

О сжатии данных массивов силовых кривых

А. И. Стефанович, Д. В. Сушко

*Институт проблем информатики Федерального исследовательского центра
«Информатика и управление» Российской академии наук, Москва, Россия*

Поступила в редколлегию 03.09.2020

Аннотация—Рассмотрена задача обратимого (без потерь) сжатия данных массивов силовых кривых, получаемых в результате исследования биологических объектов на атомно-силовом микроскопе. Построены оценки скорости кодирования при использовании для сжатия стандартных алгоритмов (DEFLATE, JPEG 2000). Предложены алгоритмы сжатия, основанные на применении универсального арифметического кодирования, и построены оценки скорости кодирования для этих алгоритмов.

КЛЮЧЕВЫЕ СЛОВА: обратимое сжатие данных; универсальное арифметическое кодирование; атомно-силовой микроскоп; массив силовых кривых

1. ВВЕДЕНИЕ

Атомно-силовой микроскоп (АСМ) был изобретен в 1986 году [1] и применяется главным образом для измерения рельефа поверхности с нанометровым разрешением. Принцип работы АСМ заключается в сканировании поверхности образца атомарно острой иглой (зондом), которая является частью гибкого кронштейна (кантилевера), закрепленного на пьезоэлектрическом двигателе. Силы, действующие на зонд со стороны поверхности, приводят к изгибу кантилевера, в результате чего происходит перераспределение лазерного сигнала на фотодетекторе. Регистрируя величину перераспределения сигнала и зная жесткость кантилевера, можно определить силу взаимодействия зонда с поверхностью, что позволяет построить изображение рельефа поверхности образца.

С момента изобретения АСМ в рамках того же базового физического принципа были разработаны различные методики исследования физических и химических свойств поверхностей, таких как механическая жесткость, электропроводность, намагниченность, температура и других. В настоящее время АСМ широко используется в качестве инструмента исследования в физике, химии, биологии, материаловедении и других науках.

Современный АСМ имеет множество различных режимов работы (сканирования). Многие приборы позволяют работать в режиме измерения силовых кривых и силовых карт. Данный режим широко используется, например, в микробиологии при исследовании клеток, вирусов, белков, нуклеиновых кислот и т.д. Силовая кривая представляет собой, по существу, график зависимости силы взаимодействия зонда и поверхности образца от расстояния между ними. При фиксированном положении зонда в плоскости образца XU снимаются две кривые: кривая подвода (зонд приближается к образцу) и кривая отвода (зонд удаляется от образца). Пары силовых кривых снимаются для множества точек в поле наблюдения. В результате формируется карта силовых кривых, которая представляет собой трехмерный массив данных. В типичном случае каждая силовая кривая включает 1024 отсчета по 2 байта (16 битов) на один отсчет, а число точек в поле наблюдения, для которых снимаются силовые кривые, составляет 128×128 . Поэтому полный объем данных одного массива силовых кривых равен $1024 \times 2 \times 2 \times 128 \times 128$ Б

~ 67 МБ, т.е. достаточно велик. Необходимость долгосрочного хранения и передачи таких данных делают задачу их сжатия весьма актуальной. При этом, поскольку получение данных связано с проведением уникального, трудоемкого и длительного по времени эксперимента, потери при сжатии недопустимы, т.е. сжатие должно быть обратимым.

Цель настоящей работы – исследовать потенциальные возможности ряда алгоритмов обратимого сжатия применительно к массивам силовых кривых, т.е. получить оценки для скорости кодирования этих массивов данных посредством таких алгоритмов. Напомним, что *скоростью кодирования* R (средней скоростью кодирования) алгоритма называется отношение длины кодового слова L (в битах), порождаемого алгоритмом для описания массива данных, к полному числу элементов (пикселей) N этого массива; единица измерения скорости кодирования – бит/пиксель (бт/п). Используемые в работе экспериментальные данные получены при сканировании мягких биологических образцов в режиме измерения силовых карт на микроскопе MultiMode V (Veeco, США). Вычисления произведены программами, написанными на языке Python.

Работа имеет следующую структуру. В разделе 2 приведены сведения об используемых экспериментальных образцах, некоторые необходимые детали, относящиеся к режиму сканирования, и описана структура подлежащих сжатию массивов силовых кривых. В разделе 3 для сжатия массивов данных применены стандартные алгоритмы (DEFLATE и JPEG 2000) и приведены полученные при этом оценки скорости кодирования. В разделе 4 рассмотрен ряд алгоритмов сжатия, основанных на применении универсального арифметического кодирования, и получены оценки для скорости кодирования рассматриваемых массивов данных этими алгоритмами.

2. МАССИВЫ СИЛОВЫХ КРИВЫХ

В качестве экспериментального материала в данной работе используются пять массивов силовых кривых (I–V), полученных в результате сканирования образцов, представляющих собой абсорбированные из раствора на твердую подложку вирусы. Первый образец (I) – это риновирус 2 на подложке из слюды, второй и третий образцы (II, III) – вирус мягкой мозаики ячменя на подложке из слюды, четвертый и пятый образцы (IV, V) – вирус табачной мозаики на стеклянной подложке.

Рассмотрим в необходимом для дальнейшего изложения объеме вопросы, связанные с процессом сканирования и структурой массивов силовых кривых. Подробное описание технологии измерений с помощью АСМ и интерпретации соответствующих данных можно найти, например, в обзоре [2].

Введем некоторые обозначения. Пусть $OXYZ$ – трехмерная декартова система координат. Твердая подложка образца располагается в горизонтальной плоскости OXY , поле наблюдения представляет собой прямоугольник $[0, \bar{x}] \times [0, \bar{y}]$ в этой плоскости. Силовые кривые измеряются в узлах (x, y) прямоугольной решетки в поле наблюдения: $(x, y) \equiv (x_i, y_j)$, где $x_i = i\Delta x$, $i = 0, 1, \dots, (I - 1)$ и $y_j = j\Delta y$, $j = 0, 1, \dots, (J - 1)$. Размеры поля наблюдения \bar{x} , \bar{y} (ширина, длина), шаги решетки Δx , Δy и числа отсчетов по ширине и длине I , J являются параметрами эксперимента. Линейные размеры поля наблюдения, шаги решетки и числа отсчетов связаны, разумеется, условиями $\bar{x} = (I - 1)\Delta x$, $\bar{y} = (J - 1)\Delta y$. Перебор узлов решетки осуществляется в следующем порядке: сначала по ширине (в направлении OY), затем по длине (в направлении OX).

