

УДК 621.391 : 519.724.2

© 2022 г. М. Ковачевич

**О МАКСИМАЛЬНОМ ЧИСЛЕ РАЗЛИЧИМЫХ СТРОК  
ПОД ДЕЙСТВИЕМ КОРОТКИХ ТАНДЕМНЫХ ДУПЛИКАЦИЙ<sup>1</sup>**

Множество всех  $q$ -ичных строк, не содержащих повторяющихся подстрок длины  $\leq 3$  (т.е. не содержащих подстрок вида  $aa$ ,  $abab$  и  $abcabc$ ), образует код, исправляющий произвольное количество мутаций типа тандемных дупликаций длины  $\leq 3$ . Иными словами, любые две такие строки различимы в том смысле, что из них не могут образоваться одинаковые строки под действием некоторого количества тандемных дупликаций длины  $\leq 3$ . Показано, что этот код асимптотически оптимален по скорости, т.е. представляет собой максимальное множество различных строк с точностью до субэкспоненциального множителя. Этот результат дает решение задачи о пропускной способности с нулевой ошибкой для последнего оставшегося случая каналов с тандемными дупликациями, удовлетворяющих свойству “единственности корневых строк”.

*Ключевые слова:* тандемная дупликация, повторение тандема, ошибка дупликации, ошибка повторения, вклейка, хранение информации на основе ДНК, исправление ошибок, пропускная способность с нулевой ошибкой, код с ограничением, строка без повторений.

**DOI:** 10.31857/S0555292322020028, **EDN:** DZAPUX

**§ 1. Введение**

Тандемные дупликации – это ошибки типа “вклейки”, естественным образом появляющиеся как мутации ДНК-строк и являющиеся поэтому потенциальным источником дефектов, возникающих в системах хранения информации на основе ДНК в живых организмах [1]. В то время как задача исправления тандемных дупликаций фиксированной и известной длины  $\ell$  хорошо изучена, причем как при конечном [2, 3], так и при бесконечном [1, 4] числе ошибок, в отношении (по-видимому, гораздо более практически значимой) задачи исправления дупликаций переменной длины известно гораздо меньше. В частности, оптимальные коды, исправляющие все конфигурации дупликаций длины  $\leq \ell$ , были найдены только в специальных случаях  $\ell = 1$  и  $\ell = 2$  [1, теорема 32].

Нашим основным результатом является доказательство того факта, что аналогичная конструкция кодов, исправляющих неограниченное число тандемных дупликаций длины  $\leq 3$  [1, теорема 27], также является асимптотически оптимальной по скорости. Этот результат дает решение задачи о пропускной способности с нулевой ошибкой для каналов с тандемными дупликациями во всех случаях, когда корневые строки (относительно дупликаций) любых строк единственны [1, теорема 40]. Однако при больших значениях  $\ell$  корневые строки относительно дупликаций длины  $\leq \ell$  перестают быть единственными [1, теорема 40], и следовательно, для решения за-

<sup>1</sup> Работа выполнена при поддержке исследовательско-инновационной программы Horizon 2020 Европейского Союза (номер гранта 856967), а также Секретариата по высшему образованию и научным исследованиям автономной провинции Воеводина, Сербия (номер проекта 142-451-2686/2021).

дачи о пропускной способности с нулевой ошибкой и смежных с ней задач в таких моделях понадобятся другие конструкции и верхние границы<sup>2</sup>.

Помимо теоретико-информационных и кодовых вопросов упомянутого типа в литературе рассматривались и некоторые другие задачи, связанные с моделями с тандемными дубликациями переменной длины (см., например, [6–8]).

**1.1. Описание модели.** Для  $q$ -ичного алфавита будем использовать обозначение  $\mathcal{A}_q := \{0, 1, \dots, q-1\}$ , а множество всех строк (или слов) над алфавитом  $\mathcal{A}_q$  обозначим через  $\mathcal{A}_q^* := \bigcup_{n=0}^{\infty} \mathcal{A}_q^n$ . Длина строки  $\mathbf{x} = x_1 \dots x_n \in \mathcal{A}_q^n$  обозначается через  $|\mathbf{x}| = n$ . Строка, полученная конкатенацией двух строк  $\mathbf{u}$  и  $\mathbf{v}$ , записывается в виде  $\mathbf{uv}$ . Строка  $\mathbf{v}$  называется подстрокой (или фрагментом) строки  $\mathbf{x}$ , если существуют (возможно, пустые) строки  $\mathbf{u}$  и  $\mathbf{w}$ , такие что  $\mathbf{x} = \mathbf{uvw}$ .

Пусть  $\ell$  – фиксированное натуральное число. В канале с тандемными ( $\leq \ell$ )-дубликациями передаваемая по каналу строка  $\mathbf{x}$  подвергается последовательному воздействию некоторого числа тандемных дубликаций длины  $\leq \ell$  каждая, где тандемная дубликация длины  $k$  – это вставка точной копии подстроки длины  $k$  рядом с исходной подстрокой (неважно, вставляется ли она слева или справа от исходной, поскольку обе операции приводят к одной и той же новой строке). Предполагается, что число возникающих дубликаций заранее не известно ни передатчику, ни приемнику и может принимать любое значение из множества натуральных чисел  $\{0, 1, 2, \dots\}$ . Более точно, канал описывается следующим образом:

- Вход:  $\mathbf{x} \equiv \mathbf{x}^{(0)}$ ;
- Из множества  $\{0, 1, 2, \dots\}$  выбирается число  $t$  (количество дубликаций);
- Для  $i = 1, \dots, t$  повторяются следующие действия:
  - Из множества  $\{1, \dots, |\mathbf{x}^{(i-1)}|\}$  произвольным образом выбирается позиция дубликации  $j$  в строке  $\mathbf{x}^{(i-1)}$ ;
  - Из множества  $\{1, \dots, \min\{j, \ell\}\}$  произвольным образом выбирается длина дубликации  $k$ ;
  - Строка  $\mathbf{x}^{(i)}$  получается вставкой копии подстроки  $x_{j-k+1}^{(i-1)} \dots x_j^{(i-1)}$  рядом с исходной подстрокой в  $\mathbf{x}^{(i-1)}$ , т.е.

$$\mathbf{x}^{(i)} = x_1^{(i-1)} \dots x_{j-k+1}^{(i-1)} \overline{x_{j-k+1}^{(i-1)} \dots x_j^{(i-1)}} x_{j-k+1}^{(i-1)} \dots x_j^{(i-1)} \overline{x_{j+1}^{(i-1)} \dots x_{|\mathbf{x}^{(i-1)}|}^{(i-1)}}$$

где исходная дублицируемая подстрока выделена чертой сверху, а вставленная копия – чертой снизу;

- Выход:  $\mathbf{y} \equiv \mathbf{x}^{(t)}$ .

