

МИНОБРНАУКИ РОССИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ПРОФЕССИОНАЛЬНОГО ОБРАЗОВАНИЯ

**«САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ ПЕТРА ВЕЛИКОГО»**

Институт Компьютерных наук и кибербезопасности

Направление 02.03.01 Математика и компьютерные науки

Высшая школа технологий искусственного интеллекта

ОТЧЕТ ПО РАБОТЕ ПО ТЕОРИИ ГРАФОВ:

**Реализация словаря на основе
хэш-таблицы и красно-черного дерева,
кодирование данных**

Вариант 29

Обучающийся: _____

Санько В. В.

Руководитель: _____

Востров А. В.

«_____» _____ 20__ г.

Санкт-Петербург, 2025

Содержание

Введение	4
1 Математическое описание	5
1.1 Хеш-таблица	5
1.2 Красно-черное дерево	7
1.2.1 Перекрашивание узлов	8
1.2.2 Вращения (Rotations)	8
1.2.3 Операция Вставки (Insertion)	8
1.2.4 Операция Удаления (Deletion)	10
1.3 Алгоритм сжатия RLE (Run-Length Encoding)	12
1.4 Математическое и алгоритмическое описание кодирования Фано	12
2 Реализация словаря на основе хеш-таблицы	14
2.1 Структура узла <code>HashNode</code>	14
2.2 Класс <code>HashTable</code>	14
2.3 Хеш-функция	15
2.4 Функция <code>add(const std::string& key, int value)</code>	15
2.5 Функция <code>get(const std::string& key)</code>	16
2.6 Функция <code>remove(const std::string& key)</code>	16
2.7 Функция <code>clear()</code>	17
2.8 Приватная функция <code>rehash()</code>	17
2.9 Публичная функция <code>print</code>	18
2.10 Публичная функция <code>visualize</code>	19
2.11 Класс <code>Dictionary</code>	19
2.12 Функция <code>addWord(const std::string& word_raw)</code>	20
2.13 Функция <code>removeWord(const std::string& word_raw)</code>	21
2.14 Функция <code>findWord(const std::string& word_raw) const</code>	21
2.15 Функция <code>loadFromFile(const std::string& filepath, bool append = false)</code>	22
3 Реализация словаря на основе Красно-Черного дерева	23
3.1 Структура данных: Красно-Черное дерево	23
3.2 Класс <code>RBTree</code>	23
3.3 Конструктор <code>RBTree()</code>	24
3.4 Функция <code>insert(const std::string& key, int value)</code>	25
3.5 Функция <code>search(const std::string& key)</code>	26
3.6 Функция <code>remove(const std::string& key)</code>	26

3.7	Приватные функции <code>leftRotate rightRotate</code>	27
3.8	Приватная функция <code>insertFixup</code>	28
3.9	Приватная функция <code>findNode</code>	30
3.10	Приватная функция <code>transplant</code>	30
3.11	Приватная функция <code>deleteFixup</code>	31
3.12	Приватная функция <code>destroyRecursive</code>	33
3.13	Приватная функция <code>inorderPrintRecursive</code>	33
3.14	Приватная функция <code>getMaxDepth</code>	34
3.15	Класс <code>Dictionary</code> (на основе <code>RBTree</code>)	34
4	Алгоритм сжатия RLE	36
4.1	Функция генерации случайного текста	36
4.2	Функция <code>RLE::advancedRleEncode</code>	38
4.3	Функция <code>RLE::advancedRleDecode</code>	39
5	Реализация кодирования Фано	42
5.1	Структура <code>FanoNode</code>	42
5.2	Функция <code>generateFanoCodes</code>	42
5.3	Функция <code>buildFanoCodes</code>	44
5.4	Функция <code>encodeFano</code>	45
5.5	Функция <code>decodeFano</code>	46
6	Вспомогательные функции обработки текста	47
6.1	Глобальная функция <code>processTextToWords</code>	47
6.2	Глобальная функция <code>readFileToString</code>	48
7	Результаты работы	49
	Заключение	53
	Список литературы	55

Введение

Одной из основных задач при работе с текстами является создание словарей. В рамках данной работы рассматриваются два подхода к построению словаря: на основе хеш-таблицы и на основе самобалансирующегося красно-черного дерева. Также необходимо уделить особое внимание выбору хеш-функции для хеш-таблицы, так как необходимо минимизировать количество коллизий. Помимо этого применить алгоритм сжатия данных (продвинутый RLE и алгоритм Фано) на текстовой информации на русском языке.

1 Математическое описание

1.1 Хеш-таблица

Хеш-таблица представляет собой структуру данных, предназначенную для эффективной реализации интерфейса ассоциативного массива, то есть отображения ключей K на значения V . Пусть имеется множество ключей U (универсум ключей). Хеш-таблица использует массив T размера m , называемый таблицей, и хеш-функцию $h : U \rightarrow \{0, 1, \dots, m - 1\}$. Хеш-функция отображает ключ $k \in U$ в индекс $h(k)$ в массиве T . Элемент $T[j]$ называется корзиной (bucket).

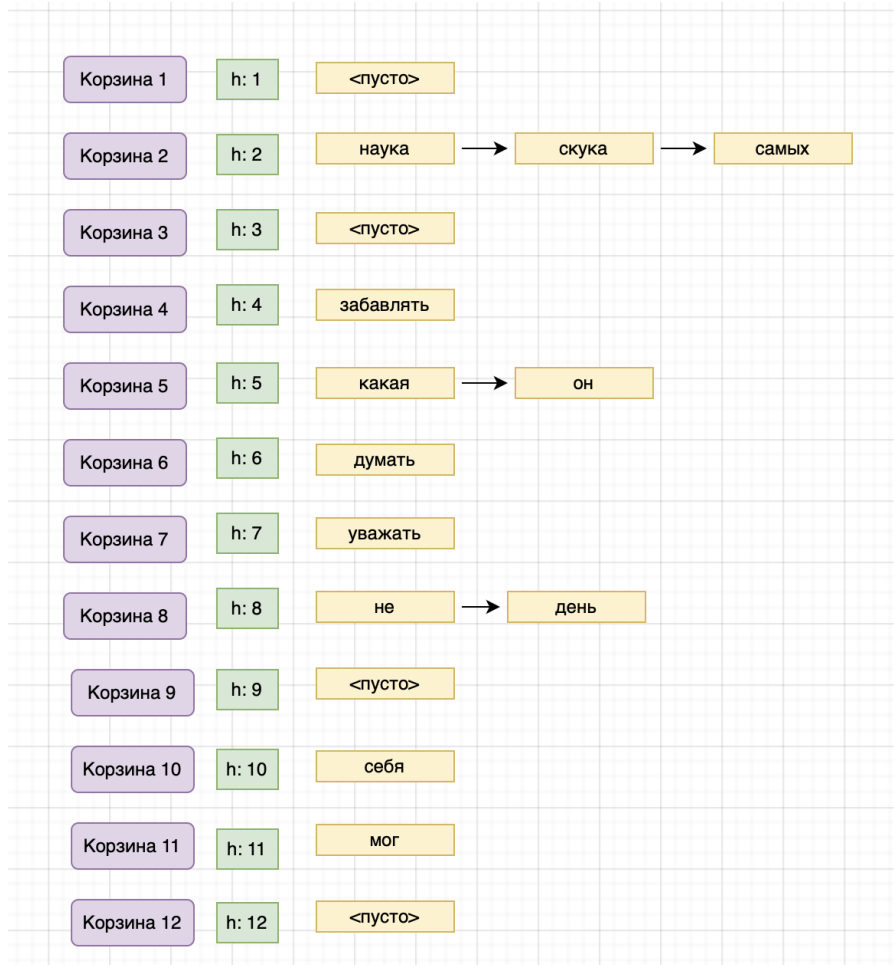


Рис. 1. Хеш-таблица.

Хеш-функция. В данной работе для строковых ключей $S = c_0c_1 \dots c_{L-1}$ (где c_i — байтовое представление i -го символа строки, L — длина строки в байтах) применяется полиномиальная хеш-функция. Значение хеш-функции вычисляется следующим образом:

$$h(S) = \left(\sum_{i=0}^{L-1} \text{ord}(c_i) \cdot p^i \right) \pmod{m}$$

где $\text{ord}(c_i)$ — числовое значение байта c_i (например, его ASCII-код или значение как ‘unsigned char’), p — некоторое простое число, выбранное в качестве основания полинома ($p = 31$), и m — текущий размер хеш-таблицы (количество корзин). Операция взятия по модулю m обеспечивает отображение результата в диапазон индексов таблицы $[0, m - 1]$. Для предотвращения переполнения при вычислении больших степеней p^i и больших сумм, модульная арифметика применяется на каждом шаге итеративного вычисления:

1. Инициализация: $\text{hash_value} = 0, p_power = 1$.
2. Для каждого символа c_i из S (от $i = 0$ до $L - 1$):

(a) В LaTeX это будет:

$$\text{hash_value} = (\text{hash_value} + \text{ord}(c_i) \cdot p_power) \pmod{m}.$$

(b) В LaTeX это будет:

$$p_power = (p_power \cdot p) \pmod{m}.$$

3. Результат: $h(S) = \text{hash_value}$.

В реализации значение $\text{value}(c_i)$ (эквивалент $\text{ord}(c_i)$ в формуле) бралось как ‘static_cast<unsigned char>(c_byte)’.

Разрешение коллизий. Коллизия возникает, если для двух различных ключей $k_1 \neq k_2$ их хеш-значения совпадают: $h(k_1) = h(k_2)$. В данной реализации для разрешения коллизий используется **метод цепочек (separate chaining)**. Каждая корзина $T[j]$ хеш-таблицы является указателем на начало связанного списка (или другой динамической структуры), хранящего все пары (k, v) , для которых $h(k) = j$. Операции вставки, поиска и удаления элемента с ключом k требуют:

1. Вычисления индекса $j = h(k)$.
2. Выполнения соответствующей операции (поиск, вставка, удаление) в связанном списке, ассоциированном с корзиной $T[j]$.

Средняя длина цепочки при равномерном распределении ключей и коэффициенте загрузки $\alpha = N/m$ (где N — количество элементов в таблице) составляет α . Таким образом, ожидаемое время выполнения операций составляет $O(1 + \alpha)$ без учета времени вычисления хеш-функции.

Коэффициент загрузки и рехеширование. Коэффициент загрузки $\alpha = N/m$ является важным показателем производительности хеш-таблицы с методом цепочек. При увеличении α средняя длина цепочек растет, что приводит к увеличению времени выполнения операций. Для поддержания эффективности, когда α достигает некоторого порогового значения α_{max} (в реализации $\alpha_{max} = 0.75$), выполняется операция рехеширования. Она заключается в создании новой

таблицы большего размера m' (обычно $m' \approx 2m$), и все N элементов из старой таблицы переносятся в новую с пересчетом их хеш-кодов по новой хеш-функции $h' : U \rightarrow \{0, 1, \dots, m' - 1\}$. Амортизированная стоимость операции вставки при использовании рехеширования остается $O(1 + \alpha)$ (при условии, что стоимость рехеширования $O(m + N)$ "размазывается" по достаточному количеству вставок).

1.2 Красно-черное дерево

Красно-черное дерево (КЧ-дерево) — это тип самобалансирующегося двоичного дерева поиска. Каждому узлу в дереве приписывается дополнительный атрибут — цвет (красный или черный). За счет поддержания инвариантов, связанных с цветами узлов и структурой дерева, КЧ-деревья гарантируют, что ни один путь от корня до любого листа не будет более чем в два раза длиннее любого другого такого пути. Это обеспечивает логарифмическую высоту дерева $h = O(\log N)$, где N — количество узлов.

