

# Лекция 11: Преобразование Бэрроуза–Уилера (BWT)

А. М. Шур

Кафедра алгебры и фундаментальной информатики УрФУ

2 мая 2020 г.

В 1994 году два ученых из Кембриджа — Дэвид Уилер и его бывший аспирант Майкл Бэрроуз — опубликовали статью

- M. Burrows, D.J. Wheeler. A block sorting lossless data compression algorithm. Tech. Report 124, Digital Equipment Corporation, 1994

о новом подходе к сжатию данных, основанном на несжимающем комбинаторном преобразовании, получившем позже название **BWT** (Burrows–Wheeler Transform)

В 1994 году два ученых из Кембриджа — Дэвид Уилер и его бывший аспирант Майкл Бэрроуз — опубликовали статью

- M. Burrows, D.J. Wheeler. A block sorting lossless data compression algorithm. Tech. Report 124, Digital Equipment Corporation, 1994

о новом подходе к сжатию данных, основанном на несжимающем комбинаторном преобразовании, получившем позже название **BWT** (Burrows–Wheeler Transform)



Слева — Майкл Бэрроуз, один из авторов первого интернет-поисковика AltaVista

Справа — Дэвид Уилер, один из пионеров программирования, изобретатель ассемблера и процедур

Перестановка множества  $\{1, \dots, n\}$  есть **биекция** этого множества на себя

- Перестановку записывают в виде строки/массива  $\theta[1..n]$  так, что  $\theta[i] = \theta(i)$ 
  - $\theta = 45321$  переводит 1 в 4, 2 в 5, 3 в 3, 4 в 2 и 5 в 1

**Перестановка** множества  $\{1, \dots, n\}$  есть **биекция** этого множества на себя

- Перестановку записывают в виде строки/массива  $\theta[1..n]$  так, что  $\theta[i] = \theta(i)$ 
  - $\theta = 45321$  переводит 1 в 4, 2 в 5, 3 в 3, 4 в 2 и 5 в 1

Перестановкой  $n$ -элементного множества можно **подействовать** на **любое  $n$ -элементное множество**, если его элементы естественным образом **линейно упорядочены**

- Например, можно действовать перестановкой
  - на строку (через номера позиций)
  - на граф (через номера вершин)
  - на матрицу (через номера строк или столбцов)
- ★ Формально, для последовательности объектов  $(y_1, \dots, y_n)$  **результат**  $\theta(y_1, \dots, y_n)$  **действия перестановки**  $\theta$  есть  $(y_{\theta(1)}, \dots, y_{\theta(n)})$

# Перестановка и действие перестановки

**Перестановка** множества  $\{1, \dots, n\}$  есть **биекция** этого множества на себя

- Перестановку записывают в виде строки/массива  $\theta[1..n]$  так, что  $\theta[i] = \theta(i)$ 
  - $\theta = 45321$  переводит 1 в 4, 2 в 5, 3 в 3, 4 в 2 и 5 в 1

Перестановкой  $n$ -элементного множества можно **подействовать** на **любое  $n$ -элементное множество**, если его элементы естественным образом **линейно упорядочены**

- Например, можно действовать перестановкой
  - на строку (через номера позиций)
  - на граф (через номера вершин)
  - на матрицу (через номера строк или столбцов)
- ★ Формально, для последовательности объектов  $(y_1, \dots, y_n)$  **результат**  $\theta(y_1, \dots, y_n)$  **действия перестановки**  $\theta$  есть  $(y_{\theta(1)}, \dots, y_{\theta(n)})$

## Пример

Если  $\theta = 769345218$ , то  $\theta(\text{ПОДСТРОКА}) = \text{КОСТРОПАД}$

Тот же результат получится при  $\theta = 729345618$

- ★ Любая перестановка превращает слово в его **анаграмму** (и обратно: любая анаграмма слова получается из него действием некоторой перестановки)

Преобразование **BWT** определяется следующим образом:

- ★ Алфавит упорядочен, строки можно сортировать лексикографически
- Взять текст  $T$  и выписать все  $|T|$  его циклических сдвигов в столбик, получив квадратную матрицу
- Отсортировать строки матрицы по возрастанию, получив матрицу  $M(T)$
- Вернуть последний столбец матрицы ( $\text{bwt}(T)$ )

Преобразование **BWT** определяется следующим образом:

- ★ Алфавит упорядочен, строки можно сортировать лексикографически
  - Взять текст  $T$  и выписать все  $|T|$  его циклических сдвигов в столбик, получив квадратную матрицу
  - Отсортировать строки матрицы по возрастанию, получив матрицу  $M(T)$
  - Вернуть последний столбец матрицы ( $\text{bwt}(T)$ )
- ★ Каждая буква в  $T$  является последней ровно для одного циклического сдвига
- ⇒  $\text{bwt}(T)$  состоит из тех же символов с теми же частотами, что и  $T$
- ⇒  $\text{bwt}(T)$  является **анаграммой**  $T$
- ⇒  $\text{bwt}$  как функция является **действием некоторой перестановки** (зависящей от  $T$ )
- ★  $\text{bwt}$  — не биекция, так как текстам, являющимся циклическими сдвигами друг друга, соответствует одна и та же матрица, а значит, один и тот же  $\text{bwt}$ -образ



Преобразование **BWT** определяется следующим образом:

- ★ Алфавит упорядочен, строки можно сортировать лексикографически
  - Взять текст  $T$  и выписать все  $|T|$  его циклических сдвигов в столбик, получив квадратную матрицу
  - Отсортировать строки матрицы по возрастанию, получив матрицу  $M(T)$
  - Вернуть последний столбец матрицы ( $\text{bwt}(T)$ )
- ★ Каждая буква в  $T$  является последней ровно для одного циклического сдвига
- ⇒  $\text{bwt}(T)$  состоит из тех же символов с теми же частотами, что и  $T$
- ⇒  $\text{bwt}(T)$  является **анаграммой**  $T$
- ⇒  $\text{bwt}$  как функция является **действием некоторой перестановки** (зависящей от  $T$ )
- ★  $\text{bwt}$  — не биекция, так как текстам, являющимся циклическими сдвигами друг друга, соответствует одна и та же матрица, а значит, один и тот же  $\text{bwt}$ -образ
- ... говорят, что Бэрроуз придумал это преобразование еще студентом, но не нашел, куда его приложить

# Преобразование Бэрроуза–Уилера – пример

Рассмотрим наш стандартный пример  $T = \text{БанБананана}\#$

- В utf-8 или cp1251 имеем  $\# < Б < а < н$ ; тогда

	#	Б	а	н	Б	а	н	а	н	а	н	а	н	а	
	Б	а	н	Б	а	н	а	н	а	н	а	н	а	н	#
	Б	а	н	а	н	а	н	а	#	Б	а	н			
	а	#	Б	а	н	Б	а	н	а	н	а	н			
	а	н	Б	а	н	а	н	а	н	а	#	Б			
$M(\text{БанБананана}\#) =$	а	н	а	#	Б	а	н	Б	а	н	а	н			
	а	н	а	н	а	#	Б	а	н	Б	а	н			
	а	н	а	н	а	н	а	#	Б	а	н	Б			
	н	а	#	Б	а	н	Б	а	н	а	н	а			
	н	а	н	а	#	Б	а	н	Б	а	н	а			
	н	а	н	а	н	а	#	Б	а	н	Б	а			
	н	Б	а	н	а	н	а	н	а	#	Б	а			

$\text{bwt}(\text{БанБананана}\#) = \text{а\#ннБннБаааа}$

# Преобразование Бэрроуза–Уилера – пример

Рассмотрим наш стандартный пример  $T = \text{БанБананана}\#$

- В utf-8 или cp1251 имеем  $\# < Б < а < н$ ; тогда

$$M(\text{БанБананана}\#) = \begin{array}{cccccccc} \# & Б & а & н & Б & а & н & а & н & а & н & а \\ Б & а & н & Б & а & н & а & н & а & н & а & \# \\ \hline Б & а & н & а & н & а & н & а & \# & Б & а & н \\ а & \# & Б & а & н & Б & а & н & а & н & а & н \\ а & н & Б & а & н & а & н & а & н & а & \# & Б \\ а & н & а & \# & Б & а & н & Б & а & н & а & н \\ а & н & а & н & а & н & а & \# & Б & а & н & Б \\ н & а & \# & Б & а & н & Б & а & н & а & н & а \\ н & а & н & а & \# & Б & а & н & Б & а & н & а \\ н & а & н & а & н & а & \# & Б & а & н & Б & а \\ н & Б & а & н & а & н & а & н & а & \# & Б & а \end{array}$$

$\text{bwt}(\text{БанБананана}\#) = \text{а}\#\text{ннБннБаааа}$

- ★ Начало каждой строки является **правым контекстом** в  $T$  для последнего символа этой строки
- ★ Символы, у которых совпадают длинные контексты, часто бывают равны
- ⇒ Если символы в  $T$  сильно зависят от контекста, в  $\text{bwt}(T)$  будут длинные последовательности одинаковых символов, которые легко сжимать

- Чтобы BWT хотя бы **теоретически** подходило для препроцессинга данных перед сжатием, нужно доказать, что  $T$  можно восстановить, зная  $\text{bwt}(T)$

Чтобы получить **практически** осмысленный алгоритм сжатия, нужно

- быстро строить  $\text{bwt}(T)$
- быстро восстанавливать  $T$  из  $\text{bwt}(T)$
- эффективно сжимать  $\text{bwt}(T)$  без потерь

- Чтобы BWT хотя бы **теоретически** подходило для препроцессинга данных перед сжатием, нужно доказать, что  $T$  можно восстановить, зная  $\text{bwt}(T)$

Чтобы получить **практически** осмысленный алгоритм сжатия, нужно

- быстро строить  $\text{bwt}(T)$
- быстро восстанавливать  $T$  из  $\text{bwt}(T)$
- эффективно сжимать  $\text{bwt}(T)$  без потерь

Вначале научимся решать задачу восстановления:

- Как уже упоминалось, разные тексты имеют одинаковый результат BWT, если являются циклическими сдвигами друг друга, а значит — строками одной матрицы
- С другой стороны, мы покажем, как по  $\text{bwt}(T)$  построить матрицу  $M(T)$
- Для однозначности восстановления  $T$  можно хранить номер строки в  $M(T)$ , которая соответствует  $T$ , **но удобнее дописать к  $T$  уникальный символ конца файла**, как мы сделали в примере; тогда номер строки  $T$  в  $M(T)$  равен позиции этого символа в  $\text{bwt}(T)$

## Восстановление матрицы $M(T)$

- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$

- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$
- Отсортировав  $\text{bwt}(T)$ , получим первый столбец  $M(T)$
- Если поставить первый столбец справа от последнего, в полученной матрице с двумя столбцами записаны все циклические подстроки  $w$  длины 2
  - Отсортировав эти подстроки, получим первые два столбца  $M(T)$
  - Приставив спереди последний столбец, получим циклические подстроки длины 3
- Будем повторять процедуру до заполнения матрицы:

# Восстановление матрицы $M(T)$

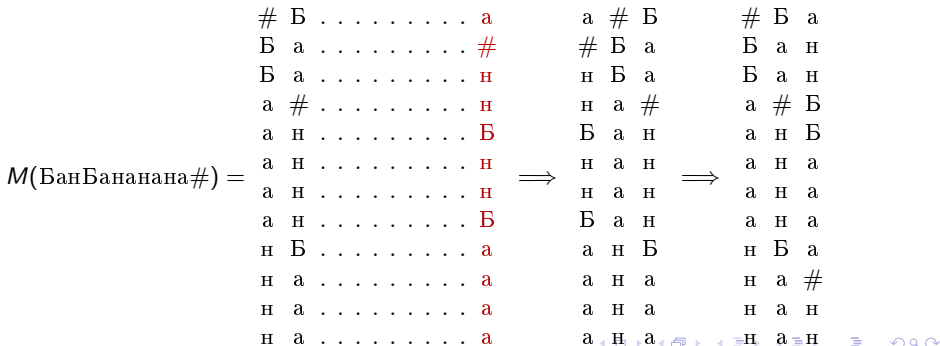
- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$
- Отсортировав  $\text{bwt}(T)$ , получим первый столбец  $M(T)$
- Если поставить первый столбец справа от последнего, в полученной матрице с двумя столбцами записаны все циклические подстроки  $w$  длины 2
  - Отсортировав эти подстроки, получим первые два столбца  $M(T)$
  - Приставив спереди последний столбец, получим циклические подстроки длины 3
- Будем повторять процедуру до заполнения матрицы:

$$\begin{array}{l}
 \# \dots\dots\dots a \\
 Б \dots\dots\dots \# \\
 Б \dots\dots\dots н \\
 а \dots\dots\dots н \\
 а \dots\dots\dots Б \\
 а \dots\dots\dots н \\
 а \dots\dots\dots н \\
 а \dots\dots\dots Б \\
 н \dots\dots\dots а \\
 н \dots\dots\dots а \\
 н \dots\dots\dots а \\
 н \dots\dots\dots а \\
 н \dots\dots\dots а
 \end{array}
 \Rightarrow
 \begin{array}{l}
 а \# \dots\dots\dots \# Б \\
 \# Б \dots\dots\dots Б а \\
 н Б \dots\dots\dots Б а \\
 н а \dots\dots\dots а \# \\
 Б а \dots\dots\dots а н \\
 н а \dots\dots\dots а н \\
 н а \dots\dots\dots а н \\
 Б а \dots\dots\dots а н \\
 а н \dots\dots\dots н Б \\
 а н \dots\dots\dots н а \\
 а н \dots\dots\dots н а \\
 а н \dots\dots\dots н а \\
 а н \dots\dots\dots а
 \end{array}
 \Rightarrow
 \begin{array}{l}
 \# Б \dots\dots\dots а \\
 Б а \dots\dots\dots а \\
 Б а \dots\dots\dots а \\
 а \# \dots\dots\dots а \\
 а н \dots\dots\dots а \\
 а н \dots\dots\dots а \\
 а н \dots\dots\dots а \\
 а н \dots\dots\dots а \\
 н Б \dots\dots\dots а \\
 н а \dots\dots\dots а \\
 н а \dots\dots\dots а \\
 н а \dots\dots\dots а \\
 а
 \end{array}$$



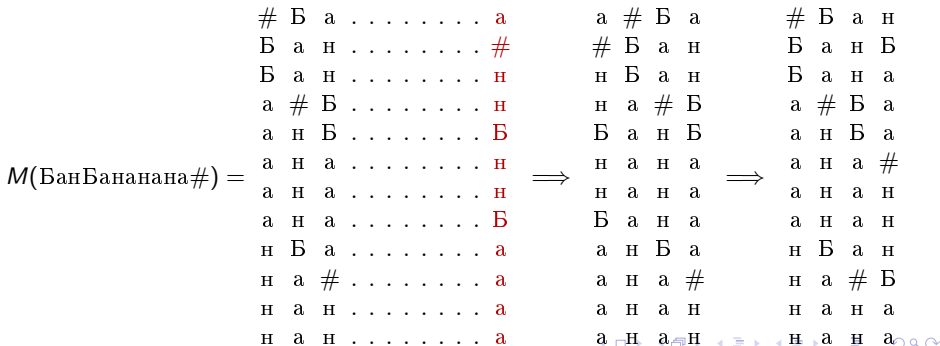
# Восстановление матрицы $M(T)$

- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$
- Отсортировав  $\text{bwt}(T)$ , получим первый столбец  $M(T)$
- Если поставить первый столбец справа от последнего, в полученной матрице с двумя столбцами записаны все циклические подстроки  $w$  длины 2
  - Отсортировав эти подстроки, получим первые два столбца  $M(T)$
  - Приставив спереди последний столбец, получим циклические подстроки длины 3
- Будем повторять процедуру до заполнения матрицы:



# Восстановление матрицы $M(T)$

- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$
- Отсортировав  $\text{bwt}(T)$ , получим первый столбец  $M(T)$
- Если поставить первый столбец справа от последнего, в полученной матрице с двумя столбцами записаны все циклические подстроки  $w$  длины 2
  - Отсортировав эти подстроки, получим первые два столбца  $M(T)$
  - Приставив спереди последний столбец, получим циклические подстроки длины 3
- Будем повторять процедуру до заполнения матрицы:



# Восстановление матрицы $M(T)$

- ★ В любом столбце матрицы  $M(T)$  записана анаграмма  $T$
- ★ В матрице, составленной из первых  $i$  столбцов  $M(T)$ , где  $1 \leq i \leq n$ , по строкам записаны, в лексикографическом порядке, все **циклические подстроки**  $T$ 
  - циклическая подстрока — это либо  $T[i..j]$ ,  $i \leq j$ , либо  $T[j..n]T[1..i]$ ,  $i < j$
- Отсортировав  $\text{bwt}(T)$ , получим первый столбец  $M(T)$
- Если поставить первый столбец справа от последнего, в полученной матрице с двумя столбцами записаны все циклические подстроки  $w$  длины 2
  - Отсортировав эти подстроки, получим первые два столбца  $M(T)$
  - Приставив спереди последний столбец, получим циклические подстроки длины 3
- Будем повторять процедуру до заполнения матрицы:

$M(\text{БанБананана}\#) =$

#	Б	а	н	Б	а	н	а	н	а	н	а
Б	а	н	Б	а	н	а	н	а	н	а	#
Б	а	н	а	н	а	н	а	#	Б	а	н
а	#	Б	а	н	Б	а	н	а	н	а	н
а	н	а	#	Б	а	н	Б	а	н	а	н
а	н	а	н	а	#	Б	а	н	Б	а	н
а	н	Б	а	н	а	н	а	н	#	Б	а
н	а	#	Б	а	н	Б	а	н	а	н	а
н	а	н	а	#	Б	а	н	Б	а	н	а
н	а	н	а	н	#	Б	а	н	Б	а	н
н	Б	а	н	а	н	а	#	Б	а	н	а

- Матрица имеет размер  $n^2$ , ее нельзя восстановить быстрее, чем за  $O(n^2)$
- Нам нужна не матрица, а одна строка; как восстановить ее за  $O(n)$ ?

- Матрица имеет размер  $n^2$ , ее нельзя восстановить быстрее, чем за  $O(n^2)$
- Нам нужна не матрица, а одна строка; как восстановить ее за  $O(n)$ ?

Сортировка строки  $w[1..n]$  над упорядоченным алфавитом — это **перестановка**  $\sigma$  такая, что строка  $\sigma(w)$  **отсортирована**:  $i < j$  влечет  $\sigma(w)[i] \leq \sigma(w)[j]$

- Если в  $w$  есть повторяющиеся буквы, то сортировок у него несколько
- Среди них одна **не меняет порядок следования одинаковых букв**
- ★ **Стабильная сортировка**  $\sigma_w$  слова  $w$  — это перестановка, удовлетворяющая условиям
  - $i < j$  влечет  $\sigma_w(w)[i] \leq \sigma_w(w)[j]$
  - $i < j$  и  $w[i] = w[j]$  влечет  $\sigma_w(i) < \sigma_w(j)$

- Матрица имеет размер  $n^2$ , ее нельзя восстановить быстрее, чем за  $O(n^2)$
- Нам нужна не матрица, а одна строка; как восстановить ее за  $O(n)$ ?

Сортировка строки  $w[1..n]$  над упорядоченным алфавитом — это **перестановка**  $\sigma$  такая, что строка  $\sigma(w)$  **отсортирована**:  $i < j$  влечет  $\sigma(w)[i] \leq \sigma(w)[j]$

- Если в  $w$  есть повторяющиеся буквы, то сортировок у него несколько
- Среди них одна **не меняет порядок следования одинаковых букв**
- ★ **Стабильная сортировка**  $\sigma_w$  слова  $w$  — это перестановка, удовлетворяющая условиям
  - $i < j$  влечет  $\sigma_w(w)[i] \leq \sigma_w(w)[j]$
  - $i < j$  и  $w[i] = w[j]$  влечет  $\sigma_w(i) < \sigma_w(j)$

Пример:  $\sigma = 419102111235678$  — стабильная сортировка для  $w = \text{а\#ннБннБаааа}$

- Матрица имеет размер  $n^2$ , ее нельзя восстановить быстрее, чем за  $O(n^2)$
- Нам нужна не матрица, а одна строка; как восстановить ее за  $O(n)$ ?

Сортировка строки  $w[1..n]$  над упорядоченным алфавитом — это **перестановка**  $\sigma$  такая, что строка  $\sigma(w)$  **отсортирована**:  $i < j$  влечет  $\sigma(w)[i] \leq \sigma(w)[j]$

- Если в  $w$  есть повторяющиеся буквы, то сортировок у него несколько
- Среди них одна **не меняет порядок следования одинаковых букв**
- ★ **Стабильная сортировка**  $\sigma_w$  слова  $w$  — это перестановка, удовлетворяющая условиям
  - $i < j$  влечет  $\sigma_w(w)[i] \leq \sigma_w(w)[j]$
  - $i < j$  и  $w[i] = w[j]$  влечет  $\sigma_w(i) < \sigma_w(j)$

Пример:  $\sigma = 419102111235678$  — стабильная сортировка для  $w = \text{а\#ннБннБаааа}$

- ★ Для быстрого восстановления текста  $T$  по  $\text{bwt}(T)$  нужна перестановка, **обратная к стабильной сортировке**  $\text{bwt}(T)$ 
  - напомним, что **обратная** к  $\sigma$  — перестановка  $\sigma^{-1}$  такая, что  $(\forall i) \sigma^{-1}(\sigma(i)) = i$