Всюду далее будем считать, что  $\ell = 3$ .

Пример 1. Примером того, как канал действует на передаваемую строку  $\mathbf{x} \in \mathcal{A}_3^8$ , является следующий список строк, каждая строка в котором получена из предыдущей путем тандемной дубликации длины  $\leq 3$ :

$$\mathbf{x} = 0\ 1\ 1\ 2\ 0\ 2\ 1\ 0, \tag{1.1a}$$

$$\mathbf{x}^{(1)} = \overline{0}\ \underline{0}\ 1\ 1\ 2\ 0\ 2\ 1\ 0, \tag{1.1b}$$

$$\mathbf{x}^{(2)} = 0\ 0\ 1\ 1\ \overline{2}\ \underline{0}\ \underline{2}\ \underline{0}\ \underline{2}\ 1\ 0, \tag{1.1c}$$

$$\mathbf{x}^{(3)} = 0\ 0\ 1\ 1\ 2\ 0\ 2\ 2\ 0\ \overline{2}\ \underline{1}\ \underline{2}\ 1\ 0, \tag{1.1d}$$

$$\mathbf{x}^{(4)} = 0\ 0\ \overline{1}\ \overline{1}\ \underline{1}\ \underline{1}\ 2\ 0\ 2\ 2\ 0\ 2\ 1\ 2\ 1\ 0. \tag{1.1e}$$

Здесь  $t = 4$ , и выходом канала является строка  $\mathbf{y} = \mathbf{x}^{(4)}$ .

<sup>2</sup> Первые конструкции кодов для таких моделей (при  $\ell \in \{4, 5, \dots\}$ ) были приведены в [5].

Будем говорить, что строка  $\mathbf{y}$  является  $t$ -потомком строки  $\mathbf{x}$ , или что  $\mathbf{x}$  является  $t$ -предком  $\mathbf{y}$ , если  $\mathbf{y}$  можно получить из  $\mathbf{x}$  путем последовательного применения  $t$  тандемных дубликаций длины  $\leq 3$ . Множество всех  $t$ -потомков строки  $\mathbf{x}$  обозначим через  $D^t(\mathbf{x})$ . Отметим, что строка может принадлежать различным множествам  $D^t(\mathbf{x})$  и  $D^s(\mathbf{x})$ ,  $s \neq t$ , поскольку в модели разрешены дубликации различных длин, т.е. множество  $D^t(\mathbf{x}) \cap D^s(\mathbf{x})$  не обязано быть пустым (например, 01111 является как 1-потомком строки 011, полученным при одной дубликации длины 2, так и 2-потомком этой же строки 011, полученным двумя дубликациями длины 1 каждая). Множество всех потомков строки  $\mathbf{x}$  обозначим через  $D^*(\mathbf{x}) := \bigcup_{t \geq 0} D^t(\mathbf{x})$ , где

$D^0(\mathbf{x}) := \{\mathbf{x}\}$ . В этих обозначениях  $D^*(\mathbf{x})$  для заданной входной строки  $\mathbf{x}$  – множество всех возможных выходов канала с тандемными ( $\leq 3$ )-дубликациями.

**1.2. Различимые строки и безошибочная передача.** Две строки  $\mathbf{x}, \mathbf{y} \in \mathcal{A}_q^*$  называются неразличимыми (confusable) в заданном канале связи, если при их передаче по этому каналу на его выходе может появиться одна и та же строка; в противном случае они называются различимыми (non-confusable). В нашей терминологии  $\mathbf{x}$  и  $\mathbf{y}$  неразличимы, если у них имеется общий потомок, т.е.  $D^*(\mathbf{x}) \cap D^*(\mathbf{y}) \neq \emptyset$ . Множество строк  $\mathcal{C} \subseteq \mathcal{A}_q^*$  называется кодом с нулевой ошибкой [9] для заданного канала, если любые два несовпадающих кодовых слова  $\mathbf{x}, \mathbf{y} \in \mathcal{C}$  различимы. Заметим, что код с нулевой ошибкой исправляет все конфигурации ошибок, реализуемые в канале, т.е. приемник может однозначно декодировать любую заданную строку на выходе канала, поскольку в  $\mathcal{C}$  имеется лишь одно кодовое слово, способное породить эту строку. Код с нулевой ошибкой  $\mathcal{C} \subseteq \mathcal{A}_q^n$  называется оптимальным, если не существует никакого другого кода с нулевой ошибкой  $\mathcal{C}' \subseteq \mathcal{A}_q^n$ , такого что  $|\mathcal{C}'| > |\mathcal{C}|$ .

Скорость кода  $\mathcal{C} \subseteq \mathcal{A}_q^n$ , выражаемая в битах на символ, определяется как  $\frac{1}{n} \log_2 |\mathcal{C}|$ . Пропускная способность с нулевой ошибкой для канала с алфавитом  $\mathcal{A}_q$  на входе определяется как  $\limsup_{n \rightarrow \infty}$  по всем скоростям оптимальных кодов с нулевой ошибкой в  $\mathcal{A}_q^n$ . Эта величина представляет собой максимальное число бит на символ, которые можно безошибочно передать по заданному каналу.

## § 2. Корневые и неприводимые строки относительно дубликаций

Последовательно применяя операцию дедупликации, т.е. удаления дублицированных подстрок длины  $\leq 3$ , каждую строку  $\mathbf{x}$  можно привести к ее *корневой* строке  $R(\mathbf{x})$ , не содержащей повторяющихся подряд подстрок длины  $\leq 3$ . Более того, как показано в [1, теорема 24], корневые строки единственны в том смысле, что независимо от порядка, в котором выполняются дедупликации, процесс гарантированно приведет к одной и той же строке. (Подчеркнем, что это “свойство единственности корневых строк” выполнено лишь для моделей с тандемными дубликациями длины (i)  $= \ell$ , (ii)  $\leq 2$  или (iii)  $\leq 3$ . Оно не выполняется, например, в моделях с тандемными дубликациями длины  $\leq \ell$  при  $\ell \in \{4, 5, \dots\}$ ; см. [1, теорема 40].)