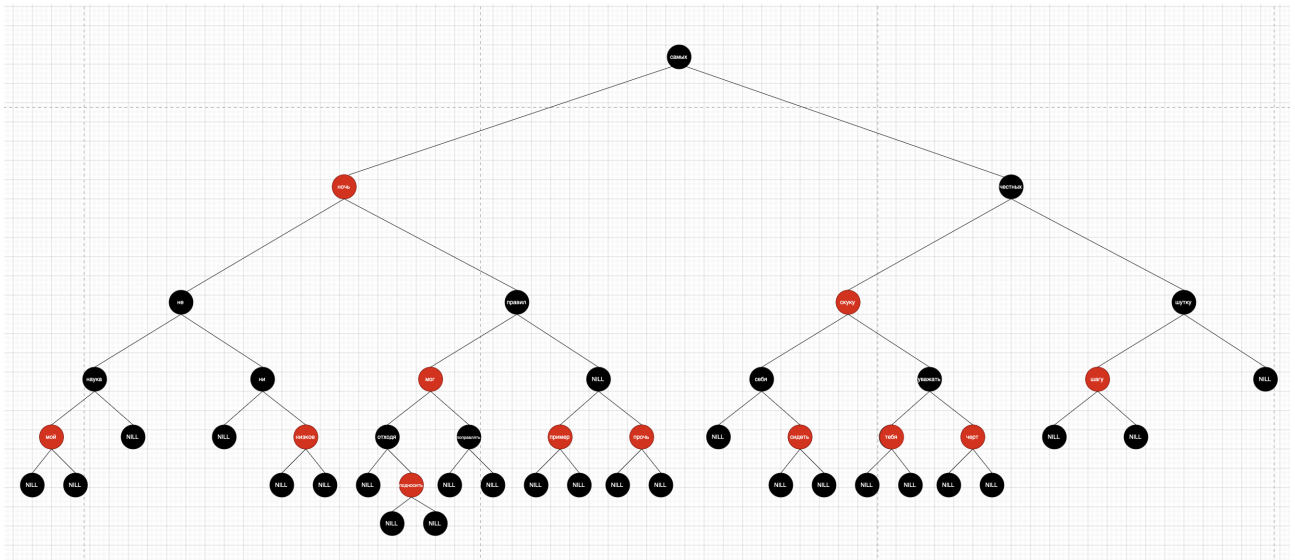


Рис. 2. Красно-черное дерево.

Свойства (инварианты) красно-черного дерева:

1. **Свойство цвета:** Каждый узел является либо красным (RED), либо черным (BLACK).
2. **Свойство корня:** Корень дерева всегда черный.
3. **Свойство листьев:** Все внешние узлы (листья, представляемые в реализации одним стражевым NIL-узлом) являются черными.
4. **Свойство красных узлов:** Если узел красный, то оба его дочерних узла (если они существуют и не являются NIL) должны быть черными. Это означает, что не может быть двух красных узлов, следующих друг за другом на пути от корня к листу.

5. **Свойство черной высоты:** Для любого узла x все простые пути от x до листовых узлов (NIL) в его поддереве содержат одинаковое количество черных узлов. Это количество называется черной высотой узла x , обозначается $bh(x)$.

Операции и балансировка. Для поддержания свойств КЧ-дерева при вставке или удалении узлов используются две основные операции: перекрашивание узлов и вращения.

1.2.1 Перекрашивание узлов

Это простейшая операция, заключающаяся в изменении цвета узла с красного на черный или с черного на красный. Сама по себе эта операция не изменяет структуру дерева, но используется в комбинации с вращениями для восстановления свойств КЧ-дерева.

1.2.2 Вращения (Rotations)

Вращения — это локальные структурные преобразования в двоичном дереве, которые изменяют связи «родитель-потомок», сохраняя при этом свойство упорядоченности двоичного дерева поиска (т.е. для любого узла x все ключи в его левом поддереве меньше $key[x]$, а все ключи в его правом поддереве больше $key[x]$). Вращения используются для изменения высот поддеревьев и помогают сбалансировать дерево. Существует два типа вращений: левое и правое. Они выполняются за константное время $O(1)$.

Левое вращение (Left Rotate) Левое вращение выполняется вокруг узла x , у которого есть правый дочерний узел y (не NIL). В результате вращения y становится новым корнем поддерева, которое ранее возглавлял x . Узел x становится левым дочерним узлом y . Левое поддерево y становится правым поддеревом x .

Правое вращение (Right Rotate) Правое вращение симметрично левому. Оно выполняется вокруг узла y , у которого есть левый дочерний узел x (не NIL). В результате x становится новым корнем поддерева, а y — его правым дочерним узлом. Правое поддерево x становится левым поддеревом y .

1.2.3 Операция Вставки (Insertion)

Вставка нового узла в КЧ-дерево включает два этапа:

1. Стандартная вставка узла как в обычном двоичном дереве поиска. Новый узел z всегда окрашивается в **красный** цвет. Его дочерние узлы (NIL) являются черными.
2. Вызов процедуры `RB-Insert-Fixup(T, z)` для восстановления свойств КЧ-дерева, если они были нарушены.

Вставка красного узла может нарушить только свойство 2 (если z — корень и он красный) или свойство 4 (если родитель z также красный). Свойство 5 (черная высота) не нарушается, так как мы добавляем красный узел, который не меняет количество черных узлов на путях.

RB-Insert-Fixup(T, z) Эта процедура восстанавливает свойства КЧ-дерева. Основной цикл продолжается, пока родитель узла z ($z.parent$) красный (что нарушает свойство 4). Обозначим $P = z.parent$ и $G = z.parent.parent$ (дедушка z). Рассматриваются три случая (и их симметричные аналоги).

Сценарий: Родитель P узла z является левым ребенком дедушки G . (Если P является правым ребенком G , случаи симметричны, с заменой «лево» на «право» и наоборот).

- **Случай 1: Дядя U узла z (правый ребенок G) красный.**

- $P.color \leftarrow BLACK$
- $U.color \leftarrow BLACK$
- $G.color \leftarrow RED$
- $z \leftarrow G$ (Проблема «два красных подряд» перемещается на два уровня вверх. Цикл продолжается для нового $z = G$).

Логика: Мы «проталкиваем» красноту дедушки вверх. Черная высота не нарушается, так как на путях через P и U один красный узел заменяется черным, а на путях через G один черный узел (старый G) заменяется красным, но его дети (P, U) стали черными, компенсируя это.

- **Случай 2: Дядя U узла z черный, и z является правым ребенком P (конфигурация «треугольник»).**

- $z \leftarrow P$
- $Left-Rotate(T, z)$ (т.е. вокруг старого P)

Логика: Это преобразование сводит ситуацию к Случаю 3. Узел z (теперь это старый P) и его новый родитель (старый z) по-прежнему красные, но теперь они образуют «линию».

- **Случай 3: Дядя U узла z черный, и z является левым ребенком P (конфигурация «линия»).**

- $P.color \leftarrow BLACK$
- $G.color \leftarrow RED$
- $Right-Rotate(T, G)$

Логика: Это преобразование устраняет нарушение свойства 4. Правое вращение вокруг G делает P новым корнем этого поддерева. P становится черным, его ребенок z остается красным, а другой ребенок G (бывший родитель P) становится красным. Поскольку P теперь черный, свойство 4 не нарушено. Цикл завершается.

После завершения цикла свойство 2 (корень должен быть черным) восстанавливается принудительным окрашиванием $T.root.color \leftarrow BLACK$.

1.2.4 Операция Удаления (Deletion)

Удаление узла из КЧ-дерева является более сложной операцией.

1. Сначала выполняется стандартная процедура удаления узла из двоичного дерева поиска. Это сводится к удалению узла y , который имеет не более одного не-NIL ребенка. Пусть x — это единственный ребенок y (или $T.NIL$, если у y нет детей), который занимает место y .
2. Сохраняется исходный цвет удаляемого узла y ($y_original_color$).
3. Если $y_original_color = RED$, то свойства КЧ-дерева не нарушаются. Красные узлы не влияют на черную высоту, и их удаление не может создать два красных узла подряд или сделать корень красным (если он не был единственным узлом).
4. Если $y_original_color = BLACK$, то удаление черного узла нарушает свойство 5 (черная высота) для всех путей, проходивших через y . Также могут быть нарушены свойства 2 или 4. Для восстановления свойств вызывается процедура $RB\text{-}Delete\text{-}Fixup(T, x)$. Узел x считается несущим «дополнительный черный цвет».

$RB\text{-}Delete\text{-}Fixup(T, x)$ Эта процедура восстанавливает свойства КЧ-дерева после удаления черного узла. Цикл продолжается, пока x не является корнем и $x.color = BLACK$. Обозначим $P = x.parent$ и w — брат узла x .

Сценарий: Узел x является левым ребенком своего родителя P .

- **Случай 1: Брат w узла x красный.**

- $w.color \leftarrow BLACK$
- $P.color \leftarrow RED$
- $Left\text{-}Rotate(T, P)$
- $w \leftarrow P.right$ (Обновляем w на нового брата x , который теперь будет черным)

Логика: Это преобразование приводит к одному из следующих случаев (2, 3 или 4), где брат x уже черный. «Дополнительный черный» на x сохраняется.

- **Случай 2:** Брат w узла x черный, и оба ребенка w черные.

- $w.color \leftarrow RED$
- $x \leftarrow P$ (Перемещаем «дополнительный черный» вверх к родителю P . Цикл продолжается для нового $x = P$).

Логика: Черный цвет «снимается» с x и w и «добавляется» к P . Пути через w теперь имеют на один черный узел меньше (сам w), что компенсирует удаленный черный узел на путях через x .

- **Случай 3:** Брат w узла x черный, левый ребенок w красный, а правый ребенок w черный.

- $w.left.color \leftarrow BLACK$
- $w.color \leftarrow RED$
- $Right\text{-}Rotate(T, w)$
- $w \leftarrow P.right$ (Обновляем w на нового брата x)

Логика: Это преобразование приводит к Случаю 4. «Дополнительный черный» на x сохраняется.

- **Случай 4:** Брат w узла x черный, и правый ребенок w красный.

- $w.color \leftarrow P.color$
- $P.color \leftarrow BLACK$
- $w.right.color \leftarrow BLACK$
- $Left\text{-}Rotate(T, P)$
- $x \leftarrow T.root$ (Устанавливаем x на корень, чтобы завершить цикл. «Дополнительный черный» устранен).

Логика: Это преобразование полностью устраняет «дополнительный черный». Цвета перераспределяются, и выполняется вращение для восстановления черных высот.

После завершения цикла, если x не NIL, его цвет устанавливается в черный ($x.color \leftarrow BLACK$) для удаления любого оставшегося «дополнительного черного» и обеспечения свойства корня.

1.3 Алгоритм сжатия RLE (Run-Length Encoding)

Алгоритм RLE — это простой метод сжатия данных без потерь, который эффективен для данных, содержащих много последовательностей повторяющихся значений.

Основной принцип RLE: Заменить последовательность одинаковых символов на пару, состоящую из самого символа и количества его повторений.

Продвинутый RLE: Для повышения эффективности на данных, где длинные серии повторов редки, используется модифицированный подход, который различает два типа последовательностей:

1. **Серия повторов (Run):** Если символ s повторяется N раз, и $N \geq \text{MIN_RUN_LENGTH}$, то эта серия кодируется как пара: счетчик N и сам символ s .
2. **Литеральная последовательность (Literal Run):** Если последовательность символов S_{lit} длиной M не содержит серий повторов длиной MIN_RUN_LENGTH или более, то она кодируется как: отрицательное значение длины $-M$, за которым следует специальный символ-разделитель '#', и затем сама последовательность S_{lit} .