В Примере  $\sigma^{-1} = 258191011123467$

- ★ Если  $\sigma$  — стабильная сортировка  $\text{bwt}(T)$ , то  $\sigma(i)$  — позиция, которую занимает в первом столбце символ  $\text{bwt}(T)[i]$ 
  - ★ Значит,  $i$ -ый символ первого столбца — это  $\text{bwt}(T)[\sigma^{-1}(i)]$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$



## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
⇒  $T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
⇒  $T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
⇒ он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
⇒  $T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
⇒  $T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \mathbf{a\#nnBnnBa\text{а\text{а\text{а}}}}$ ,  $\sigma^{-1} = \mathbf{2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7}$

$$i = 2$$

T = . . . . . #

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{nnBnnBa} \text{aaaa}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(2) = 5 \quad \text{bwt}(T)[5] = \text{B}$$

$T = \text{B} \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad . \quad \#$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\#\text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(5) = 9 \quad \text{bwt}(T)[9] = \text{a}$$

$T = \text{Б а . . . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(9) = 3 \quad \text{bwt}(T)[3] = \text{н}$$

$T = \text{Б а н . . . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(3) = 8 \quad \text{bwt}(T)[8] = \text{Б}$$

$T = \text{Б а н Б . . . . . \#}$



## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\#\text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(8) = 12 \quad \text{bwt}(T)[12] = \text{a}$$

$T = \text{Б а н Б а . . . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\#\text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(12) = 7 \quad \text{bwt}(T)[7] = \text{н}$$

$T = \text{Б а н Б а н . . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(7) = 11 \quad \text{bwt}(T)[11] = \text{a}$$

$T = \text{Б а н Б а н а . . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\#_{\text{ннБннБа}}\text{аааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(11) = 6 \quad \text{bwt}(T)[6] = \text{н}$$

$T = \text{Б а н Б а н а н . . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{ннБннБааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(6) = 10 \quad \text{bwt}(T)[10] = \text{a}$$

$T = \text{Б а н Б а н а н а . . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \text{a}\# \text{ннБннБаааа}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(10) = 4 \quad \text{bwt}(T)[4] = \text{н}$$

$T = \text{Б а н Б а н а н а н . \#}$

## Восстановление текста за линейное время (2)

Восстановим  $T$  по  $\text{bwt}(T)$  и  $\sigma^{-1}$ :

- Номер  $i$  нужной строки — это позиция  $\#$  в  $\text{bwt}(T)$
- Первый символ  $i$ -ой строки =  $i$ -й символ первого столбца =  $\sigma^{-1}(i)$ -й символ последнего столбца  
 $\Rightarrow T[1] = \text{bwt}(T)[\sigma^{-1}(i)]$
- Второй символ  $i$ -ой строки следует за первым  
 $\Rightarrow$  он является первым в той строке, в которой первый является последним (т.е. в  $\sigma^{-1}(i)$ -ой), а последним — в  $\sigma^{-1}(\sigma^{-1}(i))$ -ой (т.е. в  $\sigma^{-2}(i)$ -ой)  
 $\Rightarrow T[2] = \text{bwt}(T)[\sigma^{-2}(i)]$
- Аналогично, третий символ  $i$ -ой строки является первым в строке, где последним является второй, и т.д.  
 $\Rightarrow T[k] = \text{bwt}(T)[\sigma^{-k}(i)], k = 1, \dots, n$

Пример:  $\text{bwt}(T) = \mathbf{a\#nnBnnBa\text{aaaa}}$ ,  $\sigma^{-1} = 2\ 5\ 8\ 1\ 9\ 10\ 11\ 12\ 3\ 4\ 6\ 7$

$$\sigma^{-1}(4) = 1 \quad \text{bwt}(T)[1] = \mathbf{a}$$

$T = \mathbf{B\ a\ n\ B\ a\ n\ a\ n\ a\ n\ a\ \#}$

- Легко посчитать  $\sigma^{-1}$  за время  $O(n)$ : надо посчитать массив частот символов, для каждого символа  $x$  вычислить диапазон в упорядоченном списке и сделать один проход по  $\text{bwt}(T)$   
★ значит, текст восстанавливается по  $\text{bwt}$  за время  $O(n)$

Для построения  $\text{bwt}(T)$  нужно отсортировать лексикографически все **циклические сдвиги**  $T$

- ★ Матрицу  $M(T)$  строить не нужно, достаточно получить массив  $S$ , в котором любой элемент  $S[i]$  равен **номеру строки в  $M(T)$ , занимаемой циклическим сдвигом  $T_i = T[i..n]T[1..i-1]$** 
  - тогда  $\text{bwt}[S[i]] = T[i-1]$  (последний символ строки  $T_i$ )



Для построения  $\text{bwt}(T)$  нужно отсортировать лексикографически все **циклические сдвиги**  $T$

- ★ Матрицу  $M(T)$  строить не нужно, достаточно получить массив  $S$ , в котором любой элемент  $S[i]$  равен **номеру строки в  $M(T)$ , занимаемой циклическим сдвигом  $T_i = T[i..n]T[1..i-1]$** 
  - тогда  $\text{bwt}[S[i]] = T[i-1]$  (последний символ строки  $T_i$ )
- ★ Наличие **специального символа конца строки** упрощает задачу сортировки:
  - считаем, что специальный символ (мы его обозначаем через  $\#$ ) **меньше любого символа**
  - результат сравнения  $T_i$  и  $T_j$  будет таким же, как для  $T[i..n]$  и  $T[j..n]$
  - ★ вместо циклических сдвигов достаточно отсортировать **суффиксы текста**

Для построения  $\text{bwt}(T)$  нужно отсортировать лексикографически все **циклические сдвиги**  $T$

