

# Сжатие информации с учётом контекста

Александра Игоревна Кононова

МИЭТ

21 декабря 2020 г. — актуальную версию можно найти на  
<https://gitlab.com/illinc/otik>

# Источник с памятью

Вероятность порождения элемента определяется предысторией источника  $\rightarrow$  при сжатии необходимо учитывать контекст символа.

**источник Маркова ( $N$ -го порядка)** — состояние на  $i$ -м шаге зависит от состояний на  $N$  предыдущих шагах:

$i - 1, i - 2, \dots, i - N \rightarrow$  сжимаются не отдельные символы, а устойчивые сочетания — слова:  
коды Зива—Лемпеля (LZ77, LZ78);

**аналоговый сигнал** — источник количественных данных Маркова 1-го порядка  $\rightarrow$  кодирование длин повторений (Run Length Encoding, RLE).

# Концепция RLE

Вместо кодирования данных кодируются длины участков, на которых данные сохраняют неизменное значение.

AAAAAAAAABCCCC  $\rightarrow 8 \times A, 1 \times B, 4 \times C$

Повторение символа  $c$  подряд  $L$  раз ( $L \times c$ ),  $L > 0$  будем записывать как пару  $(L, c)$ .

Если  $c$  и  $L$  одного размера — это имеет смысл только при  $L > 2$ .

# Реализации RLE

- 1 «Сжимать» любую цепочку:  $L \times c \rightarrow (L, c)$  для любого  $L$  ( $L \geq 1$ ):

0123456789ABCDEF  $\rightarrow$  101112131415161718191A1B1C1D1E1F

увеличение объёма в наихудшем случае — в 2 раза;

- 2 Сжатые и несжатые цепочки: 
$$\begin{cases} L \times c, L \geq 3 \rightarrow (1 + L, c), \\ c_1 \dots c_L, L \geq 1 \rightarrow (0 + L, c_1 \dots c_L) \end{cases}$$

байт-октет — 7 бит на  $L$  ( $0 \dots 127$  без коррекции), 4-битный — 3 ( $0 \dots 7$ ):

0123456789AAAAAAAAABCDEF  $\rightarrow$  0123456, 789, AAAAAAA, ABCDEF  $\rightarrow$   
 $\rightarrow$  701234563789FA6ABCDEF (21) (для октетов ув. в нх. сл. на  $\frac{1}{127}$ )

Коррекция  $L$ : запись  $L - 1$  для несжатой ( $1 \dots 8$ ),  $L - 3$  для сжатой ( $3 \dots A$ ):

0123456789AAAAAAAAABCDEF  $\rightarrow$  01234567, 89, AAAAAAA, BCDEF  $\rightarrow$   
 $\rightarrow$  701234567189DA4BCDEF (20) (для октетов ув. в нх. сл. на  $\frac{1}{128}$ )

- 3 Префикс сжатой цепочки в несжатом тексте:  $L \times c \rightarrow (p, L, c)$ ,  $L \geq 4$ :

$p = C$ : 0123456789AAAAAAAAABCDEF (23)  $\rightarrow$  0123456789C8ABC1CDEF (20)

или 0123456789C8ABCODEF (19) (для октетов ув. в нх. сл. на  $\frac{2}{256}$  или  $\frac{1}{256}$ )

Коррекция: для C0 ( $p, L - 3, c$ ); для C1C 
$$\begin{cases} L \times c, c \neq p \rightarrow (p, c, L - 4) \\ L \times p \rightarrow (p, p, L - 1), \end{cases}$$

Run Length Encoding, RLE (аналоговый сигнал)

Код Зива–Лемпеля, LZ77

Семинар RLE/LZ77

Код Зива–Лемпеля, LZ78 (концепция)

Код Зива–Лемпеля–Велча, LZW

Концепция RLE

Реализации RLE

# Код Зива–Лемпеля, LZ77

1977 г., Якоб Зив (Jacob Ziv) и Абрахам Лемпель (Abraham Lempel):  
последующие вхождения цепочки заменяются ссылками на первое:  $(S, L)$

**Скользящее окно:** область перед текущей позицией кодирования,  
в которой можем искать и адресовать ( $w \leq \max(S)$ ) ссылки

---

Если  $S$  записывается одним байтом ( $n$  различных значений),  
то максимальный размер окна  $w = \max(S) = n - 1$ .

Наихудший случай (для 4-битного байта):

0123456789ABCDEF0123456789ABCDEF...

