Сжатие с учётом контекста. Словарные методы, где словарём является несжатый текст — семейство LZ77

Александра Игоревна Кононова

ТЄИМ

26 января 2021 г. — актуальную версию можно найти на https://gitlab.com/illinc/otik

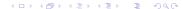


Вероятность порождения элемента определяется предысторией

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, ..., i-N o сжимаются не отдельные символы, а устойчивые сочетания — слова: коды Зива—Лемпеля (LZ77, LZ78);
- аналоговый сигнал источник количественных данных Маркова 1-го порядка o кодирование длин повторений (Run Length Encoding, RLE).

При расчёте количества информации источника Маркова N-го порядка вероятности зависят от предыдущих N символов (непостоянны).



Для компактной иллюстрации ограничений алгоритмов примем, что для устройства «доска» байт (символ кодирования) составляет не 8 бит — октет (как для Intel x86/amd64), и не 6 бит (как для IBM 7030 Stretch), а 4 бита — тетраду, или одну 16-ричную цифру:

$$\begin{array}{lll} 0 = 0000 = 0 & 8 = 1000 = 8 \\ 1 = 0001 = 1 & 9 = 1001 = 9 \\ 2 = 0010 = 2 & A = 1010 = 10 \\ 3 = 0011 = 3 & B = 1011 = 11 \\ 4 = 0100 = 4 & C = 1100 = 12 \\ 5 = 0101 = 5 & D = 1101 = 13 \\ 6 = 0110 = 6 & E = 1110 = 14 \\ 7 = 0111 = 7 & F = 1111 = 15 \end{array}$$

Символьная таблица (аналог ASCII) даётся отдельно в каждом примере, либо отсутствует.



Концепция RLE

Модель источника данных — Маркова первого порядка (аналоговый сигнал): вероятность символа зависит от предыдущего символа.

Run Length Encoding (RLE): вместо кодирования данных кодируются длины участков, на которых данные сохраняют неизменное значение.

AAAAAAAABCCCCC \rightarrow 8 × A, 1 × B, 4 × C

Повторение символа c подряд L раз $(L \times c)$ — цепочку длины L, $L \geqslant 1$ будем записывать как пару (L,c).

Если c и L одного размера — это имеет смысл (приводит к сжатию цепочки) только при L>2 ($L\geqslant 3$).



Основные варианты кодирования RLE

- Наивный любая цепочка $L \times c$ для любого L ($L \geqslant 1$) заменяется парой (L,c): $\forall L: L \times c \rightarrow (L,c)$
- Сжатые и несжатые цепочки (флаг-бит) парой (L,c) заменяются только длинные цепочки $L \geqslant 3$, короткие цепочки из одного и двух символов объединяются в более длинные несжатые цепочки, каждая из которых также предваряется длиной L; для различения сжатых и несжатых цепочек выделяется один бит поля L: $\begin{cases} L \times c, L \geqslant 3 & \to & (1+L,c), \\ c_1 \dots c_L, L \geqslant 1 & \to & (0+L,c_1 \dots c_L). \end{cases}$
- Префикс сжатой цепочки в несжатом тексте текст записывается как есть, длинные цепочки заменяются **тройками** (p, L, c), где $p \in A$;

сжатие — только при
$$L\geqslant 4$$
:
$$\begin{cases} L\times c, L\geqslant 4 & \to & (p,L,c),\\ c_1\dots c_L, L\geqslant 1 & \to & c_1\dots c_L.\\ L\times p, L\geqslant 1 & \to & (p,L,p). \end{cases}$$

Префикс p выбирается как самый редкий символ текста.



