилька2[X, A], состоящее из максимальных возможных качеств допустимых шпилек с началом A . Все тройки из множества максшпилька2[X] несравнимы, поэтому при каждом m в этом множестве не более одной тройки вида (m, C, D) . Множество максшпилька2[X] хранится в виде отсортированного (по возрастанию m) массива троек.

Затем, симметрично, рассматриваем шпильки, начало и конец которых принадлежат отрезку [НАЧАЛО, КОНЕЦ], конец петли меньше или равен X ($C' \leq X$). Определяем качество такой шпильки как пару (m, B', A') , где m - ее мощность, а B' - начало петли. Качество тем лучше, чем больше m, A' и меньше B' . Точнее, на множестве таких троек мы рассматриваем отношение частичного порядка $(m, B') \leq (m', B'')$, если $m \leq m'$ и $B' \geq B''$, $(m, B', A') \leq (m', B'', A'')$, если $(m, B') < (m', B'')$ или $(m, B') = (m', B'')$ и $A' \leq A''$. При каждом D' из отрезка [НАЧАЛО, КОНЕЦ] мы находим множество максшпилька1[X, D'] максимальных возможных качеств таких шпилек.

Теперь мы можем найти мощность[X], сравнивая шпильки обоих типов. A именно, мощность[X] равно наибольшей сумме $m+m'$, для которой *существует отрезок* [A, D'], содержащий отрезок [X, Y], со следующим свойством (m', B', A') - тройка из списка максшпилька1[D'] с наибольшим m' , для которой $B' < A$, а (m, C, D) - тройка из списка максшпилька2[A] с наибольшим m , для которой $C > D'$. Находим мощность[X] перебором отрезков [A, D']. Затем еще раз перебираем [A, D'], чтобы найти (m', A', B', D') и (m, A, C, D) .

Как при фиксированном X заполнить массив максшпилька2? Мы перебираем K из отрезка [Y, КОНЕЦ] в порядке возрастания. При каждом K рассматриваем шпильки, начало и конец которых принадлежат отрезку [НАЧАЛО, K], а начало петли больше или равно Y . Будем называть такие шпильки *K-допустимыми*. При $K=KОНЕЦ$ понятия *K-допустимости* и *допустимости* совпадают. При каждом A из отрезка [НАЧАЛО, K] мы находим множество максимальных возможных качеств *K-допустимых* шпилек с началом A и полагаем максшпилька2[A, K, X]

равным этому множеству. При $K=\text{КОНЕЦ}$ множество максшпилька2[A.K.X] будет хранить искомые тройки

Для того, чтобы прошел индуктивный переход $K \rightarrow K+1$, мы храним дополнительную информацию А именно, мы отдельно рассматриваем K -допустимые шпильки, у которых конец принадлежит отрезку [K-МАКСВЫПЯЧИВАНИЕ.K]. Назовем такие шпильки *K-продолжаемыми* (K-Продолжаемые шпильки можно продолжать, склеивая символы с номерами, большими K) *Качество K-продолжаемой шпильки* определяется как пара $(m.C)$, где m – ее мощность, а C – конец петли. Качества сравниваются покомпонентно $(m.C) \leq (m'.C')$, если $m \leq m'$ и $C \leq C'$. Информацию о K -продолжаемых шпильках мы храним в виде *очереди продолж[A]*. А именно, *продолж[A]* есть очередь длины МАКСВЫПЯЧИВАНИЕ+1 n -ый элемент которой есть множество максимальных качеств K -продолжаемых шпилек с началом А и концом K-МАКСВЫПЯЧИВАНИЕ+n-1 (Элементы в очередь добавляются с конца, а удаляются с начала, нумеруются элементы очереди с начала 1,2, ...)

Более того, мы рассматриваем и так называемые *внутренние предшпильки* (или, кратко, предшпильки). Грубо говоря, внутренняя предшпилька - это то, что останется от шпильки, если ее разделить на две части и удалить внешнюю часть. Формально (внутренняя) предшпилька отличается от шпильки тем, что в предшпильке длины внешних отрезков $\{[a_1.b_1]\}$ и $\{[c_1.d_1]\}$ могут быть меньше МИНПОДРЯД, длина внешних отрезков называется *кратностью предшпильки*. *Предшпилька K-допустима*, если ее конец равен K, а начало петли больше или равно Y. K -допустимые предшпильки нужны, поскольку их можно продолжать до шпилек, склеивая символы с номерами, большими K. Качество предшпилек определяется так же, как и качество K -продолжаемых шпилек. В массиве «предшпилька» мы храним информацию о них. А именно *предшпилька[A]* есть антиочередь длины МИНПОДРЯД-1, g -ый элемент которой есть множество максимальных качеств предшпилек кратности g (где $g=1, \dots, \text{МИНПОДРЯД}-1$) (В антиочередь элементы добавляются с начала, а удаляются с конца, нумеруются элементы антиочереди по-прежнему с начала 1,2, ...)