Такой подход позволяет избежать увеличения размера данных для коротких неповторяющихся последовательностей, что характерно для простого RLE. Использование символа-разделителя '#' между счетчиком/длиной и данными помогает декодеру однозначно парсить закодированный поток, особенно когда сами данные могут содержать цифры.

Коэффициент сжатия (CR): Определяется как отношение размера исходных данных $Size_{orig}$ к размеру сжатых данных $Size_{comp}$:

$$CR = \frac{Size_{orig}}{Size_{comp}}$$

Если $CR > 1$, данные сжаты. Если $CR < 1$, размер данных увеличился. Если $CR = 1$, размер не изменился.

Цена кодирования: Для RLE "цена кодирования" не является статистической величиной, как средняя длина кодового слова в энтропийных кодах. Фактически, это размер сжатых данных $Size_{comp}$. Он полностью зависит от структуры входных данных.

1.4 Математическое и алгоритмическое описание кодирования Фано

Алгоритм Фано предназначен для построения префиксных кодов переменной длины для символов источника с известными вероятностями их появления. Целью является минимизация средней длины кодового слова.

Пусть дан алфавит источника $S = \{s_1, s_2, \dots, s_N\}$ и соответствующие им вероятности $P = \{p_1, p_2, \dots, p_N\}$, где p_i — вероятность символа s_i , и $\sum_{i=1}^N p_i = 1$.

Процесс построения кодов Фано включает следующие шаги:

1. **Упорядочивание символов.** Все символы $s_i \in S$ сортируются в порядке невозрастания их вероятностей p_i : $p_1 \geq p_2 \geq \dots \geq p_N$.

2. **Рекурсивное разделение.** Процесс построения кодов рекурсивен. Для текущего (отсортированного) списка символов $S_{sub} = \{s_j, s_{j+1}, \dots, s_k\}$:

(a) Если список S_{sub} содержит только один символ ($j = k$), то процесс для этой ветви завершен. Накопленная последовательность битов является кодом для этого символа.

(b) Если список S_{sub} содержит более одного символа, он разделяется на два непустых подсписка $S_{sub}^{(0)} = \{s_j, \dots, s_m\}$ и $S_{sub}^{(1)} = \{s_{m+1}, \dots, s_k\}$ ($j \leq m < k$). Разделение производится таким образом, чтобы сумма вероятностей символов в $S_{sub}^{(0)}$ была как можно ближе к сумме вероятностей символов в $S_{sub}^{(1)}$. То есть, ищется такой индекс m , который минимизирует величину:

$$\left| \sum_{l=j}^m p_l - \sum_{l=m+1}^k p_l \right|$$

Или, что эквивалентно, $\sum_{l=j}^m p_l$ должна быть максимально близка к $\frac{1}{2} \sum_{l=j}^k p_l$.

(c) Всем символам из подсписка $S_{sub}^{(0)}$ к их текущему кодовому префиксу дописывается бит '0'.

(d) Всем символам из подсписка $S_{sub}^{(1)}$ к их текущему кодовому префиксу дописывается бит '1'.

(e) Процесс рекурсивно повторяется для подсписков $S_{sub}^{(0)}$ и $S_{sub}^{(1)}$.

Начальный вызов рекурсивной процедуры происходит для полного отсортированного списка символов S с пустым начальным префиксом кода.

Свойства кодов: Коды, полученные методом Фано, являются префиксными, что обеспечивает однозначность декодирования. Однако они не всегда являются оптимальными с точки зрения минимизации средней длины кодового слова. Средняя длина кодового слова L_{avg} для кодов Фано удовлетворяет неравенству Шеннона: $H(S) \leq L_{avg} < H(S) + 1$, где $H(S)$ – энтропия источника.

Декодирование: Декодирование осуществляется путем последовательного чтения битов из сжатого потока. Накопленная битовая последовательность сравнивается с префиксными кодами из построенной таблицы. При обнаружении совпадения с одним из кодов, соответствующий символ восстанавливается, и процесс продолжается для оставшейся части битового потока.

2 Реализация словаря на основе хеш-таблицы

2.1 Структура узла HashNode

```
1 namespace DictionaryWithHashTable {
2 struct HashNode {
3     std::string key;
4     int value;
5     HashNode(std::string k, int v) : key(std::move(k)), value(v) {}
6 }; }
```

Листинг 1: Структура HashNode

Описание членов структуры HashNode:

- key: Тип `std::string`. Хранит ключ элемента хеш-таблицы, а именно – слово из текста.
- value: Тип `int`. Хранит значение, ассоциированное с ключом, а именно – частоту встречаемости слова.

2.2 Класс HashTable

Хеш-таблица была реализована с использованием метода цепочек для разрешения коллизий. Основным хранилищем является вектор списков (`std::vector<std::list<HashNode>>`), где каждый элемент вектора (корзина) содержит список узлов с одинаковым хеш-значением.

```
1 class HashTable {
2 private:
3     std::vector<std::list<HashNode>> table;
4     size_t num_elements;
5     size_t table_size;
6     static constexpr double MAX_LOAD_FACTOR = 0.75;
7     size_t hashFunction(const std::string& key) const;
8     void rehash();
9 public:
10    HashTable(size_t initial_size = 101) : table_size(initial_size), num_elements(0);
11    void add(const std::string& key, int value);
12    int* get(const std::string& key);
13    const int* get(const std::string& key) const;
14    bool remove(const std::string& key);
15    void clear();
16    void print(std::ostream& os = std::cout) const;
17    void visualize(std::ostream& os = std::cout) const;
18 }; }
```

Листинг 2: Класс HashTable

Описание членов класса HashTable:

- **table:** Тип `std::vector<std::list<HashNode>`. Основное хранилище данных, массив корзин, где каждая корзина - список узлов.
- **num_elements:** Тип `size_t`. Текущее количество элементов (пар ключ-значение) в таблице.
- **table_size:** Тип `size_t`. Текущий размер таблицы.
- **MAX_LOAD_FACTOR:** Статическая константа типа `double`, равная 0.75. Пороговое значение коэффициента загрузки, при превышении которого происходит рехеширование.

2.3 Хеш-функция

- **Вход:** `key` – ключ.
- **Выход:** Значение типа `size_t`, представляющее собой хеш-код (индекс в хеш-таблице) для входного ключа.
- **Описание:** Данная функция реализует полиномиальный алгоритм хеширования для строковых ключей. Она итеративно обрабатывает каждый байт входной строки `key`.

```
1 size_t hashFunction(const std::string& key) const {
2     size_t hash_val = 0;
3     size_t p = 31;
4     size_t p_pow = 1;
5     for (char c_byte : key) {
6         hash_val = (hash_val + static_cast<unsigned char>(c_byte) * p_pow) %
7             table_size;
8         p_pow = (p_pow * p) % table_size;
9     }
10    return hash_val;
}
```

Листинг 3: Хеш-функция

2.4 Функция `add(const std::string& key, int value)`

- **Вход:** `key` – ключ (слово), `value` – значение (частота).
- **Выход:** занесенная пара ключ-значение в таблицу.
- **Описание:** Добавляет пару ключ-значение в таблицу. Если ключ уже существует, обновляет его значение. При необходимости (превышение `MAX_LOAD_FACTOR` при добавлении нового элемента) выполняет рехеширование.

```

1 void add(const std::string& key, int value) {
2     if (static_cast<double>(num_elements + 1) / table_size >= MAX_LOAD_FACTOR) {
3         rehash();
4     }
5
6     size_t index = hashFunction(key);
7     for (auto& node : table[index]) {
8         if (node.key == key) {
9             node.value = value;
10            return;
11        }
12    }
13    table[index].emplace_back(key, value);
14    num_elements++;
15 }

```

Листинг 4: Функция add в HashTable

2.5 Функция get(const std::string& key)

- **Вход:** key – ключ для поиска.
- **Выход:** Указатель на значение (int*), если ключ найден; nullptr в противном случае.
- **Описание:** Ищет элемент по ключу и возвращает указатель на его значение.

```

1 int* get(const std::string& key) {
2     size_t index = hashFunction(key);
3     for (auto& node : table[index]) {
4         if (node.key == key) {
5             return &node.value;
6         }
7     }
8     return nullptr;
9 }

```

Листинг 5: Функция get в HashTable

2.6 Функция remove(const std::string& key)

- **Вход:** key – ключ удаляемого элемента.
- **Выход:** true, если элемент удален; false, если элемент не найден.
- **Описание:** Удаляет элемент с заданным ключом из таблицы.


```

1 bool remove(const std::string& key) {
2     size_t index = hashFunction(key);
3     auto& bucket = table[index];
4     for (auto it = bucket.begin(); it != bucket.end(); ++it) {
5         if (it->key == key) {
6             bucket.erase(it);
7             num_elements--;
8             return true;
9         }
10    }
11    return false;
12 }

```

Листинг 6: Функция remove в HashTable

2.7 Функция clear()

- **Вход:** Текущая хеш-таблица.
- **Выход:** Очищенная хеш-таблица.
- **Описание:** Удаляет все элементы из хеш-таблицы, очищая все корзины и сбрасывая счетчик элементов.

```

1 void clear() {
2     for (auto& bucket : table) {
3         bucket.clear();
4     }
5     num_elements = 0;
6 }

```

Листинг 7: Функция clear в HashTable

2.8 Приватная функция rehash()

- **Вход:** Текущая хеш-таблица.
- **Выход:** Увеличенная хеш-таблица.
- **Описание:** Увеличивает размер таблицы (`table_size`), создает новую таблицу и переносит в нее все существующие элементы с пересчетом их хеш-значений.

```

1 void rehash() {
2     size_t old_table_size = table_size;
3     table_size = table_size * 2 + 1;
4     std::vector<std::list<HashNode>> old_table = std::move(table);
5
6     table.assign(table_size, std::list<HashNode>());
7     num_elements = 0;
8     for (size_t i = 0; i < old_table_size; ++i) {
9         for (const auto& node : old_table[i]) {
10             add(node.key, node.value);
11         }
12     }
13 }

```

Листинг 8: Функция rehash в HashTable

2.9 Публичная функция print

Прототип: `void print(std::ostream& os = std::cout) const`

- **Вход:** `os: std::ostream&` – ссылка на выходной поток (по умолчанию `std::cout`).
- **Выход:** вывод осуществляется в поток `os`.
- **Описание:** Выводит все элементы хеш-таблицы в указанный поток. Итерирует по каждой корзине таблицы и по каждому узлу в непустых корзинах. Пары ключ-значение выводятся в формате 'ключ': значение, разделенные запятыми, и всё содержимое заключается в фигурные скобки.

```

1 void print(std::ostream& os = std::cout) const {
2     os << "{";
3     bool first_item = true;
4     for (const auto& bucket : table) {
5         for (const auto& node : bucket) {
6             if (!first_item) {
7                 os << ", ";
8             }
9             os << ">" << node.key << ": " << node.value;
10            first_item = false;
11        }
12    }
13    os << "} ";
14 }

```

2.10 Публичная функция visualize

Прототип: `void visualize(std::ostream& os = std::cout) const`

- **Вход:** `os: std::ostream&` – ссылка на выходной поток (по умолчанию `std::cout`).
- **Выход:** осуществляется вывод в поток `os`.
- **Описание:** Выводит детальное представление хеш-таблицы. Для каждой корзины отображается ее номер и содержимое. Если корзина пуста, выводится <пусто>. Если в корзине есть элементы (цепочка коллизий), они выводятся последовательно в формате ("ключ": значение), разделенные символом ->. Также выводится общий размер таблицы и количество элементов.

```

1 void visualize(std::ostream& os = std::cout) const {
2     os << "Визуализация Хеш-таблицы (размер: " << table_size
3         << ", элементы: " << num_elements << "):" << std::endl;
4     for (size_t i = 0; i < table.size(); ++i) {
5         os << "Корзина [" << i << "]: ";
6         if (table[i].empty()) {
7             os << "<пусто>";
8         } else {
9             bool first_in_bucket = true;
10            for (const auto& node : table[i]) {
11                if (!first_in_bucket) {
12                    os << " -> ";
13                }
14                os << "(" << node.key << ": " << node.value << ")";
15                first_in_bucket = false;
16            }
17        }
18        os << std::endl;
19    }
20 }