Сравнение (без коррекции длины) І

Исходные сообщения:

- lacktriangle АААААААААААААААВВ (18) хороший случай (16 imes A и 2 imes В)
- ААВСDDEF0012 (12) нормальный случай
- 0123456789AAAAAAAABCDEF(23) нормальный случай ($8 \times A$)
- 0123456789ABCDEF (16) наихудший случай
- Наивный вариант:

```
{\tt AAAAAAAAAAAAAAABB} 
ightarrow egin{pmatrix} {\sf FA1A2B} & (6) \\ {\sf AFA1B2} & (6) \end{pmatrix} для наивного RLE допустимы оба
```

порядка полей — (L,c) или (c,L); далее используется (L,c) для единообразия $AABCDDEFO012 \rightarrow 2A1B1C2D1E1F201112 (18)$

0123456789AAAAAAABCDEF $\rightarrow 10111213141516171819$ 8A1B1C1D1E1F (32) $0123456789 \text{ABCDEF} \rightarrow 101112131415161718191 \text{A} 1 \text{B} 1 \text{C} 1 \text{D} 1 \text{E} 1 \text{F} (32)$

Увеличение объёма в наихудшем случае — в 2 раза.

Длинные цепочки встречаются редко, поэтому поле L обычно равно по длине символу c (однобайтовое)

Сравнение (без коррекции длины) ІІ

Флаг-бит: $1\leqslant L\leqslant 7$ (3 бита), $1+7\sim$ 1111 = F, $1+2\sim$ 1010 = A

 $\texttt{AAAAAAAAAAAAAAAAAABB} \to \begin{bmatrix} \texttt{AAAAAAA}, \texttt{AAAAAAA}, \texttt{AABB} \to \texttt{FAFA4AABB} \ (9) \\ \texttt{AAAAAAA}, \texttt{AAAAAAA}, \texttt{AB}, \texttt{BB} \to \texttt{FAFAAAAB} \ (8) \\ \end{bmatrix}$

но AABCDDEF0012 \rightarrow $\begin{bmatrix} \text{AABCDDE, F0012} \rightarrow \text{7AABCDDE5F0012 (14)} \\ \text{AA, BC, DD, EF, 00, 12} \rightarrow \text{AA2BCAD2EFA0212 (15)} \end{bmatrix}$

0123456789AAAAAAABCDEF \rightarrow 0123456, 789, AAAAAAA, ABCDEF \rightarrow \rightarrow 701234563789FA6ABCDEF (21)

0123456789ABCDEF \rightarrow 0123456, 789ABCD, EF \rightarrow 701234567789ABCD2EF (19)

здесь порядок полей (L,c) важен: иначе не декодируются несжатые цепочки

Увеличение объёма в наихудшем случае — для байта-тетрады на $\frac{1}{7}$ объёма файла (к каждым семи символам добавляется восьмой — флаг-бит + L), для байта-октета — на $\frac{1}{127}$ (менее процента).

《□》《圖》《意》《意》。 毫二

Префикс в несжатом тексте, $p= exttt{E}$:

Адададададададададав o EFAABB (6)

Два варианта кодирования одиночного символа $p \colon 1 \times p o igg[(p,1,p) \ (p,0) igg]$

 $\texttt{0123456789AAAAAAAAABCDEF} \rightarrow \begin{bmatrix} \texttt{0123456789E8ABCDE1EF} & (20) \\ \texttt{0123456789E8ABCDE0F} & (19) & L \neq 0 \\ \end{bmatrix}$

но только один — для двойного p: ... DEEF ightarrow ... DE2EF (а не ... DE0E0F).

Если Е1Е, то допустимо и $(p,L,c)\sim$ Е1Е/Е8А, и $(p,c,L)\sim$ ЕЕ1/ЕА8.

Но если мы хотим Е0, то только (p,L,c).

$$\texttt{0123456789ABCDEF} \rightarrow \begin{bmatrix} \texttt{0123456789ABCDE1EF} \ (18) \\ \texttt{0123456789ABCDE0F} \ (17) \\ \end{bmatrix}$$

Увеличение объёма в наихудшем случае (все символы одиночные) при (p,0) — на $\nu(p)$ от объёма файла, где $\nu(p)$ — частота p: для байта-тетрады $\nu(p)\leqslant \frac{1}{16}$, для байта-октета $\nu(p)\leqslant \frac{1}{256}<0.4\%$ (p — самый редкий символ).