В каждом узле (x, y) решетки измеряется пара силовых кривых $F_{(x,y)}^A(z)$ (кривая подвода) и $F_{(x,y)}^R(z)$ (кривая отвода), $z \equiv z_k = k\Delta z$, $k = 0, 1, \dots, (K - 1)$, причем высота z отсчитывается от поверхности образца в узле (x, y) поля наблюдения. Вертикальный размер \bar{z} области, в

которой измеряются силовые кривые, шаг измерения по вертикали Δz и число отсчетов K связаны условием $\bar{z} = (K - 1)\Delta z$ и являются параметрами эксперимента.

Значения элементов силовых кривых $F_{(x,y)}^{A,R}(z)$ пропорциональны силе, действующей на зонд в точке пространства, находящейся на расстоянии z от поверхности образца и имеющей координаты (x, y) в поле наблюдения, в процессе подвода зонда к образцу и отвода зонда от образца. Значения F записываются в виде 16 битных целых чисел (short int), т.е. целых чисел в диапазоне $[-2^{15}, (2^{15} - 1)]$. Теоретически положительные значения F соответствуют отталкиванию зонда от поверхности, отрицательные – притяжению. На практике, однако, начальная калибровка нуля выполняется не идеально и наблюдается значительный дрейф нуля в процессе измерений, поэтому точнее говорить о том, что увеличение значения F отвечает увеличению отталкивания или, что то же самое, уменьшению притяжения.

Измерение пары силовых кривых осуществляется следующим образом. Первой снимается кривая подвода. В начальном положении зонд находится далеко от поверхности образца, т.е. сила взаимодействия зонда с поверхностью заведомо равна нулю. Пьезоэлектрический двигатель осуществляет перемещение зонда вниз в направлении к поверхности с шагом Δz . При этом регистрируется сигнал на выходе фотодетектора, изменение которого пропорционально изменению силы взаимодействия зонда с поверхностью. Движение зонда к поверхности осуществляется до тех пор, пока изменение сигнала не превысит заданного порогового значения, что интерпретируется как достижение зондом поверхности образца. Кривая подвода $F_{(x,y)}^A(z_k)$, $k = 0, 1, \dots, (K - 1)$, состоит из последних K зарегистрированных значений, записанных в строку в порядке, обратном к порядку регистрации. Таким образом, нумерация элементов силовой кривой подвода осуществляется в направлении от поверхности. Возможна ситуация, когда зонд достигнет поверхности до того, как будет зарегистрировано необходимое число (K) значений. В таком случае все зарегистрированные значения записываются в начало строки, и осуществляется ее дополнение до требуемой длины (K) минимально возможным значением -2^{15} . Вслед за кривой подвода снимается кривая отвода. Зонд перемещается от поверхности вверх до заведомо «свободного» положения с шагом Δz . Кривая отвода $F_{(x,y)}^R(z_k)$, $k = 0, 1, \dots, (K - 1)$, состоит из первых K зарегистрированных значений, записанных в строку в порядке регистрации. Таким образом, нумерация элементов силовой кривой отвода также осуществляется в направлении от поверхности. Необходимость в дополнении строки в данном случае не возникает.

Для всего массива силовых кривых будем использовать обозначение $\mathbf{V} = \{V(i, j, k)\}$, $i = 0, 1, \dots, (I - 1)$, $j = 0, 1, \dots, (J - 1)$, $k = 0, 1, \dots, (2K - 1)$. При этом

$$V(i, j, k) = \begin{cases} F_{(x_i, y_j)}^A(z_k), & k = 0, 1, \dots, (K - 1); \\ F_{(x_i, y_j)}^R(z_{k-K}), & k = K, (K + 1), \dots, (2K - 1). \end{cases}$$

Одномерное упорядочивание трехмерных массивов вообще и массивов силовых кривых в частности осуществляется следующим образом: сначала по оси Z (по высоте), затем по оси Y (по длине), затем по оси X (по ширине). Для массива силовых кривых в результате упорядочивания имеем $V(n) = V(i, j, k)$, $n = 0, 1, \dots, (N - 1)$, причем $n = J \cdot 2K \cdot i + 2K \cdot j + k$, а полное число элементов массива равно $N = I \cdot J \cdot 2K$.

В табл. 1 приведены значения описанных выше параметров для используемых в работе массивов силовых кривых.

На рис. 1 представлены пары силовых кривых образца II в узлах поля наблюдения с номерами $i = 9, j = 13$ (а) и $i = 76, j = 33$ (б). Кривые подвода изображены черным цветом, кривые отвода – серым.

Таблица 1. Параметры массивов силовых кривых

Параметр	I	II	III	IV	V
Размер по X (ширина \bar{x})	755,928 нм	1000 нм	682,8 нм	1517,92 нм	1517,92 нм
Размер по Y (длина \bar{y})	755,928 нм	1000 нм	682,8 нм	1517,92 нм	1517,92 нм
Размер по Z (высота \bar{z})	80 нм	70 нм	70 нм	100 нм	100 нм
Число отсчетов по $X(I)$	64	128	128	128	128
Число отсчетов по $Y(J)$	64	128	128	128	128
Число отсчетов по $Z(K)$	4096	1024	1024	1024	1024
Полное число отсчетов (N)	2^{25}	2^{25}	2^{25}	2^{25}	2^{25}
Пороговое значение	0,1 В	0,05 В	0,05 В	0,15 В	0,15 В

На рис. 1 (б) представлена кривая подвода, при измерении которой зонд достиг поверхности образца до того, как было зарегистрировано необходимое число ($K = 1024$) значений. Поэтому в конец строки было записано нужное число минимально возможных значений (-2^{15}).

Горизонтальные в среднем участки силовых кривых (в области $k > 300$ в случае (а) и $k > 200$ в случае (б)) соответствуют отсутствию взаимодействия зонда с поверхностью. Соответствующие значения F равны приблизительно -5125 (а), -6540 (б) и значительно отличаются от нуля и между собой. Это иллюстрирует неточность начальной калибровки нуля и дрейф нуля, о чем было сказано ранее.