```

Листинг 10: Функция visualize в HashTable

2.11 Класс Dictionary

Класс `Dictionary` инкапсулирует объект `HashTable` и предоставляет интерфейс для работы со словарем.

```

1 class Dictionary {
2 private:
3     HashTable ht;
4     std::string toLowerASCII(std::string s) const;
5
6 public:
7     Dictionary(size_t initial_capacity = 101) : ht(initial_capacity) {}
8
9     void addWord(const std::string& word_raw);
10    void removeWord(const std::string& word_raw);
11    bool findWord(const std::string& word_raw) const;
12    void clear();
13    void loadFromFile(const std::string& filepath, bool append = false);
14    void print(std::ostream& os = std::cout) const;
15    void visualizeStructure(std::ostream& os = std::cout) const;
16 };

```

Листинг 11: Функция print в HashTable

2.12 Функция addWord(const std::string& word_raw)

- **Вход:** word_raw – слово для добавления или обновления частоты.
- **Выход:** слово добавлено в словарь.
- **Описание:** Приводит слово к нижнему регистру (ASCII-символы). Если слово уже есть в словаре, увеличивает его частоту. В противном случае добавляет слово с частотой 1.

```

1 void addWord(const std::string& word_raw) {
2     if (word_raw.empty()) return;
3     std::string word = toLowerASCII(word_raw);
4
5     int* current_val_ptr = ht.get(word);
6     if (current_val_ptr) {
7         (*current_val_ptr)++;
8     } else {
9         ht.add(word, 1);
10    }
11 }

```

Листинг 12: Функция addWord в Dictionary (HashTable)

2.13 Функция removeWord(const std::string& word_raw)

- **Вход:** word_raw – слово для удаления из словаря.
- **Выход:** слово удалено из словаря.
- **Описание:** Приводит слово к нижнему регистру (ASCII) и удаляет его из внутренней хеш-таблицы.

```
1 void removeWord(const std::string& word_raw) {  
2     if (word_raw.empty()) return;  
3     std::string word = toLowerASCII(word_raw);  
4     ht.remove(word);  
5 }
```

Листинг 13: Функция removeWord в Dictionary (HashTable)

2.14 Функция findWord(const std::string& word_raw) const

- **Вход:** word_raw – слово для поиска.
- **Выход:** true, если слово найдено; false в противном случае. Также выводит информацию о частоте слова, если оно найдено.
- **Описание:** Приводит слово к нижнему регистру (ASCII), ищет его в хеш-таблице и сообщает о результате.

```
1 bool findWord(const std::string& word_raw) const {  
2     if (word_raw.empty()) return false;  
3     std::string word = toLowerASCII(word_raw);  
4  
5     const int* count_ptr = ht.get(word);  
6     if (count_ptr) {  
7         std::cout << "Слово '" << word_raw << "' (ключ: '" << word  
8             << "') найдено, частота: " << *count_ptr << std::endl;  
9         return true;  
10    } else {  
11        std::cout << "Слово '" << word_raw << "' (ключ: '" << word  
12            << "') не найдено." << std::endl;  
13        return false;  
14    }  
15 }
```

Листинг 14: Функция findWord в Dictionary (HashTable)

2.15 Функция loadFromFile(const std::string& filepath, bool append = false)

- **Вход:** filepath – путь к текстовому файлу, append (опционально, по умолчанию false) – флаг, указывающий, нужно ли добавлять слова к существующему словарю или очистить его перед загрузкой.
- **Выход:** текст загружен из файла.
- **Описание:** Читает текстовый файл, разбивает его на слова и добавляет их в словарь с подсчетом частот. Использует глобальные функции readFileToString и processTextToWords.

```
1 void loadFromFile(const std::string& filepath, bool append = false) {
2     if (!append) {
3         clear();
4     }
5     try {
6         std::string content = readFileToString(filepath);
7         std::vector<std::string> words = processTextToWords(content);
8         for (const std::string& word : words) {
9             if (!word.empty()) addWord(word);
10        }
11        std::cout << "Словарь загружен/дополнен из файла '" << filepath
12                << "' (хеш-таблица)." << std::endl;
13    } catch (const std::runtime_error& e) {
14        std::cerr << "Ошибка при загрузке из файла (хеш-таблица): "
15                << e.what() << std::endl;
16    }
17 }
```

Листинг 15: Функция loadFromFile в Dictionary (HashTable)

3 Реализация словаря на основе Красно-Черного дерева

Для альтернативной реализации словаря было использовано Красно-Черное дерево (КЧ-дерево) — самобалансирующееся бинарное дерево поиска. Это обеспечивает гарантированную логарифмическую сложность $O(\log N)$ для операций вставки, удаления и поиска в худшем случае.

3.1 Структура данных: Красно-Черное дерево

```
1 namespace DictionaryWithRBTree {
2
3 enum Color { RED, BLACK };
4
5 struct Node {
6     std::string key;
7     int value;
8     Color color;
9     Node *parent, *left, *right;
10
11     Node(std::string k, int v, Color c = RED, Node* p = nullptr, Node* l = nullptr,
12         Node* r = nullptr)
13         : key(std::move(k)), value(v), color(c), parent(p), left(l), right(r) {}
14 };
```

Листинг 16: Перечисление Color

Описание членов структуры Node:

- **key:** Тип `std::string`. Хранит ключ (слово).
- **value:** Тип `int`. Хранит значение (частоту слова).
- **color:** Тип `Color`. Хранит цвет узла.
- **parent:** Указатель на родительский узел `Node*`.
- **left:** Указатель на левого дочернего узла `Node*`.
- **right:** Указатель на правого дочернего узла `Node*`.

3.2 Класс RBTree

```
1 class RBTree {
2 private:
3     Node* root;
4     Node* NIL;
```

```

5   void leftRotate(Node* x);
6   void rightRotate(Node* y);
7   void insertFixup(Node* z);
8   Node* findNode(const std::string& key) const;
9   void transplant(Node* u, Node* v);
10  Node* minimum(Node* node);
11  void deleteFixup(Node* x);
12  void destroyRecursive(Node* node) ;
13  void inorderPrintRecursive(Node* node, std::ostream& os, bool& first_item) const;
14  bool printGivenLevel(Node* node, int level, std::ostream& os, bool first_on_level
    ) const;
15  int getMaxDepth(Node* node) const;
16  void printTreeRecursive(Node* node, int space_increment, int current_space, std::
    ostream& os) const ;
17
18 public:
19     RBTREE();
20     ~RBTREE() {
21         destroyRecursive(root);
22         delete NIL;
23     }
24     RBTREE(const RBTREE&) = delete;
25     RBTREE& operator=(const RBTREE&) = delete;
26     void insert(const std::string& key, int value);
27     int* search(const std::string& key);
28     const int* search(const std::string& key) const;
29     bool remove(const std::string& key);
30     void clear();
31     void print(std::ostream& os = std::cout) const;
32     void visualize(std::ostream& os = std::cout) const;
33 };

```

Листинг 17: Перечисление Color

Описание членов класса RBTREE:

- **root:** Указатель Node* на корень дерева.
- **NIL:** Указатель Node* на специальный стражевой (sentinel) NIL-узел. Все "отсутствующие" листья и родитель корня указывают на этот узел. Он всегда черный.

3.3 Конструктор RBTREE()

- **Вход:** пустой объект RBTREE.
- **Выход:** создается объект RBTREE.

- **Описание:** Инициализирует дерево: создает стражевой NIL-узел (черного цвета, все указатели которого ссылаются на него же) и устанавливает `root` на этот NIL-узел, обозначая пустое дерево.

```

1 RBTREE() {
2     NIL = new Node("", 0, BLACK);
3     NIL->parent = NIL;
4     NIL->left = NIL;
5     NIL->right = NIL;
6     root = NIL;
7 }

```

Листинг 18: Конструктор RBTREE

3.4 Функция `insert(const std::string& key, int value)`

- **Вход:** `key` – ключ (слово), `value` – значение (частота).
- **Выход:** новый узел.
- **Описание:** Вставляет новый узел с заданным ключом и значением в дерево. Новый узел изначально красный. Если ключ уже существует, обновляет его значение. После вставки вызывается процедура `insertFixup` для восстановления свойств КЧ-дерева.

```

1 void insert(const std::string& key, int value) {
2     Node* z = new Node(key, value, RED, NIL, NIL, NIL);
3     Node* y = NIL;
4     Node* x = root;
5
6     while (x != NIL) {
7         y = x;
8         if (z->key < x->key) {
9             x = x->left;
10        } else if (z->key > x->key) {
11            x = x->right;
12        } else {
13            x->value = value;
14            delete z;
15            return;
16        }
17    }
18
19    z->parent = y;
20    if (y == NIL) {

```

```

21         root = z;
22     } else if (z->key < y->key) {
23         y->left = z;
24     } else {
25         y->right = z;
26     }
27     insertFixup(z);
28 }

```

Листинг 19: Функция insert в RBTre

3.5 Функция search(const std::string& key)

- **Вход:** key – ключ для поиска.
- **Выход:** Указатель на значение (int*), если ключ найден; nullptr в противном случае.
- **Описание:** Выполняет поиск узла по ключу, как в стандартном бинарном дереве поиска.

```

1 int* search(const std::string& key) {
2     Node* node = findNode(key);
3     return (node == NIL) ? nullptr : &node->value;
4 }
5 const int* search(const std::string& key) const {
6     Node* node = findNode(key);
7     return (node == NIL) ? nullptr : &node->value;
8 }

```

Листинг 20: Функция search в RBTre

3.6 Функция remove(const std::string& key)

- **Вход:** key – ключ удаляемого элемента.
- **Выход:** true, если элемент удален; false, если элемент не найден.
- **Описание:** Удаляет узел с заданным ключом из дерева. После физического удаления узла (или замены его преемником) вызывается процедура deleteFixup, если был удален черный узел, для восстановления свойств КЧ-дерева.

```

1 bool remove(const std::string& key) {
2     Node* z = findNode(key);
3     if (z == NIL) return false;
4
5     Node* y = z;

```

```

6      Node* x;
7      Color y_original_color = y->color;
8
9      if (z->left == NIL) {
10         x = z->right;
11         transplant(z, z->right);
12     } else if (z->right == NIL) {
13         x = z->left;
14         transplant(z, z->left);
15     } else {
16         y = minimum(z->right);
17         y_original_color = y->color;
18         x = y->right;
19         if (y->parent == z) {
20             if (x != NIL) x->parent = y;
21         } else {
22             transplant(y, y->right);
23             y->right = z->right;
24             y->right->parent = y;
25         }
26         transplant(z, y);
27         y->left = z->left;
28         y->left->parent = y;
29         y->color = z->color;
30     }
31     delete z;
32
33     if (y_original_color == BLACK) {
34         deleteFixup(x);
35     }
36     return true;
37 }

```

Листинг 21: Функция remove в RBTTree

3.7 Приватные функции leftRotate rightRotate

- **Вход:** `x: Node*` – указатель на узел, относительно которого выполняется вращение.
- **Выход:** модифицированная структура дерева.
- **Описание:** Осуществляет левое вращение дерева вокруг узла `x`. Правый дочерний узел `y` узла `x` становится новым корнем данного поддерева. Узел `x` становится левым дочерним узлом `y`. Левое поддерево `y` становится правым поддеревом `x`. Функция корректно обновляет

все указатели (`parent`, `left`, `right`) затронутых узлов и, при необходимости, указатель `root` всего дерева.

```
1  void leftRotate(Node* x) {
2      Node* y = x->right;
3      x->right = y->left;
4      if (y->left != NIL) {
5          y->left->parent = x;
6      }
7      y->parent = x->parent;
8      if (x->parent == NIL) {
9          root = y;
10     } else if (x == x->parent->left) {
11         x->parent->left = y;
12     } else {
13         x->parent->right = y;
14     }
15     y->left = x;
16     x->parent = y;
17 }
18
19 void rightRotate(Node* y) {
20     Node* x = y->left;
21     y->left = x->right;
22     if (x->right != NIL) {
23         x->right->parent = y;
24     }
25     x->parent = y->parent;
26     if (y->parent == NIL) {
27         root = x;
28     } else if (y == y->parent->right) {
29         y->parent->right = x;
30     } else {
31         y->parent->left = x;
32     }
33     x->right = y;
34     y->parent = x;
35 }
```

Листинг 22: Функция `leftRotate` в `RBTree`

3.8 Приватная функция `insertFixup`

- **Вход:** `z: Node*` – указатель на вставленный узел (изначально красный).