40 140 15 15 15 15 10 0

Коррекция длины

Максимальная длина L определяется разрядностью поля (обычно байта), минимальная — алгоритмом:

Байт	Наивный	Флаг-бит	Префикс
\overline{N} бит	$1 \leqslant L < 2^N$	$3 \leqslant L < 2^{N-1}$ (сж)	$4 \leqslant L < 2^N \ (c \neq p)$
		$1\leqslant L<2^{N-1}$ (н/сж)	$2 \leqslant L < 2^N \ (c = p)$
4 бита	$1 \leqslant L \leqslant 15$	$3 \leqslant L \leqslant 7$ (сж)	$4 \leqslant L \leqslant 15 \ (c \neq p)$
		$1\leqslant L\leqslant 7$ (н/сж)	$2 \leqslant L \leqslant 15 \ (c=p)$
8 бит	$1 \leqslant L \leqslant 255$	$3 \leqslant L \leqslant 127$ (сж)	$4 \leqslant L \leqslant 255 \ (c \neq p)$
		$1\leqslant L\leqslant 127$ (н/сж)	$2 \leqslant L \leqslant 255 \ (c=p)$
Коррекция: вместо L записывается $\widetilde{L} = L - \min L$			
4 бита	$1 \leqslant L \leqslant 16$	$3 \leqslant L \leqslant 10$ (сж)	$4 \leqslant L \leqslant 16 \ (c \neq p)$
		$1\leqslant L\leqslant 8$ (н/сж)	$2 \leqslant L \leqslant 16 \ (c=p)$
	$0 \leqslant \widetilde{L} \leqslant 15$	$0 \leqslant \widetilde{L} \leqslant 7$	$1\leqslant\widetilde{L}\leqslant15$

Цепочки с коррекцией L длиннее \Rightarrow их количество в файле меньше \Rightarrow файл после кодирования короче.

◆□▶ ◆□▶ ◆■▶ ◆■▶ ■ 990

Сравнение (с коррекцией длины)

- lacksquare Наивный: $\forall L:\ L imes c o (L-1,c)$ AAAAAAAAAAAAAABB $\rightarrow 16 \times A, 2 \times B \rightarrow FA1B$ (4)
- Флаг-бит: $egin{cases} L imes c, L\geqslant 3 & o & \Big(1+(L-3),c\Big), \\ c_1\dots c_L, L\geqslant 1 & o & \Big(0+(L-1),c_1\dots c_L\Big). \end{cases}$ AAAAAAAAAAAAAAABB $\rightarrow 10 \times \text{A}, 6 \times \text{A}, \text{BB} \rightarrow \text{FABA1BB}$ (7) 0123456789AAAAAAABCDEF \rightarrow 01234567, 89, AAAAAAA, BCDEF \rightarrow \rightarrow 701234567189DA4BCDEF (20)
- 0123456789ABCDEF o 01234567, 89ABCDEF o 701234567789ABCDEF (18)
- Префикс: $\begin{cases} L\times c, c\neq p, L\geqslant 4 & \to & (p,\ L-1,\ c) \\ L\times p, L\geqslant 2 & \to & (p,\ L-1,\ p),\ L-1\neq 0 \\ 1\times p & \to & (p,0), \end{cases}$

AAAAAAAAAAAAAABB $\rightarrow 16 \times A, BB \rightarrow EFCBB$ (5)

0123456789AAAAAAAABCDEF $\rightarrow 0123456789$ E74BCDE0F (19)

 $0123456789ABCDEF \rightarrow 0123456789ABCDEOF (17)$



Сравнение (с коррекцией длины)

цепочки в несжатом тексте: (S, L)

Код Зива-Лемпеля, LZ77/LZ1

1977 г., Якоб Зив (Jacob Ziv) и Абрахам Лемпель (Abraham Lempel): если цепочка символов (не обязательно одинаковых) встречается более одного раза, то каждое следующее вхождение цепочки заменяются ссылкой на предыдущее; ссылка состоит из относительного смещения $S\geqslant 1$ и длины L

Поиск предыдущего вхождения — основная сложность алгоритмов семейства LZ77 и неоднозначность кодирования.

Алгоритм LZSS (реализация LZ77) — дерево для ускорения поиска.