В этом контексте строка, не содержащая повторяющихся подстрок длины  $\leq 3$ , называется *неприводимой*. Другими словами, строка неприводима<sup>3</sup>, если она не содержит подстрок вида  $aa$ ,  $abab$  и  $abcabc$ , где  $a, b, c \in \mathcal{A}_q$ . Через  $\text{Irr}_q$  обозначим множество всех неприводимых строк над  $\mathcal{A}_q$ , через  $\text{Irr}_q(n)$  – множество всех неприводимых строк длины  $n$ , а через  $I_q(n)$  – мощность последнего, т.е.  $I_q(n) := |\text{Irr}_q(n)|$ . Из вышеупомянутого “свойства единственности корневых строк” следует, что любые две различные неприводимые строки  $\mathbf{x}, \mathbf{y} \in \text{Irr}_q$  различимы в канале с тандемными

<sup>3</sup> Неприводимые строки являются частным случаем строк с запрещенными конфигурациями, или строк с ограничениями [10], где множество запрещенных конфигураций имеет вид  $\{aa, abab, abcabc : a, b, c \in \mathcal{A}_q\}$ .

( $\leq 3$ )-дубликациями, т.е.  $D^*(\mathbf{x}) \cap D^*(\mathbf{y}) = \emptyset$ , и поэтому множество  $\text{Irr}_q(n)$  является кодом с нулевой ошибкой для этого канала [1, теорема 27].

Всюду далее будем предполагать, что  $q \geq 3$ , поскольку в случае двоичного алфавита рассматриваемая задача становится тривиальной. Например, над двоичным алфавитом существует лишь конечное число неприводимых строк, а именно  $\text{Irr}_2 = \{0, 1, 01, 10, 010, 101\}$ , так что пропускная способность с нулевой ошибкой для канала с тандемными ( $\leq 3$ )-дубликациями над двоичным алфавитом равна нулю.

Для нас представляет интерес асимптотическое поведение величины  $I_q(n)$  при  $n \rightarrow \infty$ , и в частности, экспонента ее скорости роста

$$\iota_q := \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 I_q(n). \quad (2.1)$$

Величину  $\iota_q$  можно охарактеризовать с помощью стандартных методов теории систем с ограничениями [10], например, как логарифм наибольшего собственного значения матрицы смежности направленного графа, представляющего собой диаграмму состояний системы, порождающей неприводимые строки. Здесь мы применим более простую характеристику из [5, предложение 2], где было показано, что для  $I_q(n)$  выполнено рекуррентное соотношение

$$I_q(n) = (q-2)I_q(n-1) + (q-3)I_q(n-2) + (q-2)I_q(n-3)$$

и поэтому

$$\iota_q = \log_2 r, \quad (2.2a)$$

где  $r$  – единственный положительный вещественный корень многочлена

$$x^3 - (q-2)x^2 - (q-3)x - (q-2),$$

т.е.  $r$  задается неявным условием

$$r^3 - (q-2)r^2 - (q-3)r - (q-2) = 0, \quad r > 0. \quad (2.2b)$$

В следующей лемме дается еще одна характеристика величины  $\iota_q$  для троичного алфавита ( $q = 3$ ), а также вытекающая из нее граница снизу на  $\iota_q$  для больших алфавитов, которая будет использована в доказательстве нашего основного результата (теорема 1).

*Лемма 1. Для любых  $q \geq 3$  и  $\beta \in [0, 1]$  справедливо неравенство*

$$\iota_q \geq \frac{H(\beta)}{1+2\beta}, \quad (2.3)$$

где  $H(\beta) := -\beta \log_2 \beta - (1-\beta) \log_2 (1-\beta)$  – функция двоичной энтропии. Равенство в (2.3) достигается тогда и только тогда, когда  $q = 3$  и  $\beta = \bar{\beta}$ , где  $\bar{\beta}$  – единственное положительное решение уравнения  $(1-x)^3 = x$ .

*Доказательство.* Докажем соотношение

$$\iota_3 = \max_{0 \leq \beta \leq 1} \frac{H(\beta)}{1+2\beta}, \quad (2.4)$$

из которого будет немедленно следовать утверждение леммы (так как  $\iota_q$  монотонно возрастает по  $q$ ). Приравняв производную функции  $\frac{H(\beta)}{1+2\beta}$  к нулю, убеждаемся, что эта функция достигает максимума в единственной положительной вещественной точке, удовлетворяющей уравнению  $(1-x)^3 = x$ , назовем ее  $\bar{\beta}$ . Тогда правую часть

равенства (2.4) можно представить в виде

$$\frac{H(\bar{\beta})}{1+2\bar{\beta}} = \log_2 \left( \bar{\beta}^{\frac{-\bar{\beta}}{1+2\bar{\beta}}} (1-\bar{\beta})^{\frac{-1+\bar{\beta}}{1+2\bar{\beta}}} \right) = -\log_2(1-\bar{\beta}). \quad (2.5)$$

С другой стороны, мы знаем, что  $t_3 = \log_2 r$ , где  $r$  – единственный положительный вещественный корень уравнения  $x^3 - x^2 - 1 = 0$  (см. (2.2)). Поэтому доказательство равенства в (2.4) равносильно доказательству того, что  $-\log_2(1-\bar{\beta}) = \log_2 r$ , т.е. что  $(1-\bar{\beta})^{-1}$  является решением уравнения  $x^3 - x^2 - 1 = 0$ . В этом можно убедиться непосредственно, подставляя  $(1-\bar{\beta})^{-1}$  вместо  $x$  и используя тот факт, что  $(1-\bar{\beta})^3 = \bar{\beta}$ . ▲

### § 3. Различимость строк в канале с тандемными ( $\leq 3$ )-дупликациями

В этом параграфе изложим несколько фактов об эволюции строк под действием тандемных дупликаций длины  $\leq 3$ , основным из которых будет вывод верхней границы на максимальное число попарно различных строк в заданном конусе потомков  $D^*(\mathbf{x})$  (предложение 1). По поводу дальнейшего изучения комбинаторных и алгоритмических аспектов (не)различимости в ( $\leq 2$ )- и ( $\leq 3$ )-каналах с тандемными дупликациями отсылаем читателя к работе [11].

В следующей лемме утверждается, что множество попарно различных строк, каждая из которых является 1-потомком данной строки  $\mathbf{x}$ , не может состоять из более чем двух элементов. В доказательстве также иллюстрируются условия, при которых две различные строки могут быть получены применением различных мутаций к одной и той же строке  $\mathbf{x}$  (см. уравнения (3.1) ниже).