3. СЖАТИЕ ПОСРЕДСТВОМ СТАНДАРТНЫХ АЛГОРИТМОВ

Динамический диапазон значений массива силовых кривых $[V_{\min}, V_{\max}]$ (т. е. диапазон значений без учета значения -2^{15} , используемого для дополнения «неполных» строк) заметно уже полного диапазона 16-битных целых чисел. Для экспериментальных массивов I–V этот динамический диапазон составляет $[-2416, -404]$, $[-7832, -4744]$, $[-6605, -1156]$, $[-3921, -311]$ и $[-3544, +187]$ соответственно, что позволяет уменьшить скорость кодирования при использовании равномерного кода. Действительно, применим к массиву следующее преобразование:

$$\mathbf{V} \rightarrow \mathbf{V}', \quad V' = \begin{cases} 0, & V = -2^{15}, \\ V - V_{\min} + 1, & V > -2^{15}. \end{cases}$$

Скорость кодирования массива \mathbf{V}' равномерным кодом равна $R = \lceil \log(V_{\max} - V_{\min} + 2) \rceil$. Здесь и далее $\log(\cdot) = \log_2(\cdot)$ – двоичный логарифм, $\lceil \cdot \rceil$ – результат округления вещественного числа до ближайшего целого вверх (ceil). В случае массивов I–V имеем для скорости кодирования 11, 12, 13 и 12 бт/п соответственно. Для восстановления исходного массива декодеру

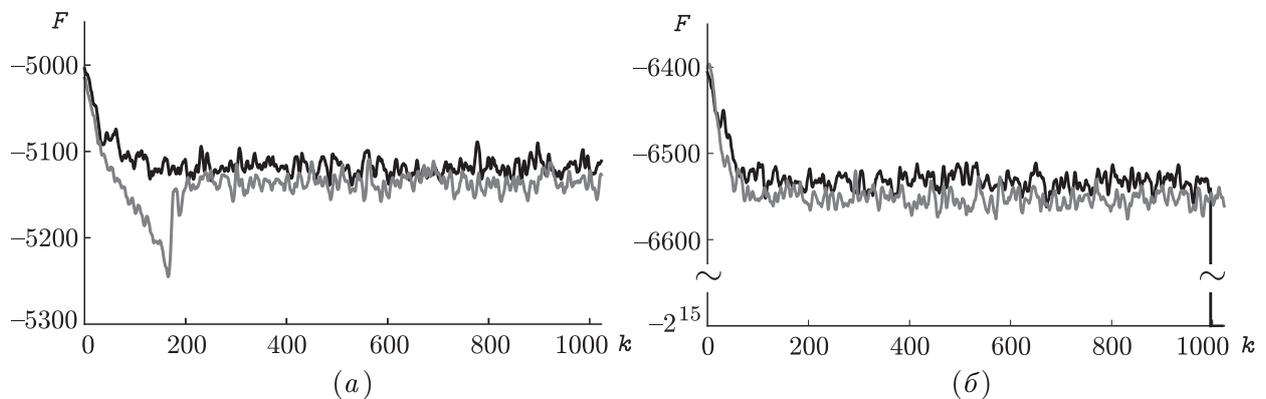


Рис. 1. Силовые кривые образца II в узлах $i = 9, j = 13$ (а) и $i = 76, j = 33$ (б)

нужно дополнительно передать значение V_{dmin} . Это требует еще 16 битов, что приводит к пренебрежимо малому увеличению скорости кодирования на величину $16/N = 2^{-21}$ бт/п.

В качестве первого стандартного алгоритма кодирования рассмотрим алгоритм кодирования DEFLATE (используемый, например, в формате ZIP), который представляет собой комбинацию алгоритмов LZ77 [3] и Хаффмана [4]. Применение реализации, входящей в дистрибутив Anaconda Python, дает представленные в первой строке табл. 2 скорости кодирования массивов I–V. (Здесь и далее значения скорости кодирования приводятся в единицах бт/п с точностью до четырех знаков после десятичной запятой.)

Таблица 2. Скорость кодирования стандартными алгоритмами

Алгоритм	I	II	III	IV	V
DEFLATE	6,3324	5,7357	5,6937	7,5882	7,3255
JPEG 2000 Z	6,6651	6,3425	6,1670	7,0269	6,5941
JPEG 2000 Y	4,1225	3,8421	3,7857	4,4637	4,3821
JPEG 2000 X	4,1023	3,7235	3,6614	4,4217	4,3363

Используем теперь алгоритм кодирования без потерь из множества алгоритмов, входящих в стандарт JPEG 2000 (подробное описание алгоритмов стандарта можно найти, например, в части II книги [5]). Данный алгоритм предназначен для кодирования изображений, т.е. по существу неотрицательных двумерных массивов, и может быть применен для кодирования трехмерных массивов силовых кривых следующим образом. Выберем опорное направление вдоль одной из осей координат (Z , Y или X), будем рассматривать значения массива как 16 битные неотрицательные целые (unsigned short int) и осуществим кодирование всех двумерных сечений массива, перпендикулярных опорной оси. Применение реализации алгоритма кодирования без потерь стандарта JPEG 2000, входящей в дистрибутив Anaconda Python, дает представленные в строках 2–4 табл. 2 скорости кодирования массивов I–V (опорная ось, перпендикулярная кодируемым сечениям, указана в заглавной колонке таблицы). Обратим внимание на общий для всех массивов характер зависимости скорости кодирования от выбора опорной оси. Минимальная скорость кодирования наблюдается для опорной оси X , несколько большая скорость – для оси Y и значительно большая скорость – для оси Z . Дело в том, что алгоритм использует, в числе прочего, статистическую зависимость значений соседних элементов, которая максимальна в случае сечений, перпендикулярных оси X и практически отсутствует в случае сечений, перпендикулярных оси Z . Это, в свою очередь, обусловлено структурой массива силовых кривых и порядком сканирования.

4. СЖАТИЕ ПОСРЕДСТВОМ АРИФМЕТИЧЕСКОГО КОДИРОВАНИЯ

4.1. Сведения об арифметическом кодировании

Традиционная схема решения задач обратимого сжатия цифровых данных заключается в применении к исходному массиву данных некоторого обратимого преобразования, обеспечивающего декорреляцию его отсчетов и/или уменьшение диапазона их значений, и последующего кодирования полученного таким образом массива как последовательности независимых отсчетов. В качестве метода кодирования в настоящее время обычно используется арифметическое кодирование (см., например, [6]). Обратимые преобразования, используемые в настоящей работе, описаны далее. Здесь мы приведем необходимые сведения, касающиеся арифметического кодирования.