- **Выход:** модифицированная структура дерева.
- **Описание:** Восстанавливает свойства красно-черного дерева после вставки узла **z**. Если родитель **z** также красный (нарушение), функция итеративно применяет вращения и перекрашивания, рассматривая три основных случая, основанных на цвете "дяди" узла **z** и конфигурации узлов (родитель **z** как левый/правый ребенок "дедушки" **z**; **z** как левый/правый ребенок своего родителя). Цель — устранить последовательность из двух красных узлов и сохранить черную высоту. В конце корень дерева всегда окрашивается в черный.

```

1 void insertFixup(Node* z) {
2     while (z->parent->color == RED) {
3         if (z->parent == z->parent->parent->left) {
4             Node* y = z->parent->parent->right;
5             if (y->color == RED) {
6                 z->parent->color = BLACK;
7                 y->color = BLACK;
8                 z->parent->parent->color = RED;
9                 z = z->parent->parent;
10            } else {
11                if (z == z->parent->right) {
12                    z = z->parent;
13                    leftRotate(z);
14                }
15                z->parent->color = BLACK;
16                z->parent->parent->color = RED;
17                rightRotate(z->parent->parent);
18            }
19        } else {
20            Node* y = z->parent->parent->left;
21            if (y->color == RED) {
22                z->parent->color = BLACK;
23                y->color = BLACK;
24                z->parent->parent->color = RED;
25                z = z->parent->parent;
26            } else {
27                if (z == z->parent->left) {
28                    z = z->parent;
29                    rightRotate(z);
30                }
31                z->parent->color = BLACK;
32                z->parent->parent->color = RED;
33                leftRotate(z->parent->parent);
34            }
35        }
36    }
37 }

```

```

36         if (z == root) break;
37     }
38     root->color = BLACK;
39 }

```

Листинг 23: Функция insertFixup в RBTre

3.9 Приватная функция findNode

- **Вход:** `key: const std::string&` – ключ (слово) для поиска.
- **Выход:** `Node*` – указатель на найденный узел, либо на `NIL`, если узел не найден.
- **Описание:** Осуществляет поиск узла с заданным ключом в красно-черном дереве. Поиск выполняется аналогично стандартному бинарному дереву поиска: начиная с корня, происходит сравнение ключей и переход в левое или правое поддерево до тех пор, пока узел не будет найден или не будет достигнут `NIL`-узел.

```

1 Node* findNode(const std::string& key) const {
2     Node* current = root;
3     while (current != NIL && current->key != key) {
4         if (key < current->key) {
5             current = current->left;
6         } else {
7             current = current->right;
8         }
9     }
10    return current;
11 }

```

Листинг 24: Функция findNode в RBTre

3.10 Приватная функция transplant

Прототип: `void transplant(Node* u, Node* v)`

- **Вход:** `u: Node*` – указатель на узел, который заменяется и `v: Node*` – указатель на узел, который ставится на место `u` (может быть `NIL`).
- **Выход:** модифицированная структура дерева.
- **Описание:** Вспомогательная функция для операции удаления. Заменяет поддерево с корнем `u` поддеревом с корнем `v`. Обновляет соответствующий дочерний указатель у родителя узла `u` (или указатель `root`, если `u` был корнем) так, чтобы он указывал на `v`. Если `v` не `NIL`, его указатель на родителя устанавливается на бывшего родителя `u`.

```

1 void transplant(Node* u, Node* v) {
2     if (u->parent == NIL) {
3         root = v;
4     } else if (u == u->parent->left) {
5         u->parent->left = v;
6     } else {
7         u->parent->right = v;
8     }
9     v->parent = u->parent;
10 }

```

Листинг 25: Функция transplant в RBTre

3.11 Приватная функция deleteFixup

Прототип: void deleteFixup(Node* x)

- **Вход:** x: Node* – узел, который условно несет "дополнительный черный" цвет после удаления другого черного узла.
- **Выход:** модифицирует цвета и структуру дерева.
- **Описание:** Восстанавливает свойства красно-черного дерева после удаления черного узла. Функция итеративно обрабатывает узел x, который считается "дважды черным для устранения нарушения черной высоты. Рассматриваются четыре случая (и их симметричные варианты) в зависимости от цвета брата x и его детей, применяя вращения и перекрашивания для перераспределения черных узлов или "проталкивания" проблемы вверх по дереву. Цикл завершается, когда "дополнительный черный" цвет устранен или достигнут корень. Корень дерева в конце всегда окрашивается в черный.

```

1 void deleteFixup(Node* x) {
2     while (x != root && x->color == BLACK) {
3         if (x == x->parent->left) {
4             Node* w = x->parent->right;
5             if (w->color == RED) {
6                 w->color = BLACK;
7                 x->parent->color = RED;
8                 leftRotate(x->parent);
9                 w = x->parent->right;
10            }
11            if (w->left->color == BLACK && w->right->color == BLACK) {
12                w->color = RED;
13                x = x->parent;

```

```

14     } else {
15         if (w->right->color == BLACK) {
16             w->left->color = BLACK;
17             w->color = RED;
18             rightRotate(w);
19             w = x->parent->right;
20         }
21         w->color = x->parent->color;
22         x->parent->color = BLACK;
23         w->right->color = BLACK;
24         leftRotate(x->parent);
25         x = root;
26     }
27 } else {
28     Node* w = x->parent->left;
29     if (w->color == RED) {
30         w->color = BLACK;
31         x->parent->color = RED;
32         rightRotate(x->parent);
33         w = x->parent->left;
34     }
35     if (w->right->color == BLACK && w->left->color == BLACK) {
36         w->color = RED;
37         x = x->parent;
38     } else {
39         if (w->left->color == BLACK) {
40             w->right->color = BLACK;
41             w->color = RED;
42             leftRotate(w);
43             w = x->parent->left;
44         }
45         w->color = x->parent->color;
46         x->parent->color = BLACK;
47         w->left->color = BLACK;
48         rightRotate(x->parent);
49         x = root;
50     }
51 }
52 }
53 x->color = BLACK;
54 }

```

Листинг 26: Функция deleteFixup в RBTree

3.12 Приватная функция destroyRecursive

Прототип: `void destroyRecursive(Node* node)`

- **Вход:** `node: Node*` – указатель на корень поддерева, которое нужно удалить.
- **Выход:** освобождает память.
- **Описание:** Рекурсивно удаляет все узлы в поддереве, корнем которого является `node`. Использует обратный обход (post-order): сначала рекурсивно обрабатывает левое и правое поддерева, а затем удаляет сам узел `node`. Базовый случай рекурсии – достижение NIL-узла. Эта функция используется для полной очистки дерева в методе `clear()` и в деструкторе.

```
1 void destroyRecursive(Node* node) {  
2     if (node != NIL) {  
3         destroyRecursive(node->left);  
4         destroyRecursive(node->right);  
5         delete node;  
6     }  
7 }
```

Листинг 27: Функция destroyRecursive в RBTre

3.13 Приватная функция inorderPrintRecursive

Прототип: `void inorderPrintRecursive(Node* node, std::ostream& os, bool& first_item)`
`const`

- **Вход:**
 - `node: Node*` – текущий узел для обработки.
 - `os: std::ostream&` – поток для вывода.
 - `first_item: bool&` – флаг, указывающий, является ли текущий выводимый элемент первым в общем списке (для корректной расстановки запятых).
- **Выход:** осуществляется вывод в поток `os`.
- **Описание:** Рекурсивно обходит дерево в инфиксном (симметричном) порядке и выводит информацию о каждом узле (ключ, значение, цвет) в поток `os`. Благодаря инфиксному обходу бинарного дерева поиска, ключи выводятся в лексикографическом (отсортированном) порядке. Флаг `first_item` используется для корректного форматирования вывода с запятыми-разделителями.

```

1 void inorderPrintRecursive(Node* node, std::ostream& os, bool& first_item) const {
2     if (node != NIL) {
3         inorderPrintRecursive(node->left, os, first_item);
4         if (!first_item) {
5             os << ", ";
6         }
7         os << "<< node->key << \": \" << node->value << \" (\" << (node->color == RED
8             ? \"R\" : \"B\") << \")\";
9         first_item = false;
10        inorderPrintRecursive(node->right, os, first_item);
11    }
12 }

```

Листинг 28: Функция `inorderPrintRecursive` в `RBTree`

3.14 Приватная функция `getMaxDepth`

Прототип: `int getMaxDepth(Node* node) const`

- **Вход:** `node: Node*` – указатель на корень поддерева, для которого вычисляется глубина.
- **Выход:** `int` – максимальная глубина поддерева (количество узлов на самом длинном пути от `node` до листового `NIL`-узла).
- **Описание:** Рекурсивно вычисляет максимальную глубину поддерева с корнем `node`. Для `NIL`-узла глубина равна 0. Для остальных узлов глубина вычисляется как $1 + \max(\text{глубина левого поддерева}, \text{глубина правого поддерева})$. Эта функция используется для определения количества уровней при визуализации дерева.

```

1 int getMaxDepth(Node* node) const {
2     if (node == NIL) return 0;
3     int left_depth = getMaxDepth(node->left);
4     int right_depth = getMaxDepth(node->right);
5     return std::max(left_depth, right_depth) + 1;
6 }

```

Листинг 29: Функция `getMaxDepth` в `RBTree`

3.15 Класс `Dictionary` (на основе `RBTree`)

Аналогичен словарию на хеш-таблице, но использует `RBTree` для хранения данных.

```

1 class Dictionary {
2 private:

```

```

3   RBTree rbt;
4   std::string toLowerASCII(std::string s) const {
5       std::transform(s.begin(), s.end(), s.begin(),
6           [](unsigned char c){ return std::tolower(c); });
7       return s;
8   }
9
10  public:
11      Dictionary() = default;
12      void addWord(const std::string& word_raw);
13      void removeWord(const std::string& word_raw);
14      bool findWord(const std::string& word_raw) const;
15      void clear();
16      void loadFromFile(const std::string& filepath, bool append = false);
17      void print(std::ostream& os = std::cout) const;
18      void visualizeStructure(std::ostream& os = std::cout) const;
19  };
20
21  }

```

Листинг 30: Класс Dictionary (на основе RBTree)

Основные функции класса Dictionary Функции Dictionary() (конструктор), addWord(), removeWord(), findWord(), clear(), loadFromFile() имеют то же назначение, входные и выходные параметры, что и в версии словаря на хеш-таблице, но для операций со словами они обращаются к соответствующим методам объекта rbt.