*Лемма 2. Рассмотрим произвольную строку  $\mathbf{x}$  и множество  $D^1(\mathbf{x})$  ее 1-потомков, и пусть  $\mathcal{C} \subseteq D^1(\mathbf{x})$  – код с нулевой ошибкой для канала с тандемными ( $\leq 3$ )-дупликациями. Тогда  $|\mathcal{C}| \leq 2$ .*

*Доказательство.* Рассмотрим  $\mathbf{x}', \mathbf{x}'' \in D^1(\mathbf{x})$  и предположим, что мутации, переводящие  $\mathbf{x}$  в  $\mathbf{x}'$  и  $\mathbf{x}''$ , применяются к различным и непересекающимся подстрокам строки  $\mathbf{x}$ . Тогда  $\mathbf{x}'$  и  $\mathbf{x}''$  неразличимы, поскольку у них имеется общий потомок; действительно, достаточно применить к  $\mathbf{x}'$  дупликацию, породившую  $\mathbf{x}''$  из  $\mathbf{x}$ , и наоборот. Теперь предположим, что дупликации, переводящие  $\mathbf{x}$  в  $\mathbf{x}'$  и  $\mathbf{x}''$ , применяются к пересекающимся подстрокам  $\mathbf{x}$ . Оказывается, что во всех возможных случаях, *кроме одного*, то же самое рассуждение, что и для непересекающихся подстрок, показывает, что  $\mathbf{x}'$  и  $\mathbf{x}''$  неразличимы (мы продемонстрируем это для случаев, когда пересечение находится с правой стороны более длинной подстроки, так как остальные случаи симметричны им):

- (i) Для случая пересекающихся подстрок длины 1 и 2 рассмотрим строку  $\mathbf{x} = \mathbf{u}abv$  и заметим, что ее потомки  $\mathbf{x}' = \mathbf{u}\overline{ab}abv$  и  $\mathbf{x}'' = \mathbf{u}a\overline{b}bv$  неразличимы, поскольку у них имеется общий потомок  $\mathbf{u}ababbbv$ ;
- (ii) Для случая пересекающихся подстрок длины 2 и 2 рассмотрим строку  $\mathbf{x} = \mathbf{u}abcv$  и заметим, что ее потомки  $\mathbf{x}' = \mathbf{u}\overline{ab}abcv$  и  $\mathbf{x}'' = \mathbf{u}a\overline{b}cbcv$  неразличимы, поскольку у них имеется общий потомок  $\mathbf{u}ababcbcv$ ;
- (iii) Для случая пересекающихся подстрок длины 2 и 3, имеющих пересечение длины 1, рассмотрим строку  $\mathbf{x} = \mathbf{u}abcdv$  и заметим, что ее потомки  $\mathbf{x}' = \mathbf{u}\overline{abc}abc dv$  и  $\mathbf{x}'' = \mathbf{u}abc\overline{d}cdv$  неразличимы, поскольку у них имеется общий потомок  $\mathbf{u}abcabc dcdv$ ;
- (iv) Для случая пересекающихся подстрок длины 2 и 3, имеющих пересечение длины 2, рассмотрим строку  $\mathbf{x} = \mathbf{u}abcv$  и заметим, что ее потомки  $\mathbf{x}' = \mathbf{u}\overline{abc}abcv$  и  $\mathbf{x}'' = \mathbf{u}abc\overline{bc}cv$  неразличимы, поскольку у них имеется общий потомок  $\mathbf{u}abcabc bcv$ ;

- (v) Для случая пересекающихся подстрок длины 3 и 3, имеющих пересечение длины 1, рассмотрим строку  $x = \underline{u}abc\underline{d}ev$  и заметим, что ее потомки  $x' = \underline{u}abc\underline{a}bc\underline{d}ev$  и  $x'' = \underline{u}abc\underline{d}ec\underline{d}ev$  неразличимы, поскольку у них имеется общий потомок  $\underline{u}abcabc\underline{d}ec\underline{d}ev$ ;
- (vi) Для случая пересекающихся подстрок длины 3 и 3, имеющих пересечение длины 2, рассмотрим строку  $x = \underline{u}abc\underline{d}v$  и заметим, что ее потомки  $x' = \underline{u}abc\underline{a}bc\underline{d}v$  и  $x'' = \underline{u}abc\underline{d}b\underline{c}d\underline{v}$  неразличимы, поскольку у них имеется общий потомок  $\underline{u}abcabc\underline{d}b\underline{c}d\underline{v}$ ;
- (vii) Для случая пересекающихся подстрок длины 1 и 3 рассмотрим строку  $x = \underline{u}abc\underline{v}$  и заметим, что ее потомки  $x' = \underline{u}abc\underline{a}bc\underline{v}$  и  $x'' = \underline{u}abc\underline{c}v$  неразличимы, поскольку у них имеется общий потомок  $\underline{u}abcabc\underline{c}v$ .

Единственный случай, не вошедший в этот список, – это случай пересекающихся подстрок длины 1 и 3, когда пересечение возникает в середине более длинной подстроки. А именно, для  $x = \underline{u}abc\underline{v}$ , где символы  $a, b, c \in \mathcal{A}_q$  различны, положим

$$x' = \underline{u}abc\underline{a}bc\underline{v}, \quad (3.1a)$$

$$x'' = \underline{u}ab\underline{b}c\underline{v}. \quad (3.1b)$$

В этом случае мы не можем применить такое же рассуждение, как и выше, чтобы показать, что строки  $x'$  и  $x''$  неразличимы, и действительно, в общем случае это не обязательно верно. Например, если обе строки  $u$  и  $v$  – пустые, то  $x'$  и  $x''$  в (3.1) различимы, поскольку символ  $a$  не может оказаться после символа  $c$  в потомках строки  $x''$ , в то время как во всех потомках слова  $x'$  символ  $a$  находится после  $c$  (аналогичный пример был приведен в [1]). Так происходит, поскольку фрагмент  $abc$ , содержащийся в исходной строке  $x$ , больше не возникнет в строке  $x''$ , так как она “разбита” вставкой еще одной копии  $b$ . Наконец, в строке  $x''$  (соответственно,  $x'$ ) можно повторить дубликацию, породившую  $x'$  (соответственно,  $x''$ ) из  $x$ , и таким образом показать, что  $x'$  и  $x''$  неразличимы, во всех случаях, кроме (3.1). Таким образом, код с нулевой ошибкой в  $D^1(x)$  может содержать не более двух кодовых слов. ▲

*Замечание 1.* Не каждый случай вида (3.1) приводит к различимым потомкам. В качестве контрпримера предположим, что  $u$  – пустая строка, а  $v = a$ , так что  $x = abca$ ,  $x' = \underline{a}bc\underline{a}bc\underline{a}$  и  $x'' = \underline{a}b\underline{b}ca$ . Тогда у  $x'$  и  $x''$  есть общий потомок  $abbcabca$ , и поэтому они неразличимы. Однако для наших целей достаточно того факта, что (3.1) является *единственным* случаем, когда два потомка *могут быть* различимыми. В частности, это факт позволит нам вывести точную *верхнюю* границу на мощность оптимальных кодов с нулевой ошибкой.