Множества из (анти)очереди предшпилька[A], продолж[A] хранятся также в виде упорядоченных по возрастанию первого компонента массивов чисел.

Как же совершить переход от K к $K+1$? Мы имеем правильно заполненные *продолж[A]*, *предшпилька[A]*, максш-

пильтка2|A| для данного K. Ясно, что если K+1 не комплементарно A, то K+1-допустимых предшпилек нет вовсе, множество K+1-допустимых шпилек совпадает с множеством K-допустимых шпилек, а множество K+1-продолжаемых шпилек равно множеству K-продолжаемых шпилек минус шпильки с концом K-МАКСВЫПЯЧИВАНИЕ. То есть, в этом случае для K+1 антиочередь предшпильтка|A| состоит из пустых множеств, множество максшпильтка2|A| не изменяется, а из очереди продолж|A| удаляется первое множество, и в ее конец вставляется пустое множество.

Более сложным образом изменяются эти массивы, если K+1 комплементарно A. В этом случае к K+1-допустимым шпилькам с началом A добавляются шпильки, полученные склеиванием A с K+1 у K-допустимых предшпилек кратности МИНПОДРЯД-1 с началом A+1, а также и шпильки, полученные склеиванием A с K+1 у шпилек с началом A+1 и концом K. Эти шпильки являются K+1-продолжаемыми, поэтому они добавляются к K+1-продолжаемым шпилькам. Из множества K-продолжаемых шпилек удаляются все шпильки с концом K-МАКСВЫПЯЧИВАНИЕ и оставшиеся шпильки вместе с вышеуказанными K+1-продолжаемыми шпильками образуют все множество K+1-продолжаемых шпилек. Множество K+1-допустимых предшпилек с началом A формируется так берутся все K-допустимые предшпильтки с началом A+1 кратности меньше МИНПОДРЯД-1, и в каждой из них склеиваются A и (K+1). Кроме этого, могут добавиться предшпильтки кратности 1.

Они берутся, во-первых, если у некоторой K-продолжаемой шпильки с началом, принадлежащим отрезку [A+1, A+МАКСВЫПЯЧИВАНИЕ+1] склеить еще A с K+1, и, во-вторых, если просто склеить A с K+1, получив предшпильтку мощности 1. Последнее возможно, если K-A находится в пределах от МИНПЕТЛЯ до МАКСПЕТЛЯ.

Итак, для получения новых значений максшпильтка2|A|, продолж|A|, предшпильтка|A| делаем следующее. Удаляем из очереди продолж|A| первое множество. Сливаем последнее множество из очереди продолж|A+1| с последним множеством из антиочередь предшпильтка|A+1| (именно в них хранится информация о K-допустимых шпильках с началом A+1 и концом K и о предшпильках кратности МИНПОДРЯД-1 с началом A+1). Об операции слияния двух множеств будет сказано позже. Увеличиваем первый компонент всех пар в полученном множестве на 1 (получим множество всех максимальных качеств продолжаемых шпилек).

лек с началом A и концом $K+1$) и ставим полученное множество в конец очереди продолж $[A]$. Теперь добавляем ко всем парам этого множества $K+1$ в качестве третьего компонента и результат вбиваем в множество максшпилька2 $[A]$

Удаляем из антиочереди предшпилька $[A+1]$ последнее множество и увеличиваем первые компоненты пар всех множеств на 1. Это будет концом антиочереди предшпилька $[A]$. Осталось найти множество, которое надо поставить в начало этой антиочереди, то есть множество максимальных качеств предшпилек кратности 1 с началом A и концом $K+1$. Для этого сливаем вместе все множества, хранимые во всех очередях продолж $[A+1]$, продолж $[A+МАКСВЫПЯЧИВАНИЕ+1]$. Получим множество максимальных качеств всех продолжаемых шпилек с началом на отрезке $[A+1, A+МАКСВЫПЯЧИВАНИЕ+1]$. Увеличим первый компонент всех пар этого множества на 1 и поставим в начало антиочереди предшпилька $[A]$. Наконец если $МИНПЕТЛЯ \leq K - A \leq МАКСПЕТЛЯ$, то добавляем в это множество еще пару $(1, K+1)$. Для того, чтобы новые значения (анти)очереди не испортили старые, *нужно перебирать A в порядке возрастания*.