Рассмотрим конечное множество $\mathfrak{A} = \{a\}$ (алфавит), буквы которого представляют собой целые числа из некоторого известного априори диапазона $[a_-, a_+]$. Пусть $\mathbf{x} = \{x_n\}_{n=0}^{N-1}$ ($N \gg 1$)

– конечная последовательность букв алфавита \mathfrak{A} , $x_n \in \mathfrak{A}$, подлежащая кодированию. Процесс кодирования заключается в последовательном просмотре всех значений x_n , вычислении *кодовой вероятности* $Q(\mathbf{x})$ (положительного вещественного числа, не превышающего единицы) для всей последовательности и формировании по этой кодовой вероятности двоичного кодового слова (результата сжатия) длины $L(\mathbf{x}) = \lceil -\log Q(\mathbf{x})/2 \rceil$ битов. Вычисление кодовой вероятности осуществляется рекуррентно. Начальная кодовая вероятность выбирается равной единице ($Q_{-1} = 1$). В момент поступления на вход кодера очередного значения x_n кодеру должно быть известно (задано заранее и/или сформировано в процессе кодирования предыдущих элементов) *условное кодовое распределение вероятностей* $\{q_n(a) = q_n(a|x_{n-1}, \dots, x_0), a \in \mathfrak{A}\}$. Условное кодовое распределение представляет собой набор неотрицательных вещественных чисел, сумма которых по всем $a \in \mathfrak{A}$ равна единице. Кроме того, равенство $q_n(a) = 0$ для некоторого конкретного значения a допустимо только в том случае, если выполнение равенства $x_n = a$ невозможно априори. Шаг рекурсии заключается в умножении текущей кодовой вероятности на условную кодовую вероятность текущего значения x_n : $Q_n = Q_{n-1} \cdot q_n(x_n)$. Результатом выполнения N шагов рекурсии является вычисление кодовой вероятности всей последовательности $Q(\mathbf{x}) = Q_{N-1}(\mathbf{x})$, длина которой равна:

$$L(\mathbf{x}) = \left\lceil -\log \left(\frac{Q_{N-1}}{2} \right) \right\rceil = \left\lceil -\log \prod_{n=0}^{N-1} q_n(x_n) + 1 \right\rceil \leq \sum_{n=0}^{N-1} -\log q_n(x_n) + 2. \quad (1)$$

Детали процедуры формирования кодового слова по кодовой вероятности не принципиальны для нашего рассмотрения и опускаются.

Восстановление исходной последовательности по кодовому слову осуществляется декодером последовательно и без задержки. В момент восстановления очередного значения x_n декодеру уже известны все предыдущие значения $\{x_0, \dots, x_{n-1}\}$ и, кроме того, должно быть известно условное кодовое распределение $\{q_n(a) = q_n(a|x_{n-1}, \dots, x_0), a \in \mathfrak{A}\}$, использованное ранее в процессе кодирования. Это позволяет декодеру восстановить значение x_n .

Таким образом, ключевую роль в процессе арифметического кодирования играют условные кодовые распределения вероятностей $\{q_n(a) = q_n(a|x_{n-1}, \dots, x_0), a \in \mathfrak{A}, n = 0, \dots, N-1\}$, выбор которых определяет длину кодового слова, т.е. степень сжатия исходных данных. При этом построение условных кодовых распределений, обеспечивающих получение возможно более коротких кодовых слов для входных данных с неизвестной (не полностью известной) статистикой, – задача универсального кодирования [7].

В работе [8] показано, что для источников без памяти (последовательности статистически независимых отсчетов) почти оптимальным является использование на n -ом шаге арифметического кодирования/декодирования следующего условного кодового распределения вероятностей:

$$q_n(x) = \frac{1 + 2\theta_n(x)}{|\mathfrak{A}| + 2n}, \quad n = 0, 1, \dots, (N-1),$$

где $\theta_n(x)$ – число вхождений буквы x (т.е. число элементов, принимающих значение x) в начальном участке последовательности \mathbf{x} до $(n-1)$ -го члена включительно, а $|\mathfrak{A}|$ – размер алфавита. Отметим, что поскольку декодирование осуществляется без задержки, на n -ом шаге декодирования все предыдущие значения уже известны, что позволяет декодеру построить условное кодовое распределение вероятностей $\{q_n(x), x \in \mathfrak{A}\}$, использованное при кодировании, без дополнительной информации.

В том случае, когда алфавит – это целые числа из априори известного диапазона $[a_-, a_+]$, а кодируемые последовательности \mathbf{x} принимают значения в диапазоне $[a_{\min}(\mathbf{x}), a_{\max}(\mathbf{x})]$, который для всех встречающихся последовательностей заметно уже полного диапазона, выгодно

модифицировать приведенный выше способ построения условных кодовых распределений следующим образом:

$$q_n(x) = \begin{cases} \frac{1 + 2\theta_n(x)}{a_{\max}(\mathbf{x}) - a_{\min}(\mathbf{x}) + 1 + 2n} & \text{при } x \in [a_{\min}(\mathbf{x}), a_{\max}(\mathbf{x})]; \\ 0 & \text{в противном случае.} \end{cases} \quad (2)$$

Это позволяет учесть то обстоятельство, что конкретная последовательность \mathbf{x} не принимает значений вне диапазона $[a_{\min}(\mathbf{x}), a_{\max}(\mathbf{x})]$. Значения границ диапазона могут быть вычислены кодером и должны быть, разумеется, переданы декодеру помимо кодового слова, что, вообще говоря, несколько увеличивает скорость кодирования. Однако встречающиеся в работе последовательности имеют длину не менее 2^{23} , а для передачи значений двух границ требуется не более 2^6 битов, поэтому увеличение скорости кодирования не превышает 2^{-17} бт/п – величина, которой можно пренебречь.

В настоящей работе условные кодовые распределения вероятностей строятся по формуле (2), а соответствующая скорость кодирования оценивается по формуле

$$R(\mathbf{x}) = \frac{1}{N}L(\mathbf{x}) = \frac{1}{N} \sum_{n=0}^{N-1} -\log q_n(x_n), \quad (3)$$

которая следует из формулы (1), если пренебречь бесконечно малым членом порядка $2N^{-1}$.

Величины $\theta(x)/N$, $x \in \mathfrak{A}$, где $\theta(x)$ – число вхождений буквы x в последовательность \mathbf{x} , образуют частотное (или эмпирическое) распределение вероятностей значений последовательности. Величина

$$H(\mathbf{x}) = \sum_{x \in \mathfrak{A}} \frac{\theta(x)}{N} \left[-\log \frac{\theta(x)}{N} \right], \quad (4)$$

(используется соглашение о том, что $0 \cdot \log 0 = 0$) называется квазиэнтропией (эмпирической энтропией) последовательности. Единица измерения квазиэнтропии – бт/п. Квазиэнтропия зависит только от самой последовательности и является нижней границей скорости арифметического кодирования (см., например, [7]). Разность $(H - R) \geq 0$ – избыточность (индивидуальная избыточность) арифметического кодирования. Величина избыточности характеризует качество решения одной из задач универсального кодирования – задачи построения условных кодовых распределений вероятностей.