Приведенное выше наблюдение справедливо в общем случае, а не только для 1-потомков строки  $x$ . А именно, если  $x', x'' \in D^*(x)$  получены применением к  $x$  двух разных конфигураций дубликаций, в каждой из этих строк можно повторить/воспроизвести дубликации, примененные ко второй, и таким образом показать, что у них имеется общий потомок. Единственный случай, когда такой процесс повторения дубликаций *может* в какой-то момент стать невозможным, это ситуация, когда к  $x'$  применяется дубликация длины 3, а к соответствующему фрагменту в  $x''$  – дубликация длины 1 (см. (3.1)), так что повторить в  $x''$  соответствующую мутацию, примененную к  $x'$ , невозможно. Так происходит из-за того, что всякий раз, когда к строке применяется дубликация длины 2 или 3, *все* фрагменты длины  $\leq 3$  исходной строки сохраняются в получившейся строке (с несколькими дополнительными подстроками, появляющимися в том месте, куда вставлялась копия). Единственный случай, в котором фрагмент длины 3 исходной строки исчезает (не появляется в получившейся строке), – это после дубликации длины 1, как показано в (3.1b). На основе этого наблюдения мы выведем верхнюю границу на мощность оп-

тимальных кодов с нулевой ошибкой в множестве всех  $t$ -потомков данной строки  $\mathbf{x}$  при любом  $t$  (предложение 1 ниже).

Пример 2. Чтобы пояснить вышесказанное, представим пример, приведенный в (1.1), в несколько ином виде (фрагмент  $\mathbf{120}$  выделен, и показаны дубликации слева (соответственно, справа) от этого фрагмента, при которых копии вставляются слева (соответственно, справа) от оригинала):

$$\mathbf{x} = 0 \mathbf{1} \mathbf{120} \mathbf{2} \mathbf{1} \mathbf{0}, \quad (3.2a)$$

$$\mathbf{x}^{(1)} = \underline{0} \bar{0} \mathbf{1} \mathbf{120} \mathbf{2} \mathbf{1} \mathbf{0}, \quad (3.2b)$$

$$\mathbf{x}^{(2)} = 0 \mathbf{0} \mathbf{1} \mathbf{1} \overline{\mathbf{202}} \underline{\mathbf{202}} \mathbf{1} \mathbf{0}, \quad (3.2c)$$

$$\mathbf{x}^{(3)} = 0 \mathbf{0} \mathbf{1} \mathbf{120} \mathbf{2} \mathbf{2} \mathbf{0} \overline{\mathbf{21}} \underline{\mathbf{21}} \mathbf{0}, \quad (3.2d)$$

$$\mathbf{x}^{(4)} = 0 \mathbf{0} \underline{\mathbf{11}} \overline{\mathbf{11}} \mathbf{202} \mathbf{2} \mathbf{0} \mathbf{2} \mathbf{1} \mathbf{2} \mathbf{1} \mathbf{0}. \quad (3.2e)$$

Пусть  $\mathbf{z} = 0 \mathbf{1} \mathbf{1220210}$ . Заметим, что в  $\mathbf{z}$  можно воспроизвести все дубликации подстрок  $\mathbf{x}$ , которые либо не пересекаются с сегментом  $\mathbf{120}$ , либо пересекаются с ним лишь частично<sup>4</sup>, как в примерах, приведенных в (3.2):

$$\mathbf{x} = 0 \mathbf{1} \mathbf{120} \mathbf{2} \mathbf{1} \mathbf{0}, \quad (3.3a)$$

$$\mathbf{z} = 0 \mathbf{1} \mathbf{1} \overline{\mathbf{2}} \underline{\mathbf{2}} \mathbf{0} \mathbf{2} \mathbf{1} \mathbf{0}, \quad (3.3b)$$

$$\mathbf{z}^{(1)} = \underline{0} \bar{0} \mathbf{1} \mathbf{1220210}, \quad (3.3c)$$

$$\mathbf{z}^{(2)} = 0 \mathbf{0} \mathbf{1} \mathbf{12} \overline{\mathbf{202}} \underline{\mathbf{202}} \mathbf{1} \mathbf{0}, \quad (3.3d)$$

$$\mathbf{z}^{(3)} = 0 \mathbf{0} \mathbf{1} \mathbf{1220220} \overline{\mathbf{21}} \underline{\mathbf{21}} \mathbf{0}, \quad (3.3e)$$

$$\mathbf{z}^{(4)} = 0 \mathbf{0} \underline{\mathbf{11}} \overline{\mathbf{11}} \mathbf{22022021210}. \quad (3.3f)$$

Поэтому любая пара строк из (3.2) и (3.3) неразличима; например, общим потомком  $\mathbf{z}$  и  $\mathbf{x}^{(3)}$  является  $\mathbf{z}^{(3)}$ . Единственная мутация, которую нельзя повторить в  $\mathbf{z}$ , – это дубликация всего фрагмента  $\mathbf{120}$ , поскольку соответствующий фрагмент в  $\mathbf{z}$  больше не существует (он был “разбит” вставленным символом 2). Например, если бы нужно было превратить  $\mathbf{x}^{(2)}$  из (3.2c) в строку

$$\mathbf{y} = 0 \mathbf{0} \mathbf{1} \overline{\mathbf{120}} \underline{\mathbf{120}} \mathbf{2} \mathbf{2} \mathbf{0} \mathbf{2} \mathbf{1} \mathbf{0} \quad (3.4)$$

вместо  $\mathbf{x}^{(3)}$ , было бы уже невозможно применить тот же процесс, что и в (3.3).