- **addWord**: вызывает rbt.search() и затем rbt.insert().
- **removeWord**: вызывает rbt.remove().
- **findWord**: вызывает rbt.search().
- **clear**: вызывает rbt.clear().
- **loadFromFile**: аналогично версии для хеш-таблицы, но использует addWord данного класса.

Код этих функций по структуре очень похож на версию для хеш-таблицы, с заменой вызовов ht на rbt.

4 Алгоритм сжатия RLE

4.1 Функция генерации случайного текста

- **Вход:** `filename: const std::string&` – имя файла для сохранения.
- **Выход:** `std::string` – сгенерированный текст.
- **Описание:** Генерирует псевдослучайный текст, имитирующий заданное распределение символов русского алфавита, цифр и спецсимволов. Для задания распределения используется весовой подход: частым символам (например, 'о', 'е', 'а', пробел) присваиваются большие веса. Выбор символа осуществляется с помощью `std::discrete_distribution`. С определенной вероятностью (в реализации 20%) функция генерирует серии повторяющихся символов (длиной от 3 до 10), что позволяет создавать данные, подходящие для тестирования RLE-сжатия. Сгенерированный текст сохраняется в файл и возвращается функцией.

```
1 const double CHAMPER_A = 1.57;
2 const double CHAMPER_M = 4.0;
3
4 const int MIN_RUN_LENGTH = 3;
5
6 double champernownePDF(double x){
7     return CHAMPER_A / (M_PI * std::cosh(CHAMPER_A * (x - CHAMPER_M)));
8 }
9
10 std::string generateRandomText(size_t approx_target_chars, const std::string&
    filename = "random_text_custom_dist.txt") {
11     std::vector<std::string> base_char_pool;
12     std::vector<std::string> RUS_LOWER = {
13         "а", "б", "в", "г", "д", "е", "ж", "з", "и", "й", "к", "л", "м",
14         "н", "о", "п", "р", "с", "т", "у", "ф", "х", "ц", "ч", "ш", "щ",
15         "ы", "э", "ю", "я"
16     };
17     std::vector<std::string> DIGITS = {"0", "1", "2", "3", "4", "5", "6", "7", "8", "9"};
18     std::vector<std::string> SPECIALS = {" ", ".", ",", "!", "?", "-", ":"};
19
20     for (const auto& s : RUS_LOWER) base_char_pool.push_back(s);
21     for (const auto& s : DIGITS) base_char_pool.push_back(s);
22     for (const auto& s : SPECIALS) base_char_pool.push_back(s);
23
24     std::vector<double> weights;
25     weights.reserve(base_char_pool.size());
26     double total_calculated_weight = 0.0;
27
```

```

28     for (size_t i = 0; i < base_char_pool.size(); ++i) {
29         double x_for_char = static_cast<double>(i);
30         double weight = champernownePDF(x_for_char);
31         weight = std::max(weight, 1e-9);
32         weights.push_back(weight);
33         total_calculated_weight += weight;
34     }
35
36     std::random_device rd;
37     std::mt19937 rng(rd());
38     std::discrete_distribution<size_t> char_dist(weights.begin(), weights.end());
39
40     std::uniform_real_distribution<double> run_probability_dist(0.0, 1.0);
41     double probability_of_run = 0.20;
42     std::uniform_int_distribution<size_t> run_length_dist(MIN_RUN_LENGTH, 10);
43
44     std::string result_text;
45     result_text.reserve(approx_target_chars * 2);
46     size_t current_char_count = 0;
47
48     while (current_char_count < approx_target_chars) {
49         size_t char_idx = char_dist(rng);
50         const std::string& selected_char_utf8 = base_char_pool[char_idx];
51
52         if (run_probability_dist(rng) < probability_of_run && (current_char_count +
53             MIN_RUN_LENGTH < approx_target_chars) ) {
54             size_t length_of_run = run_length_dist(rng);
55             length_of_run = std::min(length_of_run, approx_target_chars -
56                 current_char_count);
57
58             if (length_of_run >= MIN_RUN_LENGTH ){
59                 for (size_t k_run = 0; k_run < length_of_run; ++k_run) {
60                     result_text += selected_char_utf8;
61                 }
62                 current_char_count += length_of_run;
63             } else {
64                 result_text += selected_char_utf8;
65                 current_char_count++;
66             }
67         } else {
68             result_text += selected_char_utf8;
69             current_char_count++;
70         }
71     }

```

```

70
71     std::ofstream outfile(filename, std::ios::binary);
72     if (outfile.is_open()) {
73         outfile.write(result_text.data(), result_text.length());
74         outfile.close();
75         std::cout << "Сгенерирован файл '" << filename << "' (сгенерировано " <<
            current_char_count
76             << " UTF-8 симв., размер файла " << result_text.length() << " байт)
            с использованием champernownePDF." << std::endl;
77     } else {
78         std::cerr << "Не удалось создать файл для генерации: " << filename << std::
            endl;
79     }
80     return result_text;
81 }

```

Листинг 31: Функция generateRandomText в RLE (ключевые моменты)

4.2 Функция RLE::advancedRleEncode

- **Вход:** `input: const std::string&` – исходная строка байт для кодирования.
- **Выход:** `std::string` – RLE-закодированная строка.
- **Описание:** Реализует продвинутый алгоритм RLE. Функция последовательно обрабатывает входную строку. Если обнаруживается серия одинаковых байт длиной не менее `MIN_RUN_LENGTH`, она кодируется в формате «счетчик>#<байт>». Если длинная серия отсутствует, функция определяет максимальную возможную последовательность литералов (байтов, не образующих длинных серий) и кодирует ее в формате <длина>#<байты_литералов>». Символ '#' используется как разделитель между числовым префиксом и данными.

```

1  const int MIN_RUN_LENGTH = 3;
2
3  std::string advancedRleEncode(const std::string& input) {
4      if (input.empty()) return "";
5
6      std::stringstream encoded_ss;
7      size_t i = 0;
8      const size_t n = input.length();
9      const char SEPARATOR = '#';
10
11     while (i < n) {
12         char current_char = input[i];
13         size_t count = 1;

```

```

14     size_t j = i + 1;
15     while (j < n && input[j] == current_char) {
16         count++;
17         j++;
18     }
19
20     if (count >= MIN_RUN_LENGTH) {
21         encoded_ss << std::to_string(count) << SEPARATOR << current_char;
22         i = j;
23     } else {
24         size_t literal_start = i;
25         size_t k = i;
26         while(k < n) {
27             char check_char_literal = input[k];
28             size_t run_len_ahead = 0;
29             size_t temp_k = k;
30             while(temp_k < n && input[temp_k] == check_char_literal) {
31                 run_len_ahead++;
32                 temp_k++;
33             }
34             if (run_len_ahead >= MIN_RUN_LENGTH) {
35                 break;
36             }
37             k += run_len_ahead;
38         }
39         size_t literal_length = k - literal_start;
40         if (literal_length > 0) {
41             encoded_ss << "-" << std::to_string(literal_length) << SEPARATOR;
42             encoded_ss << input.substr(literal_start, literal_length);
43         }
44         i = k;
45     }
46 }
47 return encoded_ss.str();
48 }

```

Листинг 32: Функция advancedRleEncode в RLE

4.3 Функция RLE::advancedRleDecode

- **Вход:** encoded_input: const std::string& – RLE-закодированная строка.
- **Выход:** std::string – раскодированная (исходная) строка.
- **Описание:** Выполняет декодирование строки, сжатой функцией advancedRleEncode. Ите-

рирует по закодированной строке, ожидая определенный формат: «счетчик>#<байт>» для серий повторов или <длина>#<байты_литералов>» для последовательностей литералов. Функция считывает управляющие символы ('-' или цифру), затем число (счетчик/длинну), затем символ-разделитель '#'. После этого, в зависимости от типа сегмента, либо повторяет следующий байт указанное число раз, либо считывает указанное количество байт как литералы. Содержит проверки на корректность формата и наличие данных, выбрасывая исключения `std::runtime_error` при ошибках.

```

1 std::string advancedRleDecode(const std::string& encoded_input) {
2     if (encoded_input.empty()) return "";
3
4     std::stringstream decoded_ss;
5     size_t i = 0;
6     const size_t n = encoded_input.length();
7     const char SEPARATOR = '#';
8
9     while (i < n) {
10         char first_byte_of_seq = encoded_input[i];
11         bool is_negative_run = false;
12
13         if (first_byte_of_seq == '-') {
14             is_negative_run = true;
15             i++;
16             if (i >= n || !std::isdigit(static_cast<unsigned char>(encoded_input[i]))
17                 ) {
18                 throw std::runtime_error("RLE Decode: '-' not followed by a digit or
19                     EOF.");
20             }
21         } else if (!std::isdigit(static_cast<unsigned char>(first_byte_of_seq))) {
22             throw std::runtime_error("RLE Decode: Sequence does not start with '-' or
23                 digit. Found: '" + std::string(1, first_byte_of_seq) + "' at position
24                 " + std::to_string(i));
25         }
26
27         std::string num_str;
28         while (i < n && std::isdigit(static_cast<unsigned char>(encoded_input[i]))) {
29             num_str += encoded_input[i];
30             i++;
31         }
32
33         if (i >= n || encoded_input[i] != SEPARATOR) {
34             throw std::runtime_error("RLE Decode: Missing separator '" + std::string
35                 (1, SEPARATOR) + "' after number at pos ~" + std::to_string(i) + " (

```



```

        num_str was " + num_str + "));
31     }
32     i++;
33
34     if (num_str.empty()) {
35         throw std::runtime_error("RLE Decode: Failed to read number for length/
        count.");
36     }
37
38     int count_or_length = 0;
39     try {
40         count_or_length = std::stoi(num_str);
41     } catch (const std::out_of_range& oor) {
42         throw std::runtime_error("RLE Decode: Number '" + num_str + "' out of
        range for int.");
43     } catch (const std::invalid_argument& ia) {
44         throw std::runtime_error("RLE Decode: Number '" + num_str + "' is not a
        valid integer.");
45     }
46
47     if (count_or_length <= 0) {
48         throw std::runtime_error("RLE Decode: Invalid count/length (<=0): " +
        std::to_string(count_or_length));
49     }
50
51     if (is_negative_run) {
52         if (i + count_or_length > n) {
53             throw std::runtime_error("RLE Decode: Not enough data for literal
        sequence. Expected " + std::to_string(count_or_length) + ",
        available " + std::to_string(n-i));
54         }
55         decoded_ss << encoded_input.substr(i, count_or_length);
56         i += count_or_length;
57     } else {
58         if (i >= n) {
59             throw std::runtime_error("RLE Decode: Missing char_to_repeat after
        count and separator.");
60         }
61         char char_to_repeat = encoded_input[i];
62         i++;
63         for (int k_rep = 0; k_rep < count_or_length; ++k_rep) {
64             decoded_ss << char_to_repeat;
65         }
66     }

```

```

67     }
68     return decoded_ss.str();
69 }

```

Листинг 33: Функция `advancedRleDecode` в RLE

5 Реализация кодирования Фано

Для второго этапа сжатия (или как самостоятельный метод) был реализован алгоритм статистического кодирования Фано. Ниже приведено описание ключевых структур и функций.