Слияние всех множеств, хранимых во всех очередях продолж $[A+1]$, продолж $[A+МАКСВЫПЯЧИВАНИЕ+1]$ значительно сложнее всех остальных операций, поскольку сливается вместе $(МАКСВЫПЯЧИВАНИЕ+1)^2$ множеств. Назовем множество, полученное в результате этого, через *соседи $[A]$* . Его можно вычислить индукцией по A более простым способом. Для этого перебираем A в порядке убывания и при каждом A вычисляем антиочередь *que*е длины $МАКСВЫПЯЧИВАНИЕ+1$ n -ый элемент которой есть слияние всех множеств очереди продолж $[A+n]$. Множество *соседи $[A]$* получается в результате слияния всех элементов антиочереди *que*е (что требует слияния $МАКСВЫПЯЧИВАНИЕ+1$ множеств). При уменьшении A на единицу из антиочереди *que*е исключается последнее множество a в начало ставится слияние всех множеств из очереди продолж $[A]$ (что опять требует слияния $МАКСВЫПЯЧИВАНИЕ+1$ множеств). Таким образом для получения множества *соседи $[A]$* мы делаем два слияния $МАКСВЫПЯЧИВАНИЕ+1$ множеств вместо одного слияния $(МАКСВЫПЯЧИВАНИЕ+1)^2$ множеств.

Массив максшпилька1 заполняется симметричным образом (K перебираем в порядке убывания)

Теперь обсудим операцию слияния множеств пар вида (m, c) . Напомним, что сливаемые множества хранят попарно несравнимые пары, то есть не более одной пары с данным пер-

вым компонентом, и хранятся в виде упорядоченного массива пар. Операция слияния состоит просто в соединении двух упорядоченных массивов в новый упорядоченный массив (как в алгоритме сортировки слиянием). При этом, если окажутся рядом две пары с одним и тем же первым компонентом, то меньшая пара удаляется. Операция слияния нескольких множеств производится естественным образом: первое сливается со вторым, результат сливается с третьим и т.д. Следует заранее исключить из сливаемых множеств пустые, поскольку в типичном случае их должно быть от половины до трех четвертей (шпилька с началом А и концом К существует только если А комплементарно К).

Множества троек вида (m,C,D) сливаются точно так же. Влить множество M1 в множество M2 означает, слить M1 и M2 и результат поместить в M2.

На хранение одного множества пар вида (m,C) следует резервировать $2 \times \text{ДЛИНА_ШПИЛЬКИ}$ байт, а на хранение одного множества пар вида (m,C,D) – $3 \times \text{ДЛИНА_ШПИЛЬКИ}$ байт.

Этап 2. Сначала объясним, как, имея M (A,C,D) и Y, найти наименьшее такое $B \supseteq Y$, что есть шпилька мощности M с основаниями [A,B], [C,D] и саму шпильку, при условии, что M – наибольшая возможная мощность такой шпильки. Применяв этот алгоритм и вместе с симметричным ему, мы можем для каждого X найти K(X) (но делать этого пока не надо).

Алгоритм похож на алгоритм первого этапа. Дадим определение *внешней предшпильки*. Внешняя предшпилька отличается от шпильки тем, что мы разрешаем внутренним (а не внешним) отрезкам, как у внутренних предшпилек, иметь длину меньше МИНПОДРЯД (эту длину называем по-прежнему кратностью), а расстояние между началом и концом петли не ограничено сверху. Мы рассматриваем только предшпильки мощности не более M (более мощные предшпильки неперспективны).

Мы перебираем B в порядке возрастания, начиная с $B=A$. При фиксированном B перебираем c из отрезка [C,D] в порядке убывания и рассматриваем внешние предшпильки с основаниями [A,B] и [c,D] (назовем их B,c-допустимыми). Для каждого c из отрезка [C,D] мы вычисляем $\text{продолж}[c,B]$ = максимальная мощность B,c-допустимой предшпильки кратности не менее МИНПОДРЯД (назовем их полными), и $\text{неполн}[c,B]$ = антиочередь длины МИНПОДРЯД-1, g-ый элемент которой есть максимальная мощность B,c-допустимой предшпильки кратности g (допустимой предшпильки кратности менее МИНПОДРЯД будем называть

неполными) Если указанных предшпилек нет, то максимальная мощность считается равной минус бесконечность