Прежде чем сформулировать предложение 1, являющееся основным результатом этого параграфа, докажем одну полезную лемму.

Лемма 3. *Зафиксируем натуральные числа  $b, t, n$ , такие что  $b \leq t \leq n$ . Пусть  $\mathcal{U} \subseteq \{l_1, l_3, *\}^n$  – множество строк, удовлетворяющих следующим двум условиям:*

- (1) *Каждая строка в  $\mathcal{U}$  имеет ровно  $b$  символов  $l_3$ ,  $t - b$  символов  $l_1$ , и  $n - t$  символов  $*$ ;*
- (2) *Для любых двух различных строк  $\mathbf{u}, \mathbf{v} \in \mathcal{U}$  найдется позиция  $i \in \{1, \dots, n\}$ , в которой  $\{u_i, v_i\} = \{l_1, l_3\}$  (т.е. такая, что либо  $u_i = l_1$  и  $v_i = l_3$ , либо  $u_i = l_3$  и  $v_i = l_1$ ).*

Тогда  $|\mathcal{U}| \leq 2^{tH(b/t)}$ .

Доказательство. Рассмотрим  $n$  бросков монеты, для которой вероятность выпадения орла равна  $b/t$ , и введем следующие события, индексированные строками из множества  $\mathcal{U}$ . Для  $\mathbf{u} \in \mathcal{U}$  обозначим через  $A_{\mathbf{u}}$  событие, состоящее в том, что на  $i$ -м

<sup>4</sup> Для  $\mathbf{x}$  из (3.2a) подстроки длины  $\leq 3$ , частично пересекающиеся с подстрокой  $\mathbf{120}$ , следующие:  $\mathbf{1}, \mathbf{2}, \mathbf{0}, \mathbf{11}, \mathbf{12}, \mathbf{20}, \mathbf{02}, \mathbf{011}, \mathbf{112}, \mathbf{202}, \mathbf{021}$ .

броске монеты выпал орел, если  $u_i = l_3$ , выпала решка, если  $u_i = l_1$ , и неважно, что выпало, если  $u_i = *$ , для всех  $i = 1, \dots, n$ . Из условия (1) следует, что вероятность такого события равна

$$\Pr\{A_u\} = \left(\frac{b}{t}\right)^b \left(1 - \frac{b}{t}\right)^{t-b} = 2^{-tH(b/t)}$$

для любой строки  $u \in \mathcal{U}$ . При этом в силу условия (2) события  $A_u$  и  $A_v$  несовместны для любых  $u, v \in \mathcal{U}$ ,  $u \neq v$ . Отсюда  $|\mathcal{U}| \Pr\{A_u\} \leq 1$ , что и требовалось доказать.  $\blacktriangle$

Заметим, что приведенная верхняя граница на  $|\mathcal{U}|$  не зависит от  $n$  (т.е. от числа символов  $*$ ).

**Предложение 1.** *Рассмотрим строку  $x \in A_q^n$ , и пусть  $C \subseteq D^t(x)$  – код с нулевой ошибкой для канала с тандемными ( $\leq 3$ )-дубликациями, удовлетворяющий следующему условию: из  $t$  дубликаций, порождающих каждое кодовое слово  $y \in C$  из слова  $x$ , ровно  $b$  имеют длину 3. Тогда  $|C| \leq 2^{tH(b/t)}$ .*

**Доказательство.** Как показано выше, если две строки  $x', x'' \in D^t(x)$  различимы, то это с необходимостью означает, что к фрагменту  $abc$  в одном из предков строки  $x'$  была применена дубликация длины 3, а к среднему символу соответствующего фрагмента в некотором предке  $x''$  была применена дубликация длины 1, или наоборот. Другими словами, для любой пары кодовых слов кода с нулевой ошибкой в  $D^t(x)$  найдется фрагмент, в котором они отличаются на мутацию длины 1/длины 3. Поэтому интересующий нас вопрос состоит в том, насколько велико может быть множество строк, любые две из которых отличаются на мутацию длины 1/длины 3 в некоторой позиции. (Неважно, что это за позиции и что происходит в промежутках между ними, единственное требование состоит в том, что любая пара кодовых слов в каком-то месте отличается на мутацию длины 1/длины 3, так как это единственная возможность для того, чтобы две строки могли стать различимыми.) Поэтому  $|C|$  можно ограничить сверху максимальной мощностью множества  $\mathcal{U}$  строк над “алфавитом” {длина 1, длина 3,  $*$ }, удовлетворяющего следующим условиям: 1) любая строка в  $\mathcal{U}$  содержит ровно  $b$  символов “длина 3” и  $t - b$  символов “длина 1”; 2) любые две различные строки в  $\mathcal{U}$  отличаются в некоторой позиции на символ “длина 1”/“длина 3”. (Формальный символ  $*$  служит для заполнения пустых позиций, которые могут возникать из-за того, что различные пары кодовых слов могут отличаться на мутацию длины 1/длины 3 в различных фрагментах; см. также пример 3 ниже.) Теперь требуемая граница получается применением леммы 3.  $\blacktriangle$

**Пример 3.** Приведем пример множества строк  $\mathcal{U}$  из вышеизложенного доказательства. Рассмотрим следующую строку:

$$x = \overbrace{0123} \overbrace{4567} \overbrace{89}. \quad (3.5)$$

(Все символы в  $x$  обозначены по-разному для облегчения понимания, это никак не влияет на общность рассуждений.) Пусть применением тандемных дубликаций к четырем фрагментам в  $x$ , отмеченных скобками сверху или снизу, получены следующие потомки в  $D^3(x)$ :

$$x' = 012 \underline{012} \underline{2}3456 \underline{6}789, \quad (3.6a)$$

$$x'' = 01 \underline{1}23456 \underline{6}789 \underline{789}, \quad (3.6b)$$

$$x''' = 01 \underline{1}23456 \underline{7567}889, \quad (3.6c)$$

где подчеркнуты вставленные копии. Мутации, примененные к этим четырем фрагментам, можно описать с помощью следующих строк:

$$u' = l_3 l_1 l_1 *, \quad (3.7a)$$



$$\mathbf{u}'' = l_1 * l_1 l_3, \quad (3.7b)$$

$$\mathbf{u}''' = l_1 * l_3 l_1, \quad (3.7c)$$

где символ  $l_3$  показывает, что к соответствующему фрагменту применена дупликация длины 3, символ  $l_1$  показывает, что к среднему символу этого фрагмента применена дупликация длины 1, а  $*$  означает, что к этому фрагменту ни одна из этих мутаций не применялась. Заметим, что для любой пары строк в (3.7) найдется координата, в которой одна из них равна  $l_1$ , а другая  $l_3$ .