4.2. Кодирование разностей с предсказанием

В данном разделе в качестве преобразования, обеспечивающего декорреляцию элементов массива силовых кривых, используется переход к разности с предсказанием. Для данного массива $\mathbf{X} = \{X(n), n = 0, 1, \dots, N - 1\}$ (изначально одномерного или одномерно упорядоченного многомерного) переход к разности с простейшим предсказанием по предыдущему элементу имеет вид

$$\mathcal{D}_0 : \mathbf{X} \rightarrow \mathbf{D}_0, \quad D_0(0) = X(0), \quad D_0(n) = X(n) - X(n - 1), \quad n = 1, \dots, N - 1. \quad (5)$$

Преобразование, очевидно, является обратимым.

Применим преобразование \mathcal{D}_0 к одномерно упорядоченному массиву силовых кривых \mathbf{V} и оценим по формулам (3) и (4) скорость арифметического кодирования и квазиэнтропию. Данный вариант обработки обозначим как $\mathcal{P}(\mathcal{D}_0)$. Полученные для массивов I–V результаты приведены в соответствующей строке табл. 3. Избыточность кодирования во всех случаях составляет порядка $\sim 0,01$.

Таблица 3. Скорость кодирования и квазиэнтропия разностей с предсказанием

Вариант обработки	Скорость кодирования					Квазиэнтропия				
	I	II	III	IV	V	I	II	III	IV	V
$\mathcal{P}(\mathcal{D}_0)$	3,9631	3,7128	3,7168	4,3309	4,2647	3,9521	3,7031	3,7061	4,3198	4,2536
$\mathcal{P}(\mathcal{D}_0\mathcal{C}\mathcal{S})$	3,9528	3,7030	3,7051	4,3191	4,2546	3,9519	3,7004	3,7028	4,3174	4,2531
$\mathcal{P}(\mathcal{D}_1\mathcal{C}\mathcal{S})$	3,9504	3,6944	3,6959	4,3103	4,2463	3,9496	3,6922	3,6945	4,3091	4,2454

Попробуем уменьшить скорость кодирования за счет дополнительного применения двух преобразований предварительной обработки. Первое такое преобразование \mathcal{S} разделяет весь массив силовых кривых \mathbf{V} размерами $(I, J, 2K)$ на массив кривых подвода \mathbf{V}^A и массив кривых отвода \mathbf{V}^R размерами (I, J, K) каждый:

$$\mathcal{S} : \mathbf{V} \rightarrow \{\mathbf{V}^A, \mathbf{V}^R\}, \quad V^A(i, j, k) = V(i, j, k), \quad V^R(i, j, k) = V(i, j, K + k), \quad k = 0, 1, \dots, K - 1, \quad (6)$$

что позволяет далее работать с массивами подвода и отвода по отдельности. Преобразование \mathcal{S} разделения массива, естественно, является обратимым.

Второе преобразование \mathcal{C} сужает диапазон значений массива кривых подвода \mathbf{V}^A , непомерно широкий из-за наличия значения -2^{15} , используемого для дополнения «неполных» строк:

$$\mathcal{C} : \mathbf{V}^A \rightarrow \mathbf{C}^A, \quad C^A = \begin{cases} V_{\text{dmin}}^A - 1, & V^A = -2^{15}, \\ V^A, & V^A > -2^{15}, \end{cases} \quad (7)$$

где V_{dmin}^A – динамический минимум значений массива \mathbf{V}^A (т.е. минимум значений без учета значения -2^{15}). Теоретически нельзя исключить возможность того, что все строки в массиве подвода окажутся «полными». Поэтому декодеру, чтобы обратить преобразование сужения диапазона \mathcal{C} , нужна информация о том, встречалось ли значение -2^{15} в исходном массиве кривых подвода. Для передачи декодеру этой известной кодеру информации требуется один бит; соответствующее увеличение скорости кодирования пренебрежимо мало.

Выполним следующую последовательность преобразований. Разделим исходный массив \mathbf{V} преобразованием (6) на массивы подвода и отвода: $\mathcal{S} : \mathbf{V} \rightarrow \{\mathbf{V}^A, \mathbf{V}^R\}$. Сузим диапазон значений массива подвода преобразованием (7): $\mathcal{C} : \mathbf{V}^A \rightarrow \mathbf{C}^A$. Осуществим одномерное упорядочивание двух полученных массивов и применим к ним преобразование перехода к разности с предсказанием (5): $\mathcal{D}_0 : \mathbf{C}^A \rightarrow \mathbf{D}_0^A, \mathcal{D}_0 : \mathbf{V}^R \rightarrow \mathbf{D}_0^R$. Оценим скорость кодирования и квазиэнтропию пары массивов $\mathbf{D}_0^A, \mathbf{D}_0^R$. Поскольку массивы имеют одинаковое число элементов, общая скорость кодирования и квазиэнтропия пары массивов могут быть вычислены как средние значения скоростей кодирования и квазиэнтропий отдельных массивов:

$$R(\{\mathbf{D}_0^A, \mathbf{D}_0^R\}) = \frac{1}{2} [R(\mathbf{D}_0^A) + R(\mathbf{D}_0^R)], \quad H(\{\mathbf{D}_0^A, \mathbf{D}_0^R\}) = \frac{1}{2} [H(\mathbf{D}_0^A) + H(\mathbf{D}_0^R)], \quad (8)$$

а слагаемые в квадратных скобках вычисляются по формулам (3) и (4). Данный вариант обработки обозначим как $\mathcal{P}(\mathcal{D}_0\mathcal{C}\mathcal{S})$. Полученные для массивов силовых кривых I–V результаты представлены в соответствующей строке табл. 3.

По сравнению вариантом обработки $\mathcal{P}(\mathcal{D}_0)$ уменьшение скорости кодирования составляет приблизительно 0,01 бт/п вне зависимости от конкретного массива, а уменьшение квазиэнтропии очень мало – оно едва превышает 0,003 бт/п для массива III и еще меньше для остальных массивов. Таким образом, уменьшение скорости кодирования достигнуто за счет уменьшения избыточности кодирования, которое в данном случае составляет приблизительно 0,001–0,002 бт/п в зависимости от массива, т.е. в 5–10 раз меньше, чем в предыдущем случае. Уменьшение избыточности объясняется применением преобразования \mathcal{C} , в результате

чего условные кодовые распределения (2) для массива кривых подвода оказываются ближе к частотному распределению. В то же время применение преобразования \mathcal{S} , которое в принципе должно уменьшать квазиэнтропию, оказывается малоэффективным.

