

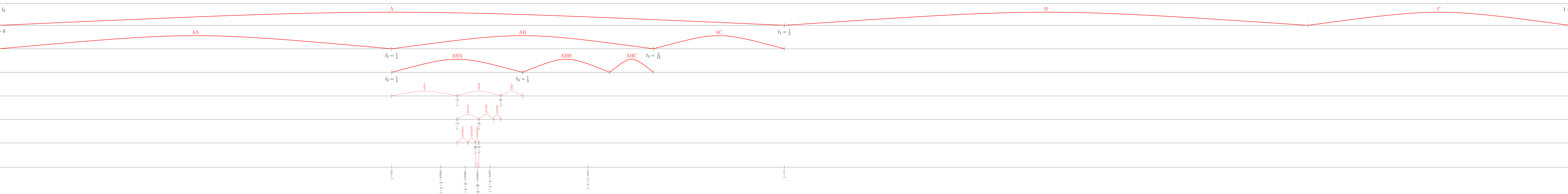
Геометрическая интерпретация арифметического кодирования сообщения $C = \text{АНАНАС}$. Все величины вещественные; расчёты абсолютно точные.

Алфавит $A = \{A, H, C\}$ упорядочен по убыванию частот $\nu = \{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}\}$ (при совпадении частот использовался бы исходный алфавитный порядок).

Начальный рабочий полуинтервал $[t_0, t_1) = [0, 1)$.

На шаге i предыдущий полуинтервал $[t_{i-1}, t_i)$ разбивается на три (по числу символов алфавита) части, длины которых пропорциональны частотам символов. В соответствии с символом c_i сообщения из этих частей выбирается следующий полуинтервал $[t_i, t_{i+1}) \subseteq [t_{i-1}, t_i)$: $[t_i, t_{i+1}) = \begin{cases} [t_{i-1}, & t_{i-1} + \frac{1}{2}\Delta_{i-1}), & c_i = A, \\ [t_{i-1} + \frac{1}{2}\Delta_{i-1}, & t_{i-1} + (\frac{1}{2} + \frac{1}{4})\Delta_{i-1}) = [t_{i-1} + \frac{1}{2}\Delta_{i-1}, & t_{i-1} + \frac{3}{4}\Delta_{i-1}), & c_i = H, \text{ где } \Delta_{i-1} = t_{i-1} - t_{i-1}, \\ [t_{i-1} + \frac{3}{4}\Delta_{i-1}, & t_{i-1} + (\frac{3}{4} + \frac{1}{8})\Delta_{i-1}) = [t_{i-1} + \frac{3}{4}\Delta_{i-1}, & t_{i-1}), & c_i = C. \end{cases}$

Получим последовательность полуинтервалов: $[0, 1) \supseteq [t_1, t_1) \supseteq [t_2, t_2) \supseteq [t_3, t_3) \supseteq [t_4, t_4) \supseteq [t_5, t_5) \supseteq [t_6, t_6)$. ($t_i \in \mathbb{R}, t_i \in \mathbb{R}$). В последнем полуинтервале $[t_6, t_6)$ выбирается точка $z \in [t_6, t_6)$ такая, что её двоичное представление конечно (я самое короткое из всех точек $[t_6, t_6)$); здесь $z = 0,0100111$. Дробная часть z (биты 0100111) – код сообщения z .



Геометрическая интерпретация целочисленной реализации арифметического кодирования сообщения $C = \text{АНАНАС}$ и соответствие целочисленным ℓ и t точкам вещественного полуинтервала $[0, 1)$. Над осью – целочисленные изображения из $[0, N)$ (используется $N = 1024$). Под осью – вещественные значения из $[0, 1)$.

Базис-линей широким линией ■ показана та часть $[Q_0, Q_N)$ вещественного полуинтервала $[0, 1)$, которая в данный момент изображается целочисленным полуинтервалом $[0, N)$. Служи цветом отмечены также биты, совпадающие для всех вещественных точек, изображаемых $[0, N)$.

Для удобства целочисленных расчётов вещественные накопленные частоты $\{\frac{0}{2}, \frac{1}{2}, \frac{1}{4}, 1\} = \{\frac{0}{2}, \frac{2}{4}, \frac{1}{2}, \frac{4}{4}\}$ заменены на совокупность целочисленных накопленных частот $\omega = \{0, 3, 5, 6\}$ и общего знаменателя (делителя) $D = 6$. Непоказанные частоты символов также целочисленные: $\{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}\} = \{\frac{3}{6}, \frac{1}{2}, \frac{1}{8}\} \Rightarrow \nu = \{3, 2, 1\}$ (я $D = \sum_j \nu_j$).

j		0	1	2	3
$\xi_j \in A$		A	H	C	
$\nu_j \in \mathbb{N}$		3	2	1	
$\omega_j \in \mathbb{N} \cup \{0\}$		0	3	5	6

$$D = \sum_{j=0}^T \nu_j = \omega_T = 6, \text{ где } T = 3 - \text{размер алфавита}; \xi_1, \dots, \xi_T \in A - \text{символы алфавита в порядке убывания частот (} \omega_1, \dots, \omega_T - \text{они же в оригинальном алфавитном порядке); } \omega_1, \dots, \omega_T \in \mathbb{N} - \text{целочисленные частоты (для больших файлов нормируются для экономии места в заголовке и достижения } N \gg 4D^2, \text{ но так, что если символ } \xi_i \text{ встретится в файле хотя бы один раз, } \nu_j > 0); \omega_0, \omega_1, \dots, \omega_T \in \mathbb{N} - \text{целочисленные накопленные частоты для расчёта границ символов: } \omega_0 = 0, \text{ даже } \omega_j = \omega_{j-1} + \nu_j.$$

$N = 1024$: $\frac{N}{3} = 512, \frac{N}{2} = 256, \frac{N}{4} = 768$

Целочисленные изображения ℓ и t не являются точными изображениями t_i и t_i вещественной реализации при $i > 1$ (погрешность накапливается), но при декодировании погрешность компенсируется (при $N = 2^n \gg 4D^2$, N должно быть одно и то же про кодирования и декодирования).

Начальный рабочий целочисленный полуинтервал $[t, t) = [0, N)$, и вначале он соответствует вещественному $[0, 1)$.

На шаге i выполняется:

– сокращение полуинтервала с учётом символа $c_i = \xi_j$ (то есть порядковый номер символа в сообщении i , а номер этого же символа в переупорядоченном по убыванию частот алфавите j) – предыдущий полуинтервал $[t, t)$ разбивается на части, длины которых пропорциональны частотам символов, и выбирается j -я, границы которой могут быть рассчитаны по формулам:

$$\Delta = t - t_i, \begin{cases} t \mapsto t + \frac{\Delta \omega_{j-1}}{D}, \\ t \mapsto t + \frac{\Delta \omega_j}{D}. \end{cases}$$

– одно или несколько (или ни одного) масштабирований (увеличений масштаба) – изменений вещественных границ $[Q_0, Q_N)$ целочисленного изображения $[0, N)$. Размер $[Q_0, Q_N)$ при каждом масштабировании уменьшается вдвое; расстояние на его изображения $[0, N)$ – вдвое увеличивается.