#### § 4. Пропускная способность с нулевой ошибкой для канала с тандемными ( $\leq 3$ )-дупликациями

Пусть  $\mathcal{C}_q^*(n) \subseteq \mathcal{A}_q^n$  – оптимальный код с нулевой ошибкой для канала с тандемными ( $\leq 3$ )-дупликациями. Для заданной неприводимой строки  $\mathbf{x} \in \text{Irr}_q$  положим

$$\mathcal{C}_q^*(n; \mathbf{x}) := \mathcal{C}_q^*(n) \cap D^*(\mathbf{x}).$$

Тогда  $\mathcal{C}_q^*(n; \mathbf{x})$  является оптимальным кодом с нулевой ошибкой на множестве всех потомков  $\mathbf{x}$  длины  $n$ . Это так, поскольку в силу свойства единственности корневой строки для тандемных дупликаций длины  $\leq 3$  любые два различных конуса потомков не пересекаются [1, следствие 26], и поэтому любые две строки, имеющие различные корневые строки, различимы. Иными словами, можно, не ограничивая общности, строить код по отдельности в конусе потомков для каждой возможной корневой/неприводимой строки. Этот факт явно сформулирован в следующей лемме; ее доказательство опущено, поскольку оно непосредственно вытекает из [1, следствие 26].

*Лемма 4. Оптимальный код длины  $n$  с нулевой ошибкой можно представить в виде несвязного объединения оптимальных кодов в каждом из конусов потомков:*

$$\mathcal{C}_q^*(n) = \bigcup_{\mathbf{x} \in \text{Irr}_q} \mathcal{C}_q^*(n; \mathbf{x}). \quad (4.1)$$

Следующее утверждение дает характеризацию экспоненты скорости роста мощности оптимальных кодов  $\mathcal{C}_q^*(n)$  или, эквивалентным образом, пропускной способности с нулевой ошибкой канала с тандемными ( $\leq 3$ )-дупликациями. В нем утверждается, что эта величина равна

$$\iota_q = \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 I_q(n)$$

(см. (2.1) и (2.2)). Другими словами, пропускная способность с нулевой ошибкой достигается на кодах  $\text{Irr}_q(n)$ , состоящих из неприводимых строк длины  $n$ .

**Теорема 1.** *Пропускная способность с нулевой ошибкой канала с тандемными ( $\leq 3$ )-дупликациями с алфавитом  $\mathcal{A}_q$ ,  $q \geq 3$ , равна  $\iota_q$ .*

*Доказательство.* Требуется показать, что

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_q^*(n)| = \iota_q. \quad (4.2)$$

Так как  $\text{Irr}_q(n) \subseteq \mathcal{C}_q^*(n)$ , то мы знаем, что

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_q^*(n)| \geq \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 I_q(n) = \iota_q$$

(см. (2.1)), поэтому достаточно доказать противоположное неравенство

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_q^*(n)| \leq \iota_q.$$

Для этого мы упростим анализ, построив достаточно большой *подкод*  $\mathcal{C}_q(n; m, t, b) \subseteq \mathcal{C}_q^*(n)$ , имеющий ту же экспоненту скорости роста, что и оптимальный код  $\mathcal{C}_q^*(n)$ , т.е.

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_q^*(n)| = \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_q(n; m, t, b)|,$$

при подходящем выборе параметров  $m, t, b$ .

Зафиксируем произвольную неприводимую строку  $\mathbf{x}$  длины  $m$ , т.е.  $\mathbf{x} \in \text{Irr}_q(m)$ , и пусть  $\mathcal{C}_q(n; \mathbf{x}, t, b) \subseteq \mathcal{C}_q^*(n; \mathbf{x})$  – код, содержащий только те кодовые слова из  $\mathcal{C}_q^*(n; \mathbf{x})$ , которые удовлетворяют следующим двум условиям: 1) каждое кодовое слово лежит в  $D^t(\mathbf{x})$ , т.е. является  $t$ -потомком строки  $\mathbf{x}$ ; 2) из  $t$  дубликаций, переводящих  $\mathbf{x}$  в заданный потомок/заданное кодовое слово, ровно  $b$  имеют длину 3. Определим такой подкод следующим образом:

$$\mathcal{C}_q(n; m, t, b) := \bigcup_{\mathbf{x} \in \text{Irr}_q(m)} \mathcal{C}_q(n; \mathbf{x}, t, b). \quad (4.3)$$

Из этой конструкции и леммы 4 следует, что

$$\mathcal{C}_q^*(n) = \bigcup_{m, t, b} \mathcal{C}_q(n; m, t, b). \quad (4.4)$$

Теперь должно быть ясно, что величина  $|\mathcal{C}_q(n; m, t, b)|$ , максимизированная по всем возможным значениям  $m, t, b$ , имеет ту же самую экспоненту скорости роста, что и  $|\mathcal{C}_q^*(n)|$  (значения  $m, t, b$  выбираются для каждого  $n$ , т.е. оптимальные значения параметров  $m, t, b$  являются, вообще говоря, функциями от длины блока  $n$ ). Это следует из равенства (4.4) и принципа Дирихле – мощность кода  $\mathcal{C}_q^*(n)$  растет экспоненциально быстро по длине блока  $n$ , а выбрать значения каждого из параметров  $m, t$  и  $b$  можно линейным количеством способов, поэтому по крайней мере для одного такого выбора код  $\mathcal{C}_q(n; m, t, b)$  будет содержать экспоненциально много кодовых слов (с тем же показателем экспоненты). Таким образом, коды  $\mathcal{C}_q(n; m, t, b)$  асимптотически оптимальны по скорости, т.е. на них достигается пропускная способность с нулевой ошибкой канала с тандемными ( $\leq 3$ )-дубликациями, когда параметры  $m, t, b$  выбираются правильным образом (так, чтобы максимизировать  $|\mathcal{C}_q(n; m, t, b)|$ ).