Скользящее окно: область перед текущей позицией кодирования, в которой можем искать и адресовать ($w \leq \max(S)$) ссылки Если S записывается N битами, то максимальный размер окна $w = \max(S) = 2^N - 1.$



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

012345678901234567890123456789012345678901234567890123 там $_{\sqcup}$ корабли $_{\sqcup}$ лавировали , $_{\sqcup}$ лавировали , $_{\sqcup}$ да $_{\sqcup}$ не $_{\sqcup}$ вылавировали $\begin{cases} S_1=23-11=12=\mathtt{C}, \\ L_1=13=\mathtt{D} \end{cases}$ принадлежит обеим цепочкам

01234567890123456789012345678901234567890123 там \sqcup корабли \sqcup лавировали, \sqcup лавировали, \sqcup да \sqcup не \sqcup вылавировали $\int S_2 = 44 - 24 = 20 > 15 = {\tt F},$ $L_2 = 10 = {\tt A}$

Простейший вариант кодирования — S и L по одному байту, тогда $1\leqslant S\leqslant 15$ и $1\leqslant L\leqslant 15$; максимальное окно $w=\max S=15$.

При программной реализации (поиск только в окне) будет найдено только одно совпадение: (S_1,L_1) . Возможен случай, когда адресуемых повторов вообще нет, даже односимвольных — как увеличить $\max S$?

там \sqcup корабли \sqcup лавировали, (S_1, L_1) да \sqcup не \sqcup вылавировали $(S_1, L_1) = \mathtt{CD} = \mathtt{op} - \mathtt{как}$ отличить ссылку от текста?

Основные варианты кодирования LZ77

- Наивный ссылки чередуются с несжатыми символами: цепочка, за которой следует c — тройка (S, L, c), несжатый символ c — тройка (0, 0, c). При необходимости текст дополняется (обычно нулями)
- Новопримитивный: ссылка (L,S), несжатый символ c как (0,c).
- LZSS (Сторер, Сжимански): флаг-бит + символ c/ссылка (S, L).
- Флаг-байт: символы c и ссылки (S,L) группируются по 8 (по числу бит в байте-октете) штук, каждая группа предваряется байтом, где каждый бит показывает тип соответствующего объекта:
 - $\left\{ egin{aligned} 0, & \text{несжатый символ} \mathsf{байт} \ c \ 1, & \mathsf{ссылка} \mathsf{два} \ \mathsf{байта} \ (S,L), 0 < S \leqslant w, L \geqslant 3 \end{aligned}
 ight.$
 - для байта-тетрады группируется, соответственно, по 4 объекта
- Префикс p как маркер ссылки: $\begin{cases} \text{ссылка } (p,S,L): \ 0 < S \leqslant w, L \geqslant 4, \\ \text{символ } p \to (p,0) \end{cases}$



```
0123456789ABCDEF0123456789ABCDEF012345
01234567890123456789012345678901234567890123
тамыкораблиылавировали, шлавировали, шдашнецвылавировали (54)
```

- Паивный: 00т00а00м00_□00к00о00р610600л00и81л61в51рС1вС1лВ1, СDдС1_□00н00е61в00ыС1а41и00р00о81а81и (84)
 Порядок полей допустим любой — и (S, L), и (L, S) — здесь (S, L).
- f 2 Новопримитивный: 0 ${\tt T0a0m0_{ll}0k0o0p0a060\pi0u0_{ll}0\pi0a0b0u0p0o0b0a0}$ 2 ${\tt B0_{ll}0m0a0b0_{ll}0m0e0_{ll}0b0u0\pi0a0b0u0p0o0b0a0\pi0u}$ (82), только (L,S)
- $footnote{3}$ Флаг-байт: 0там $_{\Box}$ 0кора0бли $_{\Box}$ 0лави0рова8ли, $_{\Box}^{CD}$ 0да $_{\Box}$ н0е $_{\Box}$ вы0лави 0рова0ли00 (55) Допустимы и (S,L), и (L,S)—в статье SS (S,L).
- Префикс $p = F = \omega$: там_корабли_павировали, вСDда_не_выОлавировали (45)
 Префикс $p - \omega$ в начале. Разный порядок полей S и L приводит к разным способам экранирования символа p. Здесь (p, S, L).