Как осуществлять переход $V \rightarrow V+1$? Пусть $c-1$ комплементарно $V+1$. Какие бывают $V+1.c-1$ -допустимые предшпильки? Полные $V+1.c-1$ -допустимые предшпильки получаются склеиванием $c-1$ и $V+1$ из полных $V.c-1$ -допустимых предшпилек и из неполных $V.c-1$ -допустимых предшпилек кратности МИНПОДРЯД-1. Неполные $V+1.c-1$ -допустимые предшпильки кратности 1 получаются из неполных $V.c-1$ -допустимых предшпилек кратности $g-1$ склеиванием $c-1$ и $V+1$. Кроме того, любую $b.c'-1$ -допустимую полную предшпильку, где b лежит на отрезке $[V-МАКСВЫПЯЧИВАНИЕ.V]$ и c' лежит на отрезке $[c.c+МАКСВЫПЯЧИВАНИЕ]$, можно продолжить, склеив $c-1$ и $V+1$.

Иными словами, если $c-1$ комплементарно $V+1$, то $\text{полн}[c-1.V+1]$ есть $1 + \max\{\text{полн}[c.V], \text{последний элемент антиочередности неполн}[c.V]\}$. Антиочередность $\text{неполн}[c-1.V+1]$ получается из антиочередности $\text{неполн}[c.V]$ удалением последнего элемента, увеличением всех элементов на 1 и добавлением в начало элемента $1 + \text{neighbor}(c-1)$, где $\text{neighbor}(c-1) =$ наибольшая мощность полной $b.c'-1$ -допустимой полной предшпильки, где b лежит на отрезке $[V-МАКСВЫПЯЧИВАНИЕ.V]$ и c' лежит на отрезке $[c.c+МАКСВЫПЯЧИВАНИЕ]$.

В начальный момент (при $V=A$) мы полагаем $\text{неполн}[C.A] = (1 \text{ минус бесконечность}, \dots \text{ минус бесконечность})$, поскольку имеется единственная $A.D$ -допустимая предшпилька, она имеет мощность 1 и кратность 1, и полагаем $\text{неполн}[c.A] = (\text{минус бесконечность}, \dots \text{ минус бесконечность})$ для всех остальных c .

Число $\text{neighbor}(c)$ вычисляется так. Пусть $\text{queue}(c)$ – антиочередность, n -ый элемент которой равен максимуму мощностей $b.c+n$ -допустимых полных предшпилек по всем b на отрезке $[V-МАКСВЫПЯЧИВАНИЕ.V]$. Тогда $\text{neighbor}(c)$ есть максимальный из элементов антиочередности $\text{queue}(c)$. А $\text{queue}(c)$ легко найти по индукции: при переходе от c к $c-1$ конец антиочередности удаляется, а в начало ставится наибольший из элементов $\text{полн}[c.V-МАКСВЫПЯЧИВАНИЕ]$, \dots , $\text{полн}[c.V-1]$, $\text{полн}[c.V]$.

Перебор V останавливаем, когда $\text{полн}[C.V]=M$ и при этом $Y \leq V$ МИНПЕТЛЯ $\leq C-V+1 \leq$ МАКСПЕТЛЯ (поскольку в этом случае имеется шпилька мощности M с основаниями $\{A.V\}, \{C.D\}$, где $Y \leq V$). Такое V обязательно найдется. Действительно, в качестве V годится любое $V \geq Y$, для которого есть шпилька мощности M с основаниями $\{A.V\}, \{C.D\}$.

Найдя V ищем самую шпильку просматривая созданные нами массивы и наращивая ее «изнутри» То есть мы будем хранить некоторую внутреннюю предшпильку - часть будущей шпильки Как хранить - неважно например ее можно задавать как список склеенных пар нуклеотидов с названием «шпилька» Кроме списка шпилька мы отдельно храним четверку m', b, c, s где m' - мощность хранимой внутренней предшпильки s - ее кратность a, b, c - ее начало и конец При этом мы поддерживаем следующий инвариант существует b, c -допустимая внешняя предшпилька кратности не менее МИНПОДРЯД- $s+1$ и мощности $M-m'+1$ (ее соединение с хранимой внешней предшпилькой даст искомую шпильку)

Начинаем с того что кладем в список шпилька пару (B, C) и полагаем $b=B, c=C, m'=1, s=1$ Затем пока $m' < M$ выполняем следующее Если $\text{полн}[c+1, b-1] = M-m'$ или в антиочереди $\text{неполн}[c+1, b-1]$ для некоторого $r \geq \text{МИНПОДРЯД}-s$ r -ый элемент равен $M-m'$ то уменьшаем b и увеличиваем c на единицу увеличивая s на единицу (то есть склеиваем s хранимой предшпильки соседние с началом и концом нуклеотиды) Иначе (это может быть только если $s \geq \text{МИНПОДРЯД}$) для некоторого c' из отрезка $[c+1, c+1+\text{МАКСВЫПЯЧИВАНИЕ}]$ и некоторого b' из отрезка $[b-1, b-1-\text{МАКСВЫПЯЧИВАНИЕ}]$ выполнено $\text{полн}[c', b'] = M-m'$