Модифицируем теперь используемое предсказание с тем, чтобы более аккуратно обрабатывать начальные элементы вертикальных строк трехмерного массива. Для данного массива $\mathbf{X} = \{X(i, j, k)\}$, $i = 0, 1, \dots, (I-1)$, $j = 0, 1, \dots, (J-1)$, $k = 0, 1, \dots, (K-1)$, определим переход к разности с предсказанием $\mathcal{D}_1 : \mathbf{X} \rightarrow \mathbf{D}_1$ следующей формулой:

$$D_1(i, j, k) = \begin{cases} X(0, 0, 0), & i = j = k = 0; \\ X(i, 0, 0) - X(i-1, 0, 0), & i = 1, \dots, I-1, j = k = 0; \\ X(i, j, 0) - X(i, j-1, 0), & i = 1, \dots, I-1, j = 1, \dots, J-1, k = 0; \\ X(i, j, k) - X(i, j, k-1), & i = 1, \dots, I-1, j = 1, \dots, J-1, k = 1, \dots, K-1. \end{cases} \quad (9)$$

Преобразование, очевидно, является обратимым.

Выполним следующую последовательность преобразований. Применим последовательно преобразования разделения массива \mathcal{S} и сужения диапазона \mathcal{C} к исходному массиву силовых кривых \mathbf{V} так, как это было сделано в предыдущем варианте обработки. Применим к двум полученным массивам преобразование перехода к разности с предсказанием (9): $\mathcal{D}_1 : \mathbf{C}^A \rightarrow \mathbf{D}_1^A$, $\mathcal{D}_1 : \mathbf{V}^R \rightarrow \mathbf{D}_1^R$. Осуществим одномерное упорядочивание двух полученных массивов и оценим скорость арифметического кодирования и квазиэнтропию. Поскольку массивы имеют одинаковое число элементов, общая скорость кодирования и квазиэнтропия пары массивов могут быть вычислены как средние значения скоростей кодирования и квазиэнтропий отдельных массивов, т.е. по формулам (8), в которых нужно лишь формально заменить нижней индекс (0 на 1). Данный вариант обработки обозначим как $\mathcal{P}(\mathcal{D}_1\mathcal{C}\mathcal{S})$. Результаты, полученные для массивов силовых кривых I–V, представлены в соответствующей строке табл. 3.

По сравнению с вариантом обработки $\mathcal{P}(\mathcal{D}_0\mathcal{C}\mathcal{S})$ уменьшение как скорости кодирования, так и энтропии составляет приблизительно 0,002 бт/п для массива I и 0,01 бт/п для остальных массивов, а избыточность кодирования составляет приблизительно 0,001–0,002 бт/п в зависимости от массива, как и ранее. Достигнутый прогресс объясняется не только модификацией преобразования перехода к разности с предсказанием ($\mathcal{D}_0 \rightarrow \mathcal{D}_1$), но и тем, что преобразование \mathcal{D}_1 применяется в совокупности с преобразованием разделения массива \mathcal{S} .

Помимо представленных выше вариантов кодирования, в ходе работы был исследован также целый ряд вариантов, использующих и другие виды предсказания, и другие преобразования предобработки. Однако результаты, достигнутые в варианте обработки $\mathcal{P}(\mathcal{D}_1\mathcal{C}\mathcal{S})$, не удалось сколько-нибудь заметно улучшить.

4.3. Кодирование компонент дискретного вейвлет-преобразования

В данном разделе рассматриваются преобразования, включающие обратимое одномерное дискретное вейвлет-преобразование (ДВП) по системе вейвлетов (5–3) [9], реализуемое посредством лифтинг-схемы [10]. Это преобразование входит в стандарт сжатия JPEG 2000 [5].

Опишем кратко, что представляет собой такое ДВП в случае последовательности. Пусть $\mathbf{x} = \{x(k)\}$, $k = 0, \dots, K-1$, – конечная целочисленная последовательность. Будем полагать, что число членов последовательности K четно и не менее четырех; в дальнейшем эти условия всегда выполнены. Преобразование заключается в разложении исходной последовательности на две последовательности (компоненты) $\mathcal{W} : \mathbf{x} \rightarrow \{(\mathcal{W}\mathbf{x})^0, (\mathcal{W}\mathbf{x})^1\} \doteq \{\mathbf{x}^0, \mathbf{x}^1\}$ с вдвое меньшим числом членов $K/2$ каждая, которое осуществляется следующим образом. Применим к

последовательности \mathbf{x} одно за другим два преобразования:

$$\mathbf{x} \rightarrow \tilde{\mathbf{y}} : \quad \tilde{y}(k) = \begin{cases} x(k), & k = 0, 2, \dots, K - 2, \\ x(k) - \lfloor (x(k-1) + x(k+1))/2 \rfloor, & k = 1, 3, \dots, K - 3, \\ x(k) - x(k-1), & k = K - 1, \end{cases}$$

и

$$\tilde{\mathbf{y}} \rightarrow \mathbf{y} : \quad y(k) = \begin{cases} x(k) + \lfloor (\tilde{y}(k+1) + 1)/2 \rfloor, & k = 0, \\ x(k) + \lfloor (\tilde{y}(k-1) + \tilde{y}(k+1) + 2)/4 \rfloor, & k = 2, 4, \dots, K - 2, \\ \tilde{y}(k), & k = 1, 3, \dots, K - 1, \end{cases}$$

где $\lfloor \cdot \rfloor$ – целая часть числа (floor). Теперь четные члены последовательности \mathbf{y} составляют последовательность \mathbf{x}^0 ($x^0(k) = y(2k)$, $k = 0, 1, \dots, K/2 - 1$), а нечетные – последовательность \mathbf{x}^1 ($x^1(k) = y(2k + 1)$, $k = 0, 1, \dots, K/2 - 1$). Описанное преобразование обратимо; вид обратного преобразования можно найти, например, в [5].

Последовательность (компонента) \mathbf{x}^0 представляет собой приближение исходной последовательности \mathbf{x} вдвое меньшего разрешения и называется приближением (масштаба 2), последовательность (компонента) \mathbf{x}^1 называется детальной составляющей (масштаба 2). Эффективность применения ДВП для сжатия данных обусловлена тем, что корреляция элементов детальной составляющей значительно меньше, чем корреляция элементов исходного сигнала, а значения детальной составляющей распределены значительно более неравномерно, чем исходные значения.