Увеличение разрядности S

Записывать S двумя байтами. Для байта-тетрады $S \leqslant 2^8 - 1 = 255$; для байта-октета $S \leqslant 2^{16} - 1 = 65\,535$, окно w обычно дополнительно ограничивается ($w << \max S$) для ускорения поиска.

Префикс: ссылка занимает 4 символа, тогда $L \geqslant 5$:

```
• \left\{ \begin{array}{ll} \mathsf{Ссылка}\; (S,L) \to (p,\;\; \log S,\;\; \operatorname{hi} S,\;\; L), \\ \mathsf{символ}\; p & \to (p,0,0),\;\; \mathsf{экранируется}\; S=0 \end{array} \right.
\{ ссылка (S,L) 	o (p,\ L,\ \log S,\ \mathrm{hi}\, S), символ p 	o (p,0),\  экранируется L=0
```

 $oldsymbol{ol}}}}}}}}}}}}}}}}$

Префикс: ссылка занимает 3 символа, $L \geqslant 4$:

```
• \left\{ \begin{array}{ll} \mathsf{Ссылка}\; (S,L) \to (p,\;\; \log S,\;\; \operatorname{hi} S + L), \\ \mathsf{символ}\; p & \to (p,0,0),\; \mathsf{экранируется}\; S = 0. \end{array} \right.
egin{array}{ll} {f C} Ссылка (S,L) 
ightarrow (p,\ L+{
m hi}\,S,\ {
m lo}\,S), \ {
m c} символ p 
ightarrow (p,0),\ {
m экранируется}\ L=0
```



- lacktriangle Наивный: $S\geqslant 0$, $L\geqslant 0\Rightarrow$ ни S, ни L невозможно скорректировать.
- $m{2}$ Флаг-байт: $S\geqslant 1$ всегда; если ссылка k символов, то $L\geqslant k+1$; значения 0 не используются как специальные: $egin{cases} S\to\widetilde{S}=S-1, \\ L\to\widetilde{L}=L-k-1. \end{cases}$
- 🗿 Префикс: $S\geqslant 1$ всегда; $L\geqslant k+1$ (k-длина ссылки, включая префикс):
 - ullet если для экранирования символа-префикса используется S=0, корректируется $L\colon L\to \widetilde L=L-k-1$ $\left(\widetilde L\geqslant 0\right).$
 - ullet если для экранирования символа-префикса используется L=0, корректируются обе: $\begin{cases} S & o & \widetilde{S}=S-1 \quad \left(\widetilde{S}\geqslant 0\right), \\ L & o & \widetilde{L}=L-k \quad \left(\widetilde{L}\geqslant 1\right). \end{cases}$



Увеличиваем $S \left\{ \begin{vmatrix} \widetilde{S} \\ \widetilde{L} \end{vmatrix} = |c|+1: \quad 0 \leqslant \widetilde{S} \leqslant 31 \right.$ для октета: $\begin{cases} 0 \leqslant \widetilde{S} \leqslant 511 \\ \widetilde{L} \leqslant \overline{L} \leqslant 127 \end{cases} \right\}$ префикс $p = F = \mathsf{I}$, $\begin{cases} \mathsf{ссылк} \ k = 3 \end{cases}$ символа $(p, \ \widetilde{L} + \mathsf{hi} \ \widetilde{S}, \ \mathsf{lo} \ \widetilde{S}), \ \widetilde{L} \geqslant 1, \end{cases}$

максимальная $\begin{cases} \widetilde{S} = S - 1: & 1 \leqslant S \leqslant 32, \text{ для} \\ \widetilde{L} = L - k = L - 3: & 4 \leqslant L \leqslant 10. \end{cases}$ октета: $\begin{cases} 1 \leqslant S \leqslant 512 \\ 4 \leqslant L \leqslant 130 \end{cases}$ окно w = 32:

там_корабли_лавировали, лавировали, лавировали

$$\begin{cases} S_1 = 12 & \rightarrow & \widetilde{S_1} = 11 & \sim & 0 \ 1011 \\ L_1 = 10 & \rightarrow & \widetilde{L_1} = 7 & \sim & 111 \end{cases} \rightarrow 0111 \ 1011$$

$$0123456789012345678901234567890123456789012345678901234567890123$$
 там $_{\sqcup}$ корабли $_{\sqcup}$ лавировали , $_{\sqcup}$ лавировали , $_{\sqcup}$ лавировали , $_{\sqcup}$ да $_{\sqcup}$ не $_{\sqcup}$ вылавировали

$$\begin{cases} S_2 = 44 - 12 = 32 & \to & \widetilde{S_2} = 31 & \sim & 1 \ 1111 \\ L_2 = 10 & \to & \widetilde{L_2} = 7 & \sim & 111 \end{cases} \to 1111 \ 1011$$

там⊔корабли⊔лавировали,ы7Ви,⊔да⊔не⊔вы0ыFВ (41)

Для октета можно увеличить
$$S$$
 на два бита:
$$\begin{cases} 0\leqslant \widetilde{S}\leqslant 1023\\ 1\leqslant \widetilde{L}\leqslant 63 \end{cases} \Rightarrow \begin{cases} 1\leqslant S\leqslant 1024\\ 4\leqslant L\leqslant 66\\ 2\leqslant C \end{cases}$$

Сравнение RLE и LZ77

Как RLE, так и LZ77 не требуют отдельного словаря.

Повторение одинаковых символов: Префикс p = E, без коррекции

RLE: $AAAAAAAA \sim 8 \times A$ $0123456789AAAAAAAABCDEF \rightarrow 0123456789E8ABCDE0F$

 $23 \rightarrow 19$

LZ77:
$$\overbrace{AAAAAAA}$$
 ~ $A+(S=1,L=7)$ Порядок (p,L,S) 0123456789 $AAAAAAAABCDEF \rightarrow$ 0123456789 $AE71BCDEOF$ 23 \rightarrow 20

Повторение подстрок общего вида:

RIF: $28 \rightarrow 24$

012345601234789AAAAAAABCDEF $\rightarrow 012345601234789$ E8ABCDEOF

1777: $28 \rightarrow 23$

 $012345601234789AAAAAAABCDEF \rightarrow 0123456E57789AE71BCDE0F$

Для большинства файлов LZ77 эффективнее RLE.



Спасибо за внимание!

ТЕИМ

http://miet.ru/

Александра Игоревна Кононова illinc@mail.ru



- Рассчитать количество информации в сообщении $C = \underbrace{\text{ум-лм-лм-...-лм}}_{80 \text{ раз}}$ (кодировка koi8-r)
- $oldsymbol{\circ}$ Закодировать сообщение C алгоритмом LZ77

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

- ① Без памяти, $A_1=\{\mathsf{л},\mathsf{s},\mathsf{-}\}$: $I_1=2\cdot 80\cdot \left(-\log_2\left(\frac{80}{239}\right)\right)+79\cdot \left(-\log_2\left(\frac{79}{239}\right)\right)=254,2\ \mathsf{бит}=31,8\ \mathsf{байт}$
- $m{2}$ Без памяти, $A_1=\{$ ля-, ля $\}$: $I_2=79\cdot\Big(-\log_2\left(rac{79}{80}
 ight)\Big)+1\cdot\Big(-\log_2\left(rac{1}{80}
 ight)\Big)=7,6$ бит =0,9 байт
- Маркова, $A_1=$ koi8-r, N=3 со следующими оценками вероятностей: $p(c_1=\mathbf{n})=p(c_2=\mathbf{s})=p(c_3=-)=\frac{1}{256}, \qquad p(-|-\mathbf{n}\mathbf{s})=\frac{79}{80}, \\ p(\mathbf{n}|\mathbf{n}\mathbf{s}-)=p(\mathbf{s}|\mathbf{s}-\mathbf{n})=1, \qquad p(eof|-\mathbf{n}\mathbf{s})=\frac{1}{80}$:

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



Задача №2

Используя 4-битный байт, закодируйте:

Оообщение (64 = 0х40 символов) — LZ77:

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

```
КОД i | 0123456789ABCDEF AVCII: a_i . аезийклмнопрсь \nu(a_i) | B1431338111C1392
```

- а) S занимает 1 байт (4 бита); б) S занимает 5 бит; в) S занимает 2 байта.
- $oxed{2}$ 16-цветный рисунок 8 imes 8 RLE, LZ77.