Пологаем $b=b', c=c', s=1$ Наконец заносим пару (b, c) в список шпилька и полагаем $m'=m'+1$

Итак после первого этапа для всех X известна максимальная мощность (мощность $[X]$) перекрытия база которого включает отрезок $[X, Y]$ (где $Y = X + \text{ДЛИНА ПЕРЕСЕЧЕНИЯ}-1$) Кроме того по любому X можно найти $K(X)$

Находим M - наибольшее значение в массиве мощность $[X]$ Находим первое в порядке возрастания X_1 для которого мощность $[X_1] = M$ Находим $K(X_1)$ Ясно что $K(X_1)$ есть первый элемент искомого последовательности Пусть $[R_1, S_1]$ - его база Ясно что если отрезок $[X, Y]$ включен в $[R_1, S_1]$ то $K(X) = K(X_1)$ Поэтому находим X_2 - наименьшее X для которого мощность $[X] = M$ и $Y(X) > S_1$ Находим $K(X_2)$ Ясно что $K(X_2) = P_2$ Пусть $[R_2, S_2]$ - его база

Так действуем до тех пор пока очередное X_i существует После этого находим M' - наибольшую мощность $[X]$ для таких X что мощность $[X] < M$ Далее делаем то же самое заменив M на M' В этом духе действуем до тех пор пока не наберется КОЛИЧЕСТВО ПЕРЕКРЫТИИ перекрытий или очередное M'' не станет равным нулю **Конец описания алгоритма.**

Приведем один пример исходной последовательности (для аттенюатора гистидинового оперона кишечной палочки Eс|hisL)

tagetaattgtacgcatgtcaatctcctcttttgtacagttcattgtacaatgatgagcgtaatta
actatttattaattagtttftagatcaaggatttftcagtgagacgaaaatccaggtctgctatttttggf
gccatcagctaagaggacagctectcttagcccccctcttccccgcctatcattaaacaatccattg
ccataaaatalataaaaaagccccttgccttctaacgtgaaagtggft

taggtfaaaagacatcagttgaataaacattcacagagacttttatgacacgcgttc.aattfaa
acaccaccateatcaccatcatcctgactagtctttcaggcgatgtgtgctggaagacattcagatctt
ccagtggtgcatgaacgcgatgagaaagccccgggaagatcacctccgggggctttttattgcgcg
gttgataacgggtcagacaggtttaaagagggaataacaaaatgacagacaacactcgttaecgata
gcatgcagaaatccggccglttaagtgatgaetcacgcgaattgctggcgcgctgtggcattaaat
taacttcacacccagcgctgategcgatggcagaaaacatgccattgatattctgcgcgctgcgtg
acgaegacattccggctctggtaatggatggcgtggttagacctgggattatcggcgaaaacgtgct
ggaagaagagctgcttaaccgcgcgcccagggtgaagatccacgctactttacctgctgctct
ggatttcggcggctgctcttctctggcaacgccggttgatgaagcctgggacggtc

3 Биологическая мотивировка

Речь идет о задаче поиска участков генома потенциально годных для образования альтернативных вторичных структур в РНК

Регуляция многих бактериальных генов осуществляется на уровне трансляции или взаимодействия процессов трансляции и транскрипции. При этом основным регуляторным сигналом является образование вторичной структуры РНК. Часто (хотя и не всегда) в рассматриваемой области могут образоваться две альтернативные структуры РНК, что и служит регуляторным переключателем. Примерами такой регуляции являются аттенюаторы оперонов синтеза аминокислот [1] и регуляторные структуры некоторых оперонов рибосомальных белков кишечной палочки [2].

Литература

1. Витресцак ИИ, Гельфанд МС. Сравнительный подход к анализу регуляции в полных геномах: аттенюаторы ароматических аминокислотных оперонов гамма-протеобактерии. Молекулярная биология 34(4) в печати.
2. Витресцак И, Бансат АК, Литов ПП, Гельфанд МС (1999) Компьютерный анализ регуляторных сигналов в полных бактериальных геномах. Инициация трансляции оперонов рибосомальных белков. Биофизика 44: 601-610.