5.1 Структура `FanoNode`

Прототип структуры:

```

1 namespace Fano {
2 struct FanoNode {
3     char symbol;
4     double probability;
5
6     FanoNode(char sym, double prob) : symbol(sym), probability(prob) {}
7 };
8 }

```

Листинг 34: Структура `FanoNode`

- **Члены структуры:**

- `symbol`: Тип `char`. Хранит символ (байт) из обрабатываемого текста.
- `probability`: Тип `double`. Хранит вычисленную вероятность (или нормализованную частоту) появления данного символа.

- **Конструктор:** `FanoNode(char sym, double prob)` инициализирует символ и его вероятность.

- **Назначение:** Данная структура используется для временного представления символов и их вероятностей, в частности, для их последующей сортировки и обработки в процессе построения кодов Шеннона-Фано.

5.2 Функция `generateFanoCodes`

- **Вход:**

- `nodes: std::vector<FanoNode*>&` – ссылка на вектор указателей на узлы `FanoNode`, отсортированные по убыванию вероятностей.
- `start, end: int` – индексы, определяющие текущий диапазон узлов в векторе `nodes` для обработки.
- `prefix: std::string` – текущий двоичный префикс кода, сформированный на предыдущих шагах рекурсии.
- `codes: std::map<char, std::string>&` – ссылка на карту, в которую будут записаны результирующие коды (символ \rightarrow строка кода).

- **Выход:** модифицированная карта `codes`.
- **Описание:** Это основная рекурсивная функция для построения кодов Шеннона-Фано. Для заданного диапазона отсортированных символов она находит точку разделения `splitIndex`, которая делит диапазон на две подгруппы с суммарными вероятностями, максимально близкими друг к другу. Затем функция рекурсивно вызывается для первой подгруппы (добавляя '0' к префиксу) и для второй подгруппы (добавляя '1' к префиксу). Базовый случай рекурсии — когда в диапазоне остается один символ; тогда текущий префикс становится его кодом. Реализована обработка случая, когда все вероятности в подгруппе равны нулю, и защита от пустых диапазонов.

```

1 void generateFanoCodes(std::vector<FanoNode*>& nodes, int start, int end, std::string
    prefix, std::map<char, std::string>& codes) {
2     if (start > end) {
3         return;
4     }
5     if (start == end) {
6         codes[nodes[start]->symbol] = (prefix.empty() && nodes.size() == 1) ? "0" :
            prefix;
7         return;
8     }
9
10    double totalWeight = 0;
11    for (int i = start; i <= end; ++i) {
12        totalWeight += nodes[i]->probability;
13    }
14
15    if (totalWeight == 0) {
16        for(int i = start; i <= end; ++i) {
17            codes[nodes[i]->symbol] = prefix + ( (i-start) % 2 == 0 ? "0" : "1");
18        }
19        return;

```

```

20 }
21
22 double partialWeight = 0;
23 int splitIndex = start;
24 for (; splitIndex < end; ++splitIndex) {
25     if (partialWeight + nodes[splitIndex]->probability > totalWeight / 2.0 &&
26         partialWeight > 0) {
27         splitIndex--;
28         break;
29     }
30     partialWeight += nodes[splitIndex]->probability;
31     if (partialWeight >= totalWeight / 2.0) {
32         break;
33     }
34 }
35 if (splitIndex < start && start < end) splitIndex = start;
36 generateFanoCodes(nodes, start, splitIndex, prefix + "0", codes);
37 if (splitIndex + 1 <= end) {
38     generateFanoCodes(nodes, splitIndex + 1, end, prefix + "1", codes);
39 }

```

Листинг 35: Функция generateFanoCodes

5.3 Функция buildFanoCodes

- **Вход:** `text: const std::string&` – входная строка, для символов которой строятся коды.
- **Выход:** `std::map<char, std::string>` – карта "символ \rightarrow код Фано".
- **Описание:** Функция оркестрирует процесс создания кодовой таблицы Шеннона-Фано. Она выполняет подсчет частот символов во входном тексте, преобразует их в вероятности, создает для каждого символа узел `FanoNode`, сортирует эти узлы по убыванию вероятностей, а затем вызывает рекурсивную функцию `generateFanoCodes` для построения самих кодов. В конце освобождается память, выделенная для узлов `FanoNode`.

```

1 std::map<char, std::string> buildFanoCodes(const std::string& text) {
2     std::map<char, double> frequencies;
3     if (text.empty()) {
4         return {};
5     }
6     for (char c : text) {
7         frequencies[c]++;
8     }

```

```

9
10     std::vector<FanoNode*> nodes;
11     nodes.reserve(frequencies.size());
12     for (auto& pair : frequencies) {
13         nodes.push_back(new FanoNode(pair.first, pair.second / text.size()));
14     }
15
16     std::sort(nodes.begin(), nodes.end(), [](FanoNode* a, FanoNode* b) {
17         return a->probability > b->probability;
18     });
19
20     std::map<char, std::string> codes;
21     if (!nodes.empty()) {
22         generateFanoCodes(nodes, 0, static_cast<int>(nodes.size()) - 1, "", codes);
23     }
24
25     for (FanoNode* node_ptr : nodes) {
26         delete node_ptr;
27     }
28
29     return codes;
30 }

```

Листинг 36: Функция buildFanoCodes

5.4 Функция encodeFano

- **Вход:**

- `text: const std::string&` – исходный текст для кодирования.
- `codes: const std::map<char, std::string>&` – таблица кодов Фано.

- **Выход:** `std::string` – строка, состоящая из символов '0' и '1', представляющая закодированный текст.

- **Описание:** Кодировать входной текст `text`, заменяя каждый символ его соответствующим кодом Шеннона-Фано из таблицы `codes`. Результатом является строка, представляющая последовательность битов. Данная реализация не выполняет упаковку битов в байты.

```

1 std::string encodeFano(const std::string& text, const std::map<char, std::string>&
  codes) {
2     std::string encoded_bit_string;
3     if (text.empty() || codes.empty()) return encoded_bit_string;
4

```

```

5   encoded_bit_string.reserve(text.length() * 4);
6
7   for (char c : text) {
8       auto it = codes.find(c);
9       if (it != codes.end()) {
10          encoded_bit_string += it->second;
11      } else {
12          std::cerr << "ПРЕДУПРЕЖДЕНИЕ (encodeFano): Код для символа '" << c
13              << "' (ASCII: " << static_cast<int>(static_cast<unsigned char
14              >(c))
15              << ") не найден!" << std::endl;
16      }
17  }
18  return encoded_bit_string;
19 }

```

Листинг 37: Функция encodeFano

5.5 Функция decodeFano

- **Вход:**

- encoded_bit_string: const std::string& – строка из '0' и '1', представляющая закодированные данные.
- codes: const std::map<char, std::string>& – таблица кодов Фано.

- **Выход:** std::string – исходный текст.

- **Описание:** Декодирует строку encoded_bit_string. Сначала создает обратную карту кодов ("код" → "символ") для удобства поиска. Затем итерирует по "битовой" строке, накапливая биты в буфер. Как только содержимое буфера совпадает с одним из кодов в обратной карте, соответствующий символ добавляется к результату, и буфер очищается.

```

1  std::string decodeFano(const std::string& encoded_bit_string, const std::map<char,
2  std::string>& codes) {
3
4      std::map<std::string, char> reversedCodes;
5      for (auto& pair : codes) {
6          if (!pair.second.empty()) {
7              reversedCodes[pair.second] = pair.first;
8          }
9      }
10 }

```

```

11     std::string decoded_text;
12     std::string current_code_buffer;
13     decoded_text.reserve(encoded_bit_string.length());
14
15     for (char bit_char : encoded_bit_string) {
16         current_code_buffer += bit_char;
17         auto it = reversedCodes.find(current_code_buffer);
18         if (it != reversedCodes.end()) {
19             decoded_text += it->second;
20             current_code_buffer.clear();
21         }
22     }
23     if (!current_code_buffer.empty()) {
24         std::cerr << "ПРЕДУПРЕЖДЕНИЕ (decodeFano): Остались необработанные биты в бу
                фере: "
25                 << current_code_buffer << std::endl;
26     }
27     return decoded_text;
28 }

```

Листинг 38: Функция decodeFano

6 Вспомогательные функции обработки текста

6.1 Глобальная функция processTextToWords

- **Вход:** `text_utf8: const std::string&` – входной текст в кодировке UTF-8.
- **Выход:** `std::vector<std::string>` – вектор строк, где каждая строка представляет собой выделенное слово.
- **Описание:** Функция выполняет токенизацию (разбиение на слова) входной текстовой строки. Разделителями слов считаются пробельные символы (проверяемые `std::isspace`) и стандартные знаки препинания ASCII (проверяемые `std::ispunct`). Функция итерирует по байтам входной строки, накапливая последовательности байтов, не являющихся разделителями, во временную строку. При встрече разделителя или по достижении конца строки, накопленная непустая строка добавляется в результирующий вектор слов. Функция корректно обрабатывает многобайтовые UTF-8 символы, сохраняя их байтовые последовательности внутри выделенных слов.

```

1 std::vector<std::string> processTextToWords(const std::string& text_utf8) {
2     std::vector<std::string> words;
3     std::string current_word;

```

```

4   for (char c_byte : text_utf8) {
5       unsigned char u_c_byte = static_cast<unsigned char>(c_byte);
6       if (std::isspace(u_c_byte) || std::ispunct(u_c_byte)) {
7           if (!current_word.empty()) {
8               words.push_back(current_word);
9               current_word.clear();
10          }
11      } else {
12          current_word += c_byte;
13      }
14  }
15  if (!current_word.empty()) {
16      words.push_back(current_word);
17  }
18  return words;
19 }

```

Листинг 39: Функция processTextToWords

6.2 Глобальная функция readFileToString

- **Вход:** filepath: const std::string& – путь к файлу.
- **Выход:** std::string – содержимое файла в виде строки.
- **Описание:** Считывает все содержимое указанного файла в строку. Файл открывается в бинарном режиме (std::ios::binary) для корректной обработки любых данных, включая UTF-8. Если открытие файла не удалось, генерируется исключение.

```

1  std::string readFileToString(const std::string& filepath) {
2      std::ifstream file(filepath, std::ios::in | std::ios::binary);
3      if (!file.is_open()) {
4          throw std::runtime_error("Не удалось открыть файл: " + filepath);
5      }
6      std::stringstream buffer;
7      buffer << file.rdbuf();
8      return buffer.str();
9  }

```

Листинг 40: Функция readFileToString

7 Результаты работы

--- Главное Меню ---

1. Работать со словарем на Хеш-таблице
2. Работать со словарем на Красно-Черном дереве
3. RLE кодирование/декодирование текста
0. Выход

Ваш выбор: 1

--- Меню Словаря (Хеш-таблица) ---

1. Добавить слово
2. Удалить слово
3. Найти слово
4. Загрузить словарь из файла (перезаписать)
5. Дополнить словарь из файла
6. Очистить словарь
7. Показать текущее содержимое словаря (стандартный print)
8. Визуализировать структуру
0. Вернуться в главное меню

Ваш выбор: █

Рис. 3. Главное меню.