- ★ Матрицу  $M(T)$  строить не нужно, достаточно получить массив  $S$ , в котором любой элемент  $S[i]$  равен **номеру строки в  $M(T)$ , занимаемой циклическим сдвигом  $T_i = T[i..n]T[1..i-1]$** 
  - тогда  $\text{bwt}[S[i]] = T[i-1]$  (последний символ строки  $T_i$ )
- ★ Наличие **специального символа конца строки** упрощает задачу сортировки:
  - считаем, что спецсимвол (мы его обозначаем через  $\#$ ) **меньше любого символа**
  - результат сравнения  $T_i$  и  $T_j$  будет таким же, как для  $T[i..n]$  и  $T[j..n]$
  - ★ вместо циклических сдвигов достаточно отсортировать **суффиксы текста**
- Итак,  $S$  — массив, в котором  $S[i]$  равен номеру суффикса  $T[i..n]$  в отсортированном списке всех суффиксов строки  $T$ 
  - такой массив называется **суффиксным массивом** строки  $T$
- Суффиксные массивы известны с 1990-х, это часто используемая в разных задачах структура данных, заслуживающая отдельной лекции
- Здесь только упомянем, что суффиксный массив можно построить за время  $O(n)$ , хотя до сих пор используют и более простые алгоритмы с временем работы  $O(n \log n)$ 
  - строить суффиксный массив умеют многие специализированные библиотеки

В предположении, что символы в тексте значительно зависят от своих контекстов,

- результат BWT состоит из большого количества кусков:
  - каждый кусок достаточно однороден (часто в нем вообще повторяется единственный символ)
  - между кусками — резкие переходы (смена контекста)
  - куски сжимаются хорошо, переходы ухудшают сжатие

В предположении, что символы в тексте значительно зависят от своих контекстов,

- результат BWT состоит из большого количества кусков:
  - каждый кусок достаточно однороден (часто в нем вообще повторяется единственный символ)
  - между кусками — резкие переходы (смена контекста)
  - куски сжимаются хорошо, переходы ухудшают сжатие
- ★ Имеет смысл кастомизировать порядок символов в алфавите, чтобы уменьшить число резких переходов
  - для текста на естественном языке символы разбивают на три группы (гласные, согласные, небуквенные символы), а внутри каждой группы упорядочивают в соответствии с особенностями языка

# Способы сжатия результата BWT

В предположении, что символы в тексте значительно зависят от своих контекстов,

- результат BWT состоит из большого количества кусков:
  - каждый кусок достаточно однороден (часто в нем вообще повторяется единственный символ)
  - между кусками — резкие переходы (смена контекста)
  - куски сжимаются хорошо, переходы ухудшают сжатие
- ★ Имеет смысл кастомизировать порядок символов в алфавите, чтобы уменьшить число резких переходов
  - для текста на естественном языке символы разбивают на три группы (гласные, согласные, небуквенные символы), а внутри каждой группы упорядочивают в соответствии с особенностями языка

## Способы сжатия:

- **MTF+ZLE+ARI:**
  - Вначале применяется еще одно обратимое несжимающее преобразование, которое преобразует последовательности равных символов в последовательности нулей и “перекашивает” статистику остальных символов;
  - На втором этапе последовательности нулей заменяются их длинами
  - На третьем применяется арифметическое кодирование в алфавите байтов
- **DC:**
  - Кодироваться расстояния между последовательными вхождениями одного символа
  - Сжатие в основном достигается за счет того, что последовательности одинаковых символов... не требуют кодирования вообще
- Про другие способы сжатия (и про названные тоже) можно прочитать, например, в статье S. Deorowicz. Second step algorithms in the Burrows–Wheeler compression algorithm. Software—Practice and Experience, 2002

# Метод стопки книг (MTF)

Преобразование **Move-To-Front** (MTF) читает входную строку слева направо и посимвольно меняет ее по следующим правилам:

- Алфавит **упорядочен**, все символы адресуются своими номерами в списке (первый элемент имеет номер 0)
- Очередной символ  $a$  заменяется на номер  $a$  в текущем списке, после чего
  - $a$  перемещается в начало списка (получает номер 0)
  - у каждого символа, **номер которого был меньше номера  $a$** , номер увеличивается на 1
- ★ Отсюда образ **стопки книг**: книгу  $a$  достали из стопки и положили наверх

Преобразование **Move-To-Front** (MTF) читает входную строку слева направо и посимвольно меняет ее по следующим правилам:

- Алфавит **упорядочен**, все символы адресуются своими номерами в списке (первый элемент имеет номер 0)
- Очередной символ  $a$  заменяется на номер  $a$  в текущем списке, после чего
  - $a$  перемещается в начало списка (получает номер 0)
  - у каждого символа, **номер которого был меньше номера  $a$** , номер увеличивается на 1
- ★ Отсюда образ **стопки книг**: книгу  $a$  достали из стопки и положили наверх
- ★ Если на вход преобразования MTF подать  $\text{bwt}(T)$ , получим два эффекта, полезных для сжатия:
  - основной: последовательности одинаковых символов заменяются на последовательности нулей
  - дополнительный: маленькие номера встречаются часто, а большие — очень редко (часто встречающиеся символы будут находиться близко к началу списка)

# Метод стопки книг (MTF)

Преобразование **Move-To-Front** (MTF) читает входную строку слева направо и посимвольно меняет ее по следующим правилам:

- Алфавит **упорядочен**, все символы адресуются своими номерами в списке (первый элемент имеет номер 0)
- Очередной символ  $a$  заменяется на номер  $a$  в текущем списке, после чего
  - $a$  перемещается в начало списка (получает номер 0)
  - у каждого символа, **номер которого был меньше номера  $a$** , номер увеличивается на 1
- ★ Отсюда образ **стопки книг**: книгу  $a$  достали из стопки и положили наверх
- ★ Если на вход преобразования MTF подать  $\text{bwt}(T)$ , получим два эффекта, полезных для сжатия:
  - основной: последовательности одинаковых символов заменяются на последовательности нулей
  - дополнительный: маленькие номера встречаются часто, а большие — очень редко (часто встречающиеся символы будут находиться близко к началу списка)

Пример:

$w =$	а	#	н	н	Б	н	н	Б	а	а	а	а
$\text{mtf}(w) =$	2	1	3	0	3	1	0	1	3	0	0	0
алфавит :	$\begin{bmatrix} \# \\ Б \\ а \\ н \end{bmatrix}$	$\begin{bmatrix} а \\ \# \\ Б \\ н \end{bmatrix}$	$\begin{bmatrix} \# \\ а \\ Б \\ н \end{bmatrix}$	$\begin{bmatrix} н \\ \# \\ а \\ Б \end{bmatrix}$	$\begin{bmatrix} н \\ \# \\ а \\ Б \end{bmatrix}$	$\begin{bmatrix} Б \\ н \\ \# \\ а \end{bmatrix}$	$\begin{bmatrix} н \\ Б \\ \# \\ а \end{bmatrix}$	$\begin{bmatrix} н \\ Б \\ \# \\ а \end{bmatrix}$	$\begin{bmatrix} Б \\ н \\ \# \\ а \end{bmatrix}$	$\begin{bmatrix} а \\ Б \\ н \\ \# \end{bmatrix}$	$\begin{bmatrix} а \\ Б \\ н \\ \# \end{bmatrix}$	$\begin{bmatrix} а \\ Б \\ н \\ \# \end{bmatrix}$



## Кодирование нулей (ZLE)

- ★ В тексте  $\text{mtf}(\text{bwt}(T))$  есть длинные повторы нулей, соответствующие контекстным зависимостям в  $T$ , в то время как повторы других символов при MTF могут возникнуть только случайно (а значит, они короткие)

На первом этапе сжатия кодируются повторы нулей

- ★ Лучше всего работает метод **Zero Length Encoding (ZLE)**; он хорошо подготавливает почву для арифметического кодирования на втором этапе

- ★ В тексте  $\text{mtf}(\text{bwt}(T))$  есть длинные повторы нулей, соответствующие контекстным зависимостям в  $T$ , в то время как повторы других символов при MTF могут возникнуть только случайно (а значит, они короткие)

На первом этапе сжатия кодируются повторы нулей

- ★ Лучше всего работает метод **Zero Length Encoding (ZLE)**; он хорошо подготавливает почву для арифметического кодирования на втором этапе

ZLE:

- сделаем 254 и 255 — два самых редких символа в  $\text{mtf}(\text{bwt}(T))$  — служебными
- каждую последовательность нулей заменим на ее длину, в качестве цифр 0 и 1 используя 254 и 255
- 254 и 255, если они встречались, заменяем на пару байт 0 254 (0 255)
- ★ **дополнительный трюк**: для последовательности из  $N$  нулей кодируем число  $N+1$  без первой цифры
  - $N+1 \geq 2$ : двоичная запись начинается с 1 и содержит еще хотя бы одну цифру
  - трюк позволяет в большинстве случаев сэкономить одну цифру по сравнению с двоичным кодом числа  $N$

## Кодирование нулей (ZLE)

- ★ В тексте  $\text{mtf}(\text{bwt}(T))$  есть длинные повторы нулей, соответствующие контекстным зависимостям в  $T$ , в то время как повторы других символов при MTF могут возникнуть только случайно (а значит, они короткие)

На первом этапе сжатия кодируются повторы нулей

- ★ Лучше всего работает метод **Zero Length Encoding (ZLE)**; он хорошо подготавливает почву для арифметического кодирования на втором этапе

ZLE:

- сделаем 254 и 255 — два самых редких символа в  $\text{mtf}(\text{bwt}(T))$  — служебными
- каждую последовательность нулей заменим на ее длину, в качестве цифр 0 и 1 используя 254 и 255
- 254 и 255, если они встречались, заменяем на пару байт 0 254 (0 255)
- ★ **дополнительный трюк**: для последовательности из  $N$  нулей кодируем число  $N+1$  без первой цифры
  - $N+1 \geq 2$ : двоичная запись начинается с 1 и содержит еще хотя бы одну цифру
  - трюк позволяет в большинстве случаев сэкономить одну цифру по сравнению с двоичным кодом числа  $N$

Пример:  $\underbrace{0 \dots 0}_{100 \text{ раз}}$  ZLE закодирует как 255 254 254 255 254 255, так как  $101_2 = \underline{1}100101$

Альтернативный вариант сжатия после BWT — **Distance Coding (DC)**

- Текст сканируется слева направо, для каждого символа записывается ссылка на следующий такой же символ

## Альтернативный вариант сжатия после BWT — Distance Coding (DC)

- Текст сканируется слева направо, для каждого символа записывается ссылка на следующий такой же символ

Более строго:

- Припишем слева к тексту  $\text{bwt}(T)$  алфавит  $A$  (обычно  $|A| = 256$ ); пусть  $X = A \cdot \text{bwt}(T)$ 
  - декодер в этот момент знает все символы  $A$  и ни одного символа  $\text{bwt}(T)$
- Для каждого символа  $a = X[i]$  слева направо записываем номер: каким по счету слева из неизвестных декодеру символов является следующий символ  $a$ 
  - если вхождений  $a$  больше нет, пишем номер 0
  - символ  $a$ , на который мы сослались, стал известен декодеру
- ★ Если кодер передает номер для символа  $X[i]$ , а на символ  $X[i+1]$  еще не было ссылки, то эта ссылка должна быть сделана сейчас, то есть  $X[i+1] = X[i]$ 
  - декодер в состоянии сам определить эту ситуацию (ему неизвестен символ, следующий за текущим) и скопировать текущий символ в следующий, **сделав следующий известным**
  - ⇒ кодер ничего не кодирует в данной ситуации
  - ★ это правило — **основной источник сжатия**
- Кодирование завершается, когда декодер знает все символы  $X$
- Номер кодируется с переполнением (см. Лекцию 6), например, по схеме  $8+16+24$  бит или  $8+16+16+\dots$
- Пример ⇒

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bana\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 1

кодер	<b>#</b>	Б	а	н	а	<b>#</b>	н	н	Б	н	н	Б	а	а	а	а
	?															
декодер	<b>#</b>	Б	а	н	.	.	.	.	.	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 1

кодер	<b>#</b>	Б	а	н	а	<b>#</b>	н	н	Б	н	н	Б	а	а	а	а
декодер	<b>#</b>	Б	а	н	.	.	.	.	.	.	.	.	.	.	.	.



Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 2

кодер	#	<b>B</b>	a	n	a	#	n	n	<b>B</b>	n	n	B	a	a	a	a
	2	?														
декодер	#	B	a	n	.	#	.	.	.	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 2

кодер	#	<b>B</b>	a	n	a	#	n	n	<b>B</b>	n	n	B	a	a	a	a
		2							4							
декодер	#	B	a	n	.	#	.	.	.	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 3

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	?													
декодер	#	Б	а	н	.	#	.	.	Б	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 3

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1													
декодер	#	Б	а	н	.	#	.	.	Б	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 4

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	?												
декодер	#	Б	а	н	а	#	.	.	Б	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Ба\#nnBnnБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 4

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1												
декодер	#	Б	а	н	а	#	.	.	Б	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 5

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	?											
декодер	#	Б	а	н	а	#	н	.	Б	.	.	.	.	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnБаааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Бапа\#nnBnnБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 5

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5											
декодер	#	Б	а	н	а	#	н	.	Б	.	.	.	.	.	.	.



Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 6

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0										
декодер	#	Б	а	н	а	#	н	.	Б	.	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 7

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0										
декодер	#	Б	а	н	а	#	н	.	Б	.	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 7

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0										
декодер	#	Б	а	н	а	#	н	н	Б	.	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 8

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		?								
декодер	#	Б	а	н	а	#	н	н	Б	.	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 8

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0				1						
декодер	#	Б	а	н	а	#	н	н	Б	.	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 9

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0	1	?								
декодер	#	Б	а	н	а	#	н	н	Б	н	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 9

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2							
декодер	#	Б	а	н	а	#	н	н	Б	н	.	.	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 10

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2							
декодер	#	Б	а	н	а	#	н	н	Б	н	.	Б	а	.	.	.



Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 10

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2							
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Ба\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 11

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2		0					
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#nnBnnBaааа}$
- Алфавит  $\{\#, \text{B}, \text{a}, \text{n}\}$
- $X = \text{\#Bана\#nnBnnBaааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 12

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2		0	0				
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	.	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 13

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2		0	0				
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	.	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Бана\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 14

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2		0	0				
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	.

Рассмотрим наш сквозной пример:

- Текст  $\text{bwt}(T) = \text{a\#ннБннБаааа}$
- Алфавит  $\{\#, \text{Б}, \text{а}, \text{н}\}$
- $X = \text{\#Ба\#ннБннБаааа}$
- Символы, неизвестные декодеру, обозначены точками

Итерация 15

кодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а
	2	4	1	1	5	0		1	2		0	0				
декодер	#	Б	а	н	а	#	н	н	Б	н	н	Б	а	а	а	а

$\text{dc}(\text{bwt}(T)) = 2411501200$

★ Результат DC можно еще немного сжать Хаффманом или арифметикой

- ★ Изначально BWT-алгоритмы позиционировались как “золотая середина” между быстрыми, но слабыми словарными алгоритмами и мощным, но медленным PPM
- ★ На самом деле, они лучше, чем среднее арифметическое LZ и PPM
  - Самый известный архиватор общего назначения, основанный на BWT — bzip2 — сжимает многократно упоминавшийся датасет английской википедии в 3.9 раза, а, например, архиватор M03 — в 6.1 раза, обгоняя 7-Zip и WinRar на их максимальных настройках и по коэффициенту сжатия, и по скорости

- ★ Изначально BWT-алгоритмы позиционировались как “золотая середина” между быстрыми, но слабыми словарными алгоритмами и мощным, но медленным PPM
- ★ На самом деле, они лучше, чем среднее арифметическое LZ и PPM
  - Самый известный архиватор общего назначения, основанный на BWT — bzip2 — сжимает многократно упоминавшийся датасет английской википедии в 3.9 раза, а, например, архиватор M03 — в 6.1 раза, обгоняя 7-Zip и WinRar на их максимальных настройках и по коэффициенту сжатия, и по скорости
  - BWT имеет и другие приложения; самое известное — индексирование данных (FM-индекс и его производные)
  - ... ну и это просто математически красиво!