в пределах окна нет ни одного совпадения.

Увеличить размер окна  $w$  можно, записывая  $S$  двумя байтами.

# Пример №1

Alternative vexillum codicis inf. interpretatio (AVCII): ,абвгдеиклмнорты  
0123456789ABCDEF

01234567890123456789012345678901234567890123

там корабли лавировали, лавировали, да не вылавировали

$\begin{cases} S_1 = 23 - 11 = 12 = C, \\ L_1 = 13 = D \end{cases}$  — кодирование 4-битными байтами возможно

012345678901234567890123456789012345678901234567890123

там корабли лавировали, лавировали, да не вылавировали

$S_2 = 20 > F$  — кодирование 4-битным байтом невозможно

При программной реализации (поиск только в окне) при  $w = 15$  будет найдено только одно совпадение  $(S_1, L_1)$ .

там корабли лавировали,  $(S_1, L_1)$  да не вылавировали

$(S_1, L_1) = CD = \text{ор}$

# Пример №2, Михаил Быстрынцев

The compression and the decompression leave an impression. Hahahahaha!

The compression and t[he ]de[compression ]leave[ an] i[mpression]. Hah[ahahaha]!

The compression and the decompression leave an impression. Hahahahaha!

Абсолютные ссылки  $(P, L)$ , где  $P$  — позиция первого вхождения цепочки,  $L$  — её длина:

The compression and t(2,3)de(5,12)leave(16,3) i(8,7). Hah(61,7)!

Замена абсолютных позиций  $P$  на относительные смещения  $S$ :

The compression and t(20,3)de(22,12)leave(28,3) i(42,7). Hah(2,7)!

## Как разделить несжатый текст от ссылок?

# Реализации LZ77

- 1 Несжатый символ  $c$  как «ссылка»  $(0, c)$  — увеличение объёма в  $n$ х. сл. вдвое
- 2 Цепочка, за которой следует  $c$  — тройка  $(S, L, c)$ , несжатый символ  $c$  — тройка  $(0, 0, c)$  — увеличение объёма в наихудшем случае втрое

- 3 Быстрынцев: символы и цепочки группируются по 8, тип — байт-флаги:  
 The<sub>⌊</sub>c o mp 0b 0000 0000 The comp — увеличение объёма  
 ressi o n<sub>⌊</sub> 0b 0000 0000 ressi on в наихудшем случае на  $\frac{1}{8}$   
 and<sub>⌊</sub>t(20,3)de 0b 0000 0100 and t 

20	03
----	----

 de

- 4 Префикс  $p$  как маркер ссылки:  $\begin{cases} \text{ссылка } (p, S, L) : 0 < S \leq w, L \geq 4, \\ \text{символ } p \rightarrow (p, 0) \end{cases}$   
 — увеличение объёма в наихудшем случае на  $\frac{1}{256}$

$p = F = \text{ы}; \quad S_1 = 12 = C, \quad L_1 = 13 = D$

там корабли лавировали, **ы**CD да не **вы**0 лавировали ( $L$  без коррекции)

там корабли лавировали, **ы**C9 да не **вы**0 лавировали ( $L \rightarrow L - 4$ )



# Сравнение с RLE

Префикс  $p = C$ :

Повторение одинаковых символов:

RLE:

0123456789AAAAAAAABCDEF  $\rightarrow$  0123456789C8ABCDEF 23  $\rightarrow$  19

LZ77:

0123456789AAAAAAAABCDEF  $\rightarrow$  0123456789AC17BCDEF 23  $\rightarrow$  20

Повторение подстрок:

RLE:

28  $\rightarrow$  24

012345601234789AAAAAAAABCDEF  $\rightarrow$  012345601234789C8ABCDEF

LZ77:

28  $\rightarrow$  23

012345601234789AAAAAAAABCDEF  $\rightarrow$  0123456C75789AC17BCDEF

# Задача №1

- 1 Рассчитать количество информации в сообщении  $C = \underbrace{\text{«ля-ля-ля-...-ля»}}_{80 \text{ раз}}$  (кодировка koï8-r)
- 2 Закодировать сообщение  $C$  алгоритмом LZ77

# Информация и модели источника

Заменяя вероятности символов на их оценки по модели  $X$ , получаем **оценку** количества информации:

- 1 Без памяти,  $A_1 = \{\text{л}, \text{я}, -\}$ :

$$I_1 = 2 \cdot 80 \cdot \left( -\log_2 \left( \frac{80}{239} \right) \right) + 79 \cdot \left( -\log_2 \left( \frac{79}{239} \right) \right) = 254,2 \text{ бит} = 31,8 \text{ байт}$$