--- Меню Словаря (Хеш-таблица) ---

1. Добавить слово
2. Удалить слово
3. Найти слово
4. Загрузить словарь из файла (перезаписать)
5. Дополнить словарь из файла
6. Очистить словарь
7. Показать текущее содержимое словаря (стандартный print)
8. Визуализировать структуру
0. Вернуться в главное меню

Ваш выбор: 3

Введите слово для поиска: мой

Слово 'мой' (ключ: 'мой') найдено, частота: 1

Рис. 4. Поиск слова.

```

--- Меню Словаря (Хеш-таблица) ---
1. Добавить слово
2. Удалить слово
3. Найти слово
4. Загрузить словарь из файла (перезаписать)
5. Дополнить словарь из файла
6. Очистить словарь
7. Показать текущее содержимое словаря (стандартный print)
8. Визуализировать структуру
0. Вернуться в главное меню
Ваш выбор: 8
Визуализация Хеш-таблицы (размер: 101, элементы: 56):
Корзина [0]: <пусто>
Корзина [1]: ("Печально": 1)
Корзина [2]: <пусто>
Корзина [3]: ("забавлять": 1)
Корзина [4]: <пусто>
Корзина [5]: <пусто>
Корзина [6]: ("Ему": 1)
Корзина [7]: ("себя": 2)
Корзина [8]: <пусто>
Корзина [9]: ("наука": 1)
Корзина [10]: <пусто>
Корзина [11]: <пусто>
Корзина [12]: <пусто>
Корзина [13]: <пусто>
Корзина [14]: ("Его": 1) -> ("низкое": 1)
Корзина [15]: ("Мой": 1)
Корзина [16]: ("Он": 1)
Корзина [17]: ("какая": 1) -> ("Не": 1) -> ("думать": 1)
Корзина [18]: <пусто>
Корзина [19]: <пусто>
Корзина [20]: ("уважать": 1)
Корзина [21]: <пусто>
Корзина [22]: ("Вздохнуть": 1)
Корзина [23]: ("заставил": 1)
Корзина [24]: <пусто>
Корзина [25]: <пусто>
Корзина [26]: <пусто>
Корзина [27]: ("день": 1)
Корзина [28]: <пусто>
Корзина [29]: ("скука": 1)
Корзина [30]: ("отходя": 1)
Корзина [31]: <пусто>
Корзина [32]: ("больным": 1)
Корзина [33]: ("мог": 1) -> ("подушки": 1)
Корзина [34]: <пусто>
Корзина [35]: <пусто>
Корзина [36]: ("самых": 1)
Корзина [37]: ("шагу": 1)

```

Рис. 5. Отображение хэш-таблицы.

8. Визуализировать структуру

0. Вернуться в главное меню

Ваш выбор: 8

Визуализация Красно-Черного Древа:

```

      ("шутку":1 B)
    ("честных":1 B)
      ("шагу":1 R)
        ("черт":1 R)
        ("уважать":1 B)
        ("тебя":1 R)
      ("скука":1 R)
        ("сидеть":1 R)
        ("себя":2 B)
    ("самых":1 B)
      ("про":1 B)
      ("правил":1 B)
        ("прочь":1 R)
        ("пример":1 R)
        ("поправлять":1 B)
        ("подушки":1 R)
        ("подносить":1 R)
        ("отходя":1 B)
    ("ночь":1 R)
      ("ни":1 B)
      ("не":2 B)
      ("наука":1 B)
      ("мой":1 R)
    ("мог":1 B)
      ("лучше":1 R)
      ("лекарство":1 B)
      ("коварство":1 R)
      ("какая":1 R)
      ("и":3 B)
    ("заставил":1 B)
      ("занемог":1 R)
      ("забавлять":1 B)
      ("же":1 R)
    ("дядя":1 B)
      ("думать":1 R)
      ("другим":1 B)
      ("день":1 R)
      ("выдумать":1 B)
      ("возьмет":1 R)
    ("в":1 B)
      ("больным":1 B)
      ("боже":1 R)
      ("С":1 R)
      ("Полуживого":1 B)
      ("Печально":1 R)
    ("Он":1 R)
      ("Но":1 B)
      ("Не":1 R)
      ("Мой":1 B)
      ("Когда":2 B)
      ("Какое":1 R)
      ("И":1 R)
      ("Ему":1 R)
      ("Его":1 B)
      ("Вздыхать":1 R)
```

Конец визуализации.

Рис. 6. Отображение КЧ-дерева.

```

--- Главное Меню ---
1. Работать со словарем на Хеш-таблице
2. Работать со словарем на Красно-Черном дереве
3. RLE кодирование/декодирование текста
0. Выход
Ваш выбор: 3

--- Меню RLE и Статистического Сжатия ---
1. Одноступенчатый RLE (генерация текста)
2. Двухступенчатый RLE -> RLE (генерация текста)
3. Двухступенчатый RLE -> Фано (генерация текста)
4. RLE для файла 'sample_text_rus.txt'
5. Фано для файла 'sample_text_rus.txt'
0. Вернуться в главное меню
Ваш выбор: 1
Сгенерирован файл 'debug_rle_generated.txt' (сгенерировано 10000 UTF-8 симв., размер файла 10000 байт) с использованием champernownePDF.

--- ОТЛАДКА RLE: ОДНОСТУПЕНЧАТЫЙ ---
Размер исходного текста: 10000 байт (10000 UTF-8 символов было запрошено).
Исходный сгенерированный текст:

Размер закодированного текста: 7877 байт.
Закодированный текст:
коэф сжатия : 1.26952

Размер декодированного текста: 10000 байт.
Проверка RLE: Декодирование ВЕРНО.
--- КОНЕЦ ОТЛАДКИ RLE ---

--- Меню RLE и Статистического Сжатия ---
1. Одноступенчатый RLE (генерация текста)
2. Двухступенчатый RLE -> RLE (генерация текста)
3. Двухступенчатый RLE -> Фано (генерация текста)
4. RLE для файла 'sample_text_rus.txt'
5. Фано для файла 'sample_text_rus.txt'
0. Вернуться в главное меню
Ваш выбор: 2
Сгенерирован файл 'random_rle_test_2stage.txt' (сгенерировано 10000 UTF-8 симв., размер файла 10000 байт) с использованием champernownePDF.
Исходный текст: дддгддддддддгзددггдгдддддддгдддддддддгдддддддддддгггг...
1-й этап RLE: 3#д-1#г10#д-8#гзددггдг6#д-5#еддгг8#д-1#е3#д9#г-1#е...
коэф сжатия 1: 1.27992
2-й этап RLE: -7813#3#д-1#г10#д-8#гзددггдг6#д-5#еддгг8#д-1#е3#д9...
коэф сжатия 2: 0.999233
Размер исходного: 10000, после 1-го этапа: 7813, после 2-го этапа: 7819
коэф сжатия общ: 1.27894
Проверка двухступенчатого RLE: Декодирование ВЕРНО.

--- Меню RLE и Статистического Сжатия ---
1. Одноступенчатый RLE (генерация текста)
2. Двухступенчатый RLE -> RLE (генерация текста)
3. Двухступенчатый RLE -> Фано (генерация текста)
4. RLE для файла 'sample_text_rus.txt'
5. Фано для файла 'sample_text_rus.txt'
0. Вернуться в главное меню
Ваш выбор: 3
Сгенерирован файл 'debug_rle_fano.txt' (сгенерировано 10000 UTF-8 симв., размер файла 10000 байт) с использованием champernownePDF.

--- Тест: Двухступенчатый RLE -> Фано ---
Размер исходного текста: 10000 байт.
Размер после RLE: 8002 байт.
Применение Фано к результату RLE...
Размер после Фано (длина битовой строки): 29128 символов ('0'/'1').
Условный коэфф. сжатия Ф поверх RLE (биты RLE / биты Ф): 2.20
Общий УСЛОВНЫЙ коэфф. сжатия (биты оригинала / биты Ф): 0.69
Декодирование...
Декодирование Фано -> RLE: ВЕРНО.
Полное декодирование RLE -> Фано -> RLE -> Оригинал: ВЕРНО.

--- Меню RLE и Статистического Сжатия ---
1. Одноступенчатый RLE (генерация текста)
2. Двухступенчатый RLE -> RLE (генерация текста)
3. Двухступенчатый RLE -> Фано (генерация текста)
4. RLE для файла 'sample_text_rus.txt'
5. Фано для файла 'sample_text_rus.txt'
0. Вернуться в главное меню
Ваш выбор:

```

Рис. 7. Операции с функциями сжатия.

Заключение

В ходе выполнения работы были успешно реализованы структуры данных хеш-таблица и красно-черное дерево, а также словари на их основе с поддержкой всех заявленных операций. Реализован продвинутый алгоритм сжатия RLE.

Анализ RLE показал, что его эффективность сильно зависит от характера входных данных: для текстов с большим количеством повторов он может дать сжатие, тогда как для случайных данных или текстов естественного языка с низкой избыточностью повторов он может приводить к увеличению размера.

Достоинства

Главными достоинствами разработанной программы являются:

- **Реализация ключевых структур данных с нуля:** Программа демонстрирует глубокое понимание принципов работы хеш-таблиц и красно-черных деревьев, так как эти структуры были реализованы самостоятельно без использования стандартных контейнеров STL (за исключением базовых, таких как `std::vector` и `std::list` для внутренних нужд хеш-таблицы).
- **Интерактивный консольный интерфейс:** Наличие меню позволяет пользователю выбирать структуру данных для работы, выполнять операции со словарями и алгоритмом RLE, а также визуализировать внутреннее устройство хеш-таблицы и КЧ-дерева, что улучшает наглядность и удобство использования.
- **Обработка русского текста (UTF-8):** Приняты меры для базовой обработки текстов на русском языке в кодировке UTF-8 при чтении файлов и генерации случайного текста.
- **Структурированный и модульный код:** Код организован с использованием пространств имен для различных частей программы, что способствует его читаемости и поддержке.

Недостатки

К недостаткам текущей реализации можно отнести:

- **Ограниченная нормализация текста:** Приведение слов к нижнему регистру в словарях реализовано только для ASCII-символов. Для полноценной регистронезависимой обработки русского текста в UTF-8 потребовались бы более сложные решения или использование специализированных библиотек.
- **Консольная визуализация:** Хотя визуализация структур данных реализована, ее наглядность в консольном режиме ограничена, особенно для больших и сложных деревьев. Создание действительно "графического" представления дерева в консоли является нетривиальной задачей.

- **Производительность файловых операций и генерации текста:** Для очень больших файлов чтение всего содержимого в строку (`readFileToString`) или генерация больших случайных текстов символ за символом может быть не оптимальной с точки зрения производительности и использования памяти по сравнению с потоковыми методами обработки.
- **Двухэтапный алгоритм сжатия RLE:** Малоэффективен так как после первого этапа сжатия остается небольшой процент символов, которые можно сжать.

Масштабируемость:

Генерация текста: Текущий генератор случайного текста может быть расширен путем добавления новых символов, изменения их весов или логики генерации серий без изменения основной структуры функции. Производительность генерации очень больших текстов может потребовать оптимизации.

Возможные улучшения:

- Реализация полноценной поддержки UTF-8 для операций со строками в словарях (например, корректное приведение к нижнему регистру для русского языка).
- Рассмотрение возможности использования более эффективных форматов для хранения RLE-закодированных данных (например, бинарного формата вместо строкового представления чисел).
- Оптимизация файловых операций и генерации текста для работы с очень большими объемами данных (потоковая обработка).
- Улучшение консольной визуализации, например, путем более сложного форматирования или использования псевдографики для отображения ветвей дерева.
- Добавление измерения времени выполнения ключевых операций для сравнительного анализа производительности словарей на хеш-таблице и КЧ-дереве.

Список материалов

- [1] Новиков, Ф.А. Дискретная математика для программистов. – СПб.: Питер, 2008. – С. 384.
(дата обращения: 03.06.2025)
<https://stugum.files.wordpress.com/2014/03/novikov.pdf>

- [2] Шапоров С.Д. "Дискретная математика"(дата обращения: 03.06.2025)
<https://wiki.fenix.help/informatika/sdnf?ysclid=lr6bu660ud240614897>

- [3] Хаггарт Р. "Дискретная математика для программистов"(дата обращения: 24.05.2025)
<https://studizba.com/show/1019108-1-gorbatov-va-fundamentalnye-osnovy.html>