Одномерное ДВП трехмерного массива $\mathbf{X} = \{X(i, j, k)\}$ размерами I, J, K вдоль третьего измерения (Z) состоит в применении описанного преобразования \mathcal{W} к строкам $\mathbf{x}_{i,j} = \{X_{i,j}(k)\}$ этого массива ($X_{i,j}(k) = X(i, j, k)$). Результатом является разложение исходного массива на два массива размерами $I, J, K/2$:

$$\mathcal{W}_Z : \mathbf{X} \rightarrow \{(\mathcal{W}_Z \mathbf{X})^0, (\mathcal{W}_Z \mathbf{X})^1\} \doteq \{\mathbf{X}^0, \mathbf{X}^1\},$$

где

$$X^0(i, j, k) = (\mathcal{W}\mathbf{x}_{i,j})^0(k), \quad X^1(i, j, k) = (\mathcal{W}\mathbf{x}_{i,j})^1(k).$$

Массивы \mathbf{X}^0 и \mathbf{X}^1 – приближение и детальная составляющая массива \mathbf{X} .

Рассмотрим массив силовых кривых \mathbf{V} и используем для декорреляции его элементов следующее обратимое преобразование. Сначала применим к массиву одномерное ДВП и разложим его на компоненты $\mathcal{W}_Z : \mathbf{V} \rightarrow \{\mathbf{V}^0, \mathbf{V}^1\}$, после чего для приближения \mathbf{V}^0 осуществим переход к разности с предсказанием (9): $\mathcal{D}_1 : \mathbf{V}^0 \rightarrow \mathbf{D}_1^0$. Осуществим одномерное упорядочивание полученных массивов и оценим скорость арифметического кодирования и квазиэнтропию. Поскольку массивы имеют одинаковое число элементов, общая скорость кодирования и квазиэнтропия пары массивов могут быть вычислены как средние значения скоростей кодирования и квазиэнтропий отдельных массивов:

$$R(\{\mathbf{D}_1^0, \mathbf{V}^1\}) = \frac{1}{2} [R(\mathbf{D}_1^0) + R(\mathbf{V}^1)], \quad H(\{\mathbf{D}_1^0, \mathbf{V}^1\}) = \frac{1}{2} [H(\mathbf{D}_1^0) + H(\mathbf{V}^1)],$$

а слагаемые в квадратных скобках вычисляются по формулам (3) и (4). Данный вариант обработки обозначим как $\mathcal{P}(\mathcal{D}_1 \mathcal{W}_Z)$. Результаты, полученные для массивов I–V, представлены в соответствующей строке табл. 4.

Сравним полученные результаты с лучшими результатами предыдущего раздела, достигнутыми в варианте обработки $\mathcal{P}(\mathcal{D}_1 \mathcal{CS})$. По сравнению с вариантом $\mathcal{P}(\mathcal{D}_1 \mathcal{CS})$ заметное уменьшение как квазиэнтропии, так и скорости кодирования наблюдается для всех массивов. Среднее

Таблица 4. Скорость кодирования и квазиэнтропия компонент ДВП

Вариант обработки	Скорость кодирования					Квазиэнтропия				
	I	II	III	IV	V	I	II	III	IV	V
$\mathcal{P}(\mathcal{D}_1\mathcal{W}_z)$	3,8969	3,4781	3,4502	4,2646	4,2160	3,8781	3,4615	3,4319	4,2455	4,1970
$\mathcal{P}(\mathcal{D}_1\mathcal{W}_z\mathcal{C}\mathcal{S})$	3,8714	3,4392	3,4096	4,2217	4,1742	3,8701	3,4364	3,4069	4,2198	4,1724

(по массивам) уменьшение квазиэнтропии составляет приблизительно 0,14 бт/п, минимальное (массив V) – 0,05 бт/п, максимальное (массив III) – 0,26 бт/п. Уменьшение квазиэнтропии связано с применением ДВП и демонстрирует его эффективность. Среднее (по массивам) уменьшение скорости кодирования составляет приблизительно 0,12 бт/п, минимальное (массив V) – 0,03 бт/п, максимальное (массив III) – 0,25 бт/п. Избыточность кодирования в данном варианте обработки составляет порядка 0,02 бт/п для всех массивов. Это предсказуемо хуже, чем для варианта обработки $\mathcal{P}(\mathcal{D}_1\mathcal{C}\mathcal{S})$, и связано с тем, что в данном случае отсутствует преобразование сужения диапазона \mathcal{C} .

Попробуем уменьшить скорость кодирования за счет применения предварительной обработки, состоящей в разделении массива силовых кривых на массивы кривых подвода и отвода и сужении диапазона значений массива подвода. Выполним следующую последовательность преобразований. Разделим исходный массив \mathbf{V} преобразованием (6) на массивы подвода и отвода: $\mathcal{S} : \mathbf{V} \rightarrow \{\mathbf{V}^A, \mathbf{V}^R\}$. Сузим диапазон значений массива подвода преобразованием (7): $\mathcal{C} : \mathbf{V}^A \rightarrow \mathbf{C}^A$. Применим к полученным массивам $\mathbf{C}^A, \mathbf{V}^R$ одномерное ДВП:

$$\mathcal{W}_Z : \mathbf{C}^A \rightarrow \{\mathbf{C}^{A,0}, \mathbf{C}^{A,1}\}, \quad \mathcal{W}_Z : \mathbf{V}^R \rightarrow \{\mathbf{V}^{R,0}, \mathbf{V}^{R,1}\}.$$

Применим к нулевым компонентам (приближениям) полученных разложений преобразование перехода к разности с предсказанием (9): $\mathcal{D}_1 : \mathbf{C}^{A,0} \rightarrow \mathbf{D}_1^{A,0}, \mathcal{D}_1 : \mathbf{V}^{R,0} \rightarrow \mathbf{D}_1^{R,0}$. Осуществим одномерное упорядочивание четырех полученных массивов и оценим скорость арифметического кодирования и квазиэнтропию. Поскольку массивы имеют одинаковое число элементов, общая скорость кодирования и квазиэнтропия всех четырех массивов могут быть вычислены как средние значения скоростей кодирования и квазиэнтропий каждого из четырех массивов по отдельности:

$$R(\{\mathbf{D}_1^{A,0}, \mathbf{C}^{A,1}, \mathbf{D}_1^{R,0}, \mathbf{V}^{R,1}\}) = \frac{1}{4} [R(\mathbf{D}_1^{A,0}) + R(\mathbf{C}^{A,1}) + R(\mathbf{D}_1^{R,0}) + R(\mathbf{V}^{R,1})],$$

$$H(\{\mathbf{D}_1^{A,0}, \mathbf{C}^{A,1}, \mathbf{D}_1^{R,0}, \mathbf{V}^{R,1}\}) = \frac{1}{4} [H(\mathbf{D}_1^{A,0}) + H(\mathbf{C}^{A,1}) + H(\mathbf{D}_1^{R,0}) + H(\mathbf{V}^{R,1})],$$

