

Преобразование Барроуза-Уиллера

5 января 2023 г. 1:27

Пусть имеется последовательность на выходе источника.

Берем все возможные циклические сдвиги (т.е. берем последнюю букву и ставим в начало и т.д.), сортируем их в лексикографическом порядке и формируем из них таблицу.

Результатом преобразования являются:

1. Последний столбец таблицы
2. Номер строки, которая соответствует исходной последовательности

0	A	NIF WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE	C
1	A	NNOT DO AS WE WOULD WE SHOULD DO AS WE CANIF WE	C
2	A	S WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO	-
3	A	S WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO	-
4	C	ANIF WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE	-
5	C	ANNOT DO AS WE WOULD WE SHOULD DO AS WE CANIF WE	-
6	D	O AS WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD	-
7	D	O AS WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT	-
8	D	DO AS WE CANIF WE CANNOT DO AS WE WOULD WE SHOU	L
9	D	WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE WOU	L
10	E	CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO AS	W
11	E	CANNOT DO AS WE WOULD WE SHOULD DO AS WE CANIF	W
12	E	SHOULD DO AS WE CANIF WE CANNOT DO AS WE WOULD	W
13	E	WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO AS	W
14	F	WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE CAN	I*
15	H	OULD DO AS WE CANIF WE CANNOT DO AS WE WOULD WE	S
16	I	F WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE CA	N
17	L	D DO AS WE CANIF WE CANNOT DO AS WE WOULD WE SHO	U
18	L	D WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE WO	U
19	N	IF WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE C	A
20	N	OT DO AS WE WOULD WE SHOULD DO AS WE CANIF WE C	A
21	N	OT DO AS WE WOULD WE SHOULD DO AS WE CANIF WE CA	N
22	O	T DO AS WE WOULD WE SHOULD DO AS WE CANIF WE CAN	N
23	O	ULD DO AS WE CANIF WE CANNOT DO AS WE WOULD WE S	H
24	O	ULD WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE	W
25	O	AS WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD	D
26	O	AS WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT	D
27	S	HOULD DO AS WE CANIF WE CANNOT DO AS WE WOULD WE	-
28	S	WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO	A
29	S	WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO	A
30	T	DO AS WE WOULD WE SHOULD DO AS WE CANIF WE CANN	O
31	U	LD DO AS WE CANIF WE CANNOT DO AS WE WOULD WE SH	O
32	U	LD WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE W	O
33	W	E CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO AS	-
34	W	E CANNOT DO AS WE WOULD WE SHOULD DO AS WE CANIF	-
35	W	E SHOULD DO AS WE CANIF WE CANNOT DO AS WE WOULD	-
36	W	E WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO AS	-
37	W	OULD WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE	-
38	-	AS WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD D	O
38	-	AS WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT D	O
40	-	CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO AS W	E
41	-	CANNOT DO AS WE WOULD WE SHOULD DO AS WE CANIF W	E
42	-	DO AS WE CANIF WE CANNOT DO AS WE WOULD WE SHOUL	D
43	-	DO AS WE WOULD WE SHOULD DO AS WE CANIF WE CANN	O
44	-	SHOULD DO AS WE CANIF WE CANNOT DO AS WE WOULD W	E
45	-	WE CANIF WE CANNOT DO AS WE WOULD WE SHOULD DO A	S
46	-	WE CANNOT DO AS WE WOULD WE SHOULD DO AS WE CANI	F
47	-	WE SHOULD DO AS WE CANIF WE CANNOT DO AS WE WOULD	D
48	-	WE WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO A	S
49	-	WOULD WE SHOULD DO AS WE CANIF WE CANNOT DO AS W	E

После формирования таблицы мы получили номер строки (16) и последний столбец:

CC____LLWWWWISNUUAANNHWDD_AA000____OOEEDTESFDSE

Видно, что формировалась последовательность с явными повторениями символов - это можно учесть при сжатии.

В итоге передается этот последний столбец и номер строки, содержащий исходную последовательность

Как декодировать этот последний столбец?

Если отсортировать его в алфавитном порядке, то получим первый столбец.

Таким образом, известен номер строки (16), первый и последний столбец таблицы.

Надо найти второй столбец.

Смотрим строку 16, в первом столбце видим I.

Ищем I в последнем столбце.

Нашли строку, снова смотрим первый её столбец.

И т.д.

Но в последнем столбце присутствуют одинаковые символы - какой выбрать?

Считаем, какой по счету этот символ в последнем столбце, в первом столбце ищем такой же по счету этот символ.

Кодирование последнего столбца "стопкой книг"

- CC____LLWWWWISNUUAANNHWDD_AA000____OOEEDTESFDSE
- Будем писать `esc` каждый раз, когда появляется новый символ и число различных букв между предыдущим и текущим появлением буквы.
- Получим: `esc 0 esc 0 0 0 0 0 esc 0 esc 0 0 0 esc esc esc esc 0 esc 0 2 0 esc 6 esc 0 9 5 0 esc 0 0 2 0 0 0 0 1 0 esc 0 4 esc 2 10 esc 4 2 3`
- Применяем нумерационное кодирование:
 - ▶ Кодирование композиции: 35 бит.
 - ▶ Кодирование последовательности: 94 бита.

- Передача символов, соответствующих `esc`: $8 \times 15 = 120$ бит.
- Передача индекса исходной строки в преобразовании: 6 бит.

Итого: 255 бит.