Теперь вычислим скорость построенных кодов. Согласно (4.3) и предложению 1 (в котором утверждается, что  $|\mathcal{C}_q(n; \mathbf{x}, t, b)| \leq 2^{tH(b/t)}$ ), мощность кода  $\mathcal{C}_q(n; m, t, b)$  можно оценить сверху как

$$|\mathcal{C}_q(n; m, t, b)| \leq I_q(m) \cdot 2^{tH(b/t)}, \quad (4.5)$$

а длину этого кода можно оценить снизу как

$$n \geq m + 3b + (t - b) = m + t + 2b \quad (4.6)$$

(начальная неприводимая строка имеет длину  $m$ , и ровно  $b$  дубликаций, порождающих ее потомков, имеют длину 3). Следовательно,

$$\frac{1}{n} \log_2 |\mathcal{C}_q(n; m, t, b)| \leq \frac{\log_2 I_q(m) + tH(b/t)}{m + t + 2b}. \quad (4.7)$$

Чтобы найти асимптотику этой величины при  $n \rightarrow \infty$ , нужно рассмотреть два случая, соответствующие разным выборам параметров  $m, t, b$ :

- $m = o(t)$ . Пусть  $\lim_{t \rightarrow \infty} \frac{b}{t} = \beta \in [0, 1]$ . Тогда

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_2 |C_q(n; m, t, b)| \leq \frac{H(\beta)}{1 + 2\beta} \leq \iota_q, \quad (4.8)$$

где первое неравенство следует из (4.7), а второе совпадает с (2.3);

- $t = \mathcal{O}(m)$ . Пусть  $\liminf_{m \rightarrow \infty} \frac{t}{m} = \tau \geq 0$  и  $\lim_{t \rightarrow \infty} \frac{b}{t} = \beta \in [0, 1]$ . Тогда

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_2 |C_q(n; m, t, b)| \leq \frac{\iota_q + \tau H(\beta)}{1 + \tau(1 + 2\beta)} \leq \iota_q. \quad (4.9)$$

Здесь снова первое неравенство вытекает из (4.7), а второе равносильно (2.3).

Итак, все выборы значений параметров  $m, t, b$  дают асимптотическую скорость кодов  $C_q(n; m, t, b)$ , не превосходящую  $\iota_q$ . Поскольку все эти коды оптимальны по скорости, как указано во втором абзаце этого доказательства, равенство (4.2) тем самым доказано. ▲

## § 5. Заключение

Эволюция строк под действием тандемных дупликаций – интересная и нетривиальная задача, важная в нескольких областях исследований. В настоящей статье исследована различимость строк при тандемных дупликациях переменной длины – задача, инспирированная исправлением ошибок в каналах связи, в которых передаваемые сообщения подвергаются мутациям такого типа. А именно, для случая дупликаций длины  $\leq 3$  получена верхняя граница на максимальную мощность множества попарно различных строк, которая вместе с конструкцией из [1] определяет максимальную скорость, достижимую кодами, исправляющими произвольное число таких дефектов.

В случаях, когда корневые строки относительно дупликаций не единственны, например, для модели ( $\leq \ell$ )-тандемных дупликаций с параметром  $\ell$ , большим чем 3, максимальные достижимые скорости остаются неизвестными. Благодаря свойству “неединственности корневых строк” анализ эволюции и различимости строк в таких моделях более сложен, и поэтому для решения задачи о пропускной способности с нулевой ошибкой и смежных с ней задач потребуются дальнейшие исследования и, возможно, другие методы.

## СПИСОК ЛИТЕРАТУРЫ

1. Jain S., Farnoud F., Schwartz M., Bruck J. Duplication-Correcting Codes for Data Storage in the DNA of Living Organisms // IEEE Trans. Inform. Theory. 2017. V. 63. № 8. P. 4996–5010. <https://doi.org/10.1109/TIT.2017.2688361>
2. Kovačević M., Tan V.Y.F. Asymptotically Optimal Codes Correcting Fixed-Length Duplication Errors in DNA Storage Systems // IEEE Commun. Lett. 2018. V. 22. № 11. P. 2194–2197. <https://doi.org/10.1109/LCOMM.2018.2868666>
3. Lenz A., Jünger N., Wachter-Zeh A. Bounds and Constructions for Multi-Symbol Duplication Error Correcting Codes, <https://arXiv.org/abs/1807.02874v3> [cs.IT], 2018.
4. Kovačević M. Zero-Error Capacity of Duplication Channels // IEEE Trans. Commun. 2019. V. 67. № 10. P. 6735–6742. <https://doi.org/10.1109/TCOMM.2019.2931342>
5. Chee Y.M., Chrisnata J., Kiah H.M., Nguyen T.T. Efficient Encoding/Decoding of GC-Balanced Codes Correcting Tandem Duplications // IEEE Trans. Inform. Theory. 2020. V. 66. № 8. P. 4892–4903. <https://doi.org/10.1109/TIT.2020.2981069>

6. *Farnoud F., Schwartz M., Bruck J.* The Capacity of String-Duplication Systems // IEEE Trans. Inform. Theory. 2016. V. 62. № 2. P. 811–824. <https://doi.org/10.1109/TIT.2015.2505735>
7. *Jain S., Farnoud F., Bruck J.* Capacity and Expressiveness of Genomic Tandem Duplication // IEEE Trans. Inform. Theory. 2017. V. 63. № 10. P. 6129–6138. <https://doi.org/10.1109/TIT.2017.2728079>
8. *Leupold P., Martín-Vide C., Mitrana V.* Uniformly Bounded Duplication Languages // Discrete Appl. Math. 2005. V. 146. № 3. P. 301–310. <https://doi.org/10.1016/j.dam.2004.10.003>
9. *Shannon C.E.* The Zero Error Capacity of a Noisy Channel // IRE Trans. Inform. Theory. 1956. V. 2. № 3. P. 8–19. <https://doi.org/10.1109/TIT.1956.1056798>
10. *Marcus B.H., Roth R.M., Siegel P.H.* An Introduction to Coding for Constrained Systems (unpublished manuscript), 5th ed., 2001. Available online at <http://www.math.ubc.ca/~marcus/Handbook/>.
11. *Chee Y.M., Chrisnata J., Kiah H.M., Nguyen T.T.* Deciding the Confusability of Words under Tandem Repeats in Linear Time // ACM Trans. Algorithms. 2019. V. 15. № 3. Art. 42 (22 pp.). <https://doi.org/10.1145/3338514>

*Ковачевич Младен*  
Факультет техничких наук,  
Универзитет г. Нови-Сад, Србија  
[kmladen@uns.ac.rs](mailto:kmladen@uns.ac.rs)

Поступила в редакцию  
04.10.2021  
После доработки  
20.04.2022  
Принята к публикации  
21.04.2022