а слагаемые в квадратных скобках вычисляются по формулам (3) и (4). Данный вариант обработки обозначим как $\mathcal{P}(\mathcal{D}_1\mathcal{W}_Z\mathcal{C}\mathcal{S})$. Результаты, полученные для массивов силовых кривых I–V, представлены в соответствующей строке табл. 4.

Сравним эти результаты с результатами предыдущего варианта обработки $\mathcal{P}(\mathcal{D}_1\mathcal{W}_Z)$. Для всех массивов наблюдается уменьшение квазиэнтропии и скорости кодирования. Уменьшение квазиэнтропии составляет приблизительно 0,008 бт/п для массива I и 0,025 бт/п для остальных массивов и связано с тем, что проведенная предварительная обработка (преобразования \mathcal{S} и \mathcal{C}) увеличивают эффективность последующих декорреляционных преобразований (\mathcal{W}_Z и \mathcal{D}_1). Уменьшение скорости кодирования составляет приблизительно 0,025 бт/п для массива I и 0,04 бт/п для остальных массивов и связано с уменьшением избыточности кодирования. В данном варианте обработки избыточность кодирования составляет порядка 0,001–0,003 бт/п в зависимости от массива, т.е. в 7–20 раз меньше, чем в предыдущем варианте обработки. Уменьшение избыточности объясняется применением преобразования \mathcal{C} сужения диапазона значений к массиву кривых подвода \mathbf{V}^A .

В ходе работы был исследован целый ряд вариантов обработки, основанных на применении других видов ДВП. Были рассмотрены, в частности, варианты применения двумерного ДВП к сечениям массива, перпендикулярным выбранному измерению, и применение трехмерного ДВП к массиву в целом. Также были рассмотрены всевозможные варианты применения ДВП, основанного на использовании системы вейвлетов (2–2) Хаара (см., например, [5]). Во всех случаях результаты оказались хуже результатов, полученных в варианте обработки $\mathcal{P}(\mathcal{D}_1\mathcal{W}_Z\mathcal{CS})$.

5. ЗАКЛЮЧЕНИЕ

В работе проведено исследование потенциальных возможностей алгоритмов обратимого сжатия применительно к массивам силовых кривых, получаемых в результате исследования мягких биологических объектов на атомно-силовом микроскопе. Получены оценки скорости кодирования при использовании стандартных алгоритмов сжатия (DEFLATE, JPEG 2000). Эти оценки (см. табл. 2) показывают, что правильное применение алгоритма кодирования без потерь из стандарта JPEG 2000 обеспечивает приблизительно четырехкратное сжатие массивов силовых кривых. Построен ряд алгоритмов сжатия на основе универсального арифметического кодирования, для которых также получены оценки скорости кодирования (см. табл. 3, 4). Все предложенные алгоритмы также обеспечивают приблизительно четырехкратное сжатие. Более детальное сравнение результатов показывает, что наиболее эффективный из предложенных алгоритмов (вариант обработки $\mathcal{P}(\mathcal{D}_1\mathcal{W}_Z\mathcal{CS})$) сжимает массивы силовых кривых на 4–8% лучше, чем это может быть сделано алгоритмом стандарта JPEG 2000.

Авторы благодарят А. Д. Протопопову за предоставленный экспериментальный материал и полезные обсуждения, касающиеся работы АСМ.

СПИСОК ЛИТЕРАТУРЫ

1. Binnig G., Quate C. F., Gerber Ch. Atomic Force Microscope. *Phys. Rev. Lett.*, 1986, vol. 56, no. 9, pp. 930–933.
2. Butt H.-J., Cappella B., Kappl M. Force measurements with the atomic force microscope: Technique, interpretation and applications. *Surface Science Reports*, 2005, vol. 59, pp. 1–152.
3. Ziv J., Lempel A. A Universal Algorithm for Sequential Data Compression. *IEEE Transactions On Information Theory*, 1977, vol. IT-23, no. 3, pp. 337–343.
4. Huffman D. A. A Method for the Construction of Minimum Redundancy Codes. *Proc. of IRE*, 1952, vol. 40, no. 9, pp. 1098–1101.
5. Taubman D. S., Marcellin M. W. *JPEG2000: Image Compression Fundamentals, Standards, and Practice*. New York: Springer Science + Business Media, 2002.
6. Witten I. H., Neal R. M., Cleary J. G. Arithmetic Coding for Data Compression. *Commun. of the ACM*, 1987, vol. 30, no. 6, pp. 520–540.
7. Штарьков Ю. М. *Универсальное кодирование. Теория и алгоритмы*. М.: ФИЗМАТЛИТ, 2013.
8. Штарьков Ю. М. Универсальное последовательное кодирование отдельных сообщений. *Проблемы передачи информации*, 1987, т. 23, № 3, стр. 3–17.
9. Le Gall D., Tabatabai A. Sub-band Coding of Digital Images Using Symmetric Short Kernel Filters and Arithmetic Coding Techniques. *IEEE Conference on Acoustics, Speech, and Signal Processing Proceedings*, 1988, pp. 761–765.
10. Sweldens W. The lifting scheme: A custom-design construction of biorthogonal wavelets. *Applied and Computational Harmonic Analysis*, 1996, vol. 3, no 2, pp. 186–200.

Статью представил к публикации член редколлегии Н. А. Кузнецов

On data compression of force volumes**Stefanovich A. I., Sushko D. V.**

We consider the problem of reversible (lossless) data compression of force volumes obtained as a result of studying biological samples using an atomic force microscope. We construct the bit rate estimates for standard algorithms (DEFLATE, JPEG 2000) used for the compression. We propose compression algorithms based on universal arithmetic coding, and construct the bit rate estimates for these algorithms.

KEYWORDS: reversible data compression, universal arithmetic coding, atomic force microscope, force volume.