- 2 Без памяти,  $A_1 = \{\text{ля-}, \text{ля}\}$ :

$$I_2 = 79 \cdot \left( -\log_2 \left( \frac{79}{80} \right) \right) + 1 \cdot \left( -\log_2 \left( \frac{1}{80} \right) \right) = 7,6 \text{ бит} = 0,9 \text{ байт}$$

- 3 Маркова,  $A_1 = \text{ko}i8\text{-r}$ ,  $N = 3$  со следующими оценками вероятностей:

$$\begin{aligned} p(c_1 = \text{л}) &= p(c_2 = \text{я}) = p(c_3 = -) = \frac{1}{256}, & p(-|\text{ля}) &= \frac{79}{80}, \\ p(\text{л}|\text{ля-}) &= p(\text{я}|\text{я-л}) = 1, & p(\text{eof}|\text{ля}) &= \frac{1}{80}: \end{aligned}$$

$$I_3 = 3 \cdot \log_2 256 + 79 \cdot \left( -\log_2 \left( \frac{79}{80} \right) \right) + 1 \cdot \left( -\log_2 \left( \frac{1}{80} \right) \right) = 31,6 \text{ бит} = 3,9 \text{ байт}$$

# Задача №2

Используя 4-битный байт, закодируйте:

- 1 Сообщение (64 = 0x40 символов) — LZ77:

	0123456789ABCDEF	0123456789ABCDEF	0123456789ABCDEF	0123456789ABCDEF
00	косил косой косарь	за косарь косой	косой косой космеи	на косе при осе.

код $i$	0123456789ABCDEF
AVCII: $a_i$	.aeзийклмнопрсь
$\nu(a_i)$	B1431338111C1392

а)  $S$  занимает 1 байт (4 бита); б)  $S$  занимает 5 бит; в)  $S$  занимает 2 байта.

- 2 16-цветный рисунок  $8 \times 8$  — RLE, LZ77.

# Код Зива–Лемпеля, LZ78 (концепция)

1978 г., Якоб Зив (Jacob Ziv) и Абрахам Лемпель (Abraham Lempel):

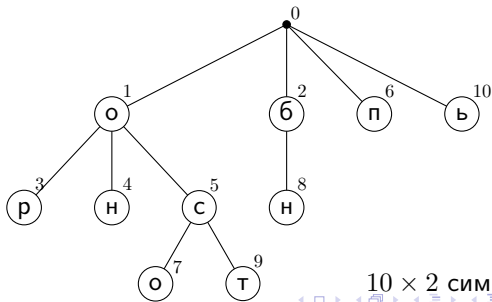
- 1 Скользящее окно не используем — кодируем в один проход вперёд  $\Rightarrow$  высокая скорость кодирования-декодирования.
- 2 Словарь = дерево, узел — номер и символ  $(n, c)$ , корень —  $(0, \text{пустая строка})$ , слово читается от корня.
- 3 Вначале словарь пуст (только корень).
- 4 На каждом шаге
  - к словарю добавляется узел (лист);
  - в выходной поток — номер родителя и символ нового листа  $(P, c)$ .
- 5 Когда кончается ёмкость номера листа, дерево:
  - либо уничтожается и растится заново;
  - либо ветви уничтожаются выборочно (сложно);
  - либо фиксируется и не растёт (нет прикорневого узла  $\Rightarrow$  сбои);
  - либо увеличивается разрядность номера.
- 6 При необходимости вх-й поток дополняется (либо конец обр-ся особо).



# «Обороноспособность» (18, 8 разных)

- 1 Вначале словарь = корень (пустая строка),  $n = 0$ ,  $i = 0$ .
- 2  $P = 0$  (текущий узел — корень),  $c_i$  (текущий символ входного потока);
- 3 Если  $c_i$  — дочерний  $P$ ,  $P = c_i$  и читаем  $c_{i+1}$  ( $++i$ )
- 4 Если  $c_i$  нет в дочерних узлах  $P$ :
  - добавляем  $P$  дочерний узел  $(n, c_i)$ ,  $++n$  и читаем  $c_{i+1}$  ( $++i$ );
  - в выходной поток пишем  $(P, c_i)$ .

1	(0,о)	о
2	(0,б)	б
3	(1,р)	ор
4	(1,н)	он
5	(1,с)	ос
6	(0,п)	п
7	(5,о)	осо
8	(2,н)	бн
9	(5,т)	ост
10	(0,ь)	ь



10 × 2 СИМВОЛОВ

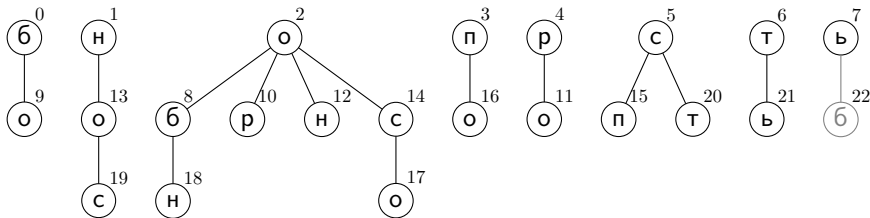


# Код Зива–Лемпеля–Велча, LZW

1984 г., Терри Велч (Terry Welch) по концепции LZ78:

- 1 Вначале словарь = первый уровень (все одиночные символы,  $N$  штук). Тогда корень можно не нумеровать (прикорневые нумеруем с нуля).
- 2 При добавлении  $P$  дочернего узла  $(n, c_i)$ :
  - оставляем  $c_i$  во входном потоке;
  - в выходной поток пишем  $(P)$ .
- 3 При декодировании узла  $n$ :
  - в выходной поток пишем всю ветвь;
  - в дерево добавляем только прикорневой узел.
- 4  $|n| \gg |c|$ , во многих реализациях увеличивается по битам.
- 5 Дерево часто разворачивается в таблицу.
- 6 Вх-й поток всегда дополняется как минимум одним незначащим символом.

# «Обороноспособность» (18, алфавит из 8)



8	(2, б)	об	2	16	(3, о)	по	3
9	(0, о)	бо	0	17	(14, о)	осо	14
10	(2, р)	ор	2	18	(8, н)	обн	8
11	(4, о)	ро	4	19	(13, с)	нос	13
12	(2, н)	он	2	20	(5, т)	ст	5
13	(1, о)	но	1	21	(6, ь)	ть	6
14	(2, с)	ос	2	22	(7, б)	ьб	7
15	(5, п)	сп	5	15 СИМВОЛОВ			



# Декодирование (замечания)

- 1 При декодировании на  $i$ -м шаге (входной номер  $P_i$ ) в выходной поток добавляется строка  $C_i$ , соответствующая узлу  $P_i$  (то есть на один символ короче, чем была на  $i$ -м шаге при кодировании),  
а в дереве словаря от самого  $P_i$  должен отрасти дочерний узел с номером  $i$  и неизвестным символом  $c_i$ .
- 2 Символ  $c_i$  узла  $i$  становится известным только на шаге  $i + 1$  ( $c_i$  — это первый символ подстроки  $C_{i+1}$  шага  $i + 1$ ),  
поэтому на практике узел  $i$  добавляется в словарь на шаге  $i + 1$ .
- 3 Если на  $i + 1$  шаге получаем ссылку на ещё не добавленный узел  $P_{i+1} = i$ , то всё равно  $c_i$  — это первый символ подстроки  $C_{i+1}$ ,  
а первый символ  $C_{i+1}$  мы знаем даже при  $P_{i+1} = i$   
(как и все до предпоследнего включительно)!  
При  $P_{i+1} = i$  последний символ строки  $C_{i+1}$  совпадает с первым:  
 $C_{i+1} = axx \dots xa$ .

# Декодирование

$\textcircled{6}^0$      $\textcircled{н}^1$      $\textcircled{о}^2$      $\textcircled{п}^3$      $\textcircled{р}^4$      $\textcircled{с}^5$      $\textcircled{т}^6$      $\textcircled{ь}^7$

2 о  
 0 6 8 (2, 6)  
 2 о 9 (0, о)  
 4 р 10 (2, р)  
 2 о 11 (4, о)  
 1 н 12 (2, н)  
 2 о 13 (0, о)  
 5 с 14 (2, с)

3 п 15 (5, с)  
 14 ос 16 (3, о)  
 8 об 17 (14, о)  
 13 но 18 (8, н)  
 5 с 19 (13, с)  
 6 т 20 (5, т)  
 7 ь 21 (6, ь)

# Спасибо за внимание!

МИЭТ

<http://miet.ru/>

Александра Игоревна Кононова

[illinc@mail.ru](mailto:illinc@mail.ru)