

2. Dodds P. S., Rothman D. H. Scaling, universality and geomorphology // *Annu. Rev. Earth Planet. Sci.* - 2000. - Vol. 28. – p. 571 – 610.
3. Гришанин К.В. Устойчивость русел рек и каналов - Л: Гидрометеиздат, 1974. – 143 с.
4. Komarova N. L., Hulscher S. J. M. H. Linear instability mechanisms for sand wave formation // *Journal of Fluid Mech.* - 2000. - Vol. 413. – p. 219 – 246.
5. Петров А.Г., Потапов И.И. Постановка и решение задачи об устойчивости несвязного дна канала // *ПМТФ.* 2010. - Т. 51. № 1. – с. 62 – 74.
6. Петров А.Г., Потапов И.И. О развитии возмущений песчаного дна канала // *Доклады Академии наук.* 2010. - Т. 431. № 2. – с. 191 – 195.
7. Крат Ю.Г., Потапов И.И. Влияние нерегулярного возмущения потока на образование донных волн: препринт №175. - Хабаровск: Вычислительный центр ДВО РАН, 2012. – 22 с.
8. Guy N.P., Simons D.B., Richardson E.V. Summary of alluvial channel data from flume experiments. *Geol. Survey Profess. Paper 462-I.* Washington, 1967. – p.96.
9. Stephen E.Coleman, Juan J.Fedele, Marcelo H.Garcia. Closed-conduit bed-forms initiation and development // *Journal of Hydraulic Engineering.* - 2003. – Vol. 129, №12. – p. 956 – 965.
10. Jeremy G. Venditti, Michael Church. Morphodynamics of small-scale superimposed sand waves over migrating dune bed forms // *Water Resources Research* - 2005. - Vol.41. – p. 14.
11. Tjerry S. Morphological Calculations of Dunes in Alluvial Rivers // *Ph.D.-thesis.* / ISVA, Technical University of Denmark. 1995. pp. 193.
12. Sanne L.N. Modeling of sand dunes in steady and tidal flow // *Ph.D.-thesis.* / MEK-DTU, Technical University of Copenhagen, Denmark. 2003. pp. 185.

## **ИНФОРМАЦИОННЫЕ ТЕХНОЛОГИИ КЛАССИФИКАЦИИ И КОДИРОВАНИЯ СИМВОЛЬНЫХ ДАННЫХ**

Иванов В.Г., Ломоносов Ю.В., Любарский М.Г.

*Украина, Харьков, Национальный университет «Юридическая академия Украины им.  
Ярослава Мудрого»*

Предложен и исследован новый метод сжатия изображений текста на основе эффективного разделяющего правила связанных символов при их автоматической классификации. Обосновано применение метода обработки графического словаря символьных данных на основе формирования и классификации вертикальных элементов строки. Предложенный новый 2-х этапный алгоритм сжатия изображения текста имеет преимущество около 25% в степени компрессии в сравнении с известным алгоритмом JB2 формата DjVu.

### **Information technologies of classification and encoding of character data. Ivanov V.G., Lomonosov U.V., Lyubarsky M.G.**

The new method of compression of images of the text on the basis of an effective dividing rule of coherent symbols is offered and investigated at their automatic classification. Application of a method of processing of the graphic dictionary of symbolical data on the basis of formation and classification of vertical elements of a line is proved. Offered new 2th a stage algorithm the algorithm of compression of the image of the text has advantage about 25 % in compression degree in comparison with known algorithm JB2 of format DjVu.

Теория и практика сжатия данных является основным и эффективным инструментом формирования и архивации цифрового контента огромного числа различных текстовых документов: книг, журналов, технической документации и т.д. При этом становится очевидным, что наибольшая степень сжатия этого типа изображений может быть достигнута только с использованием методов классификации и распознавания образов [1 – 7]. В частности, на классификации изображений символов, составляющих изображение текста, основан наиболее эффективный алгоритм JB2, который как составная часть входит в графический формат DjVu [8, 9].

В идеальном случае после проведения классификации при таком подходе, все изображения одного и того же символа в тексте должны попадать в один класс. В таком случае, можно заменить все изображения из каждого класса одним представителем, например, усредненным изображением данного символа. Изображение текста, после классификации символов, можно представить в виде «графического словаря» – набора усредненных изображений каждого символа, и «картой регионов» – описанием положения каждого символа в тексте. В итоге результирующий размер файла изображения текста существенно уменьшится [10]. В реальности не удается получить для каждого символа только один класс его изображений и соответственно иметь для него в «графическом словаре» только одного представителя. Чем выше уровень шумов (при одном и том же качестве классификации), тем больше вариантов изображения одного и того же символа (классов) находится в словаре и тем больший объем имеет «графический словарь». Это существенно снижает общую степень сжатия всего изображения.

Целью настоящей работы является построение двухэтапного алгоритма сжатия текстовых изображений: выделение связанных символов и их автоматическая классификация – первый этап; дополнительное сжатие «графического словаря» после разбиения его представителей на вертикальные элементы строки и их последующей классификации – второй этап.

Предлагаемый алгоритм сжатия текстовых изображений схематично представлен на рис. 1.

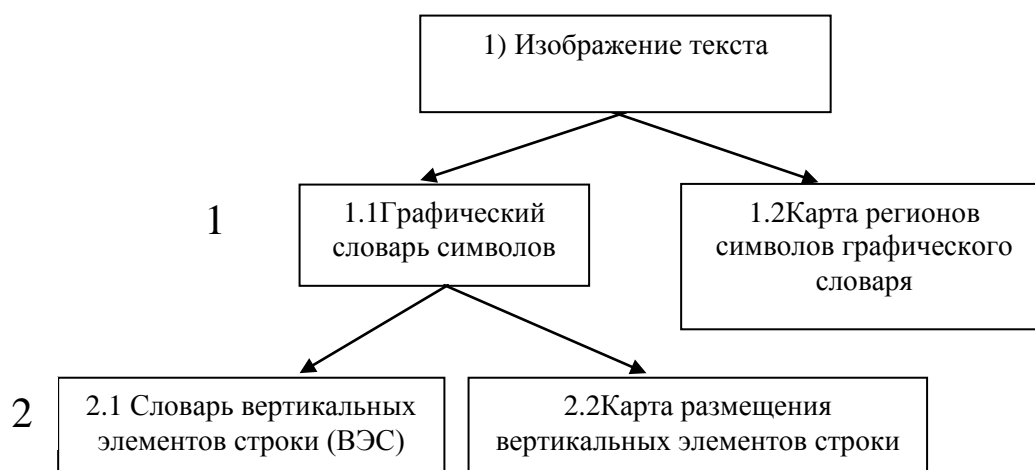


Рис.1. Схема двух этапной обработки изображения текста.

На первом этапе, при создании «графического словаря» и «карты регионов» используется технология выделения связанных символов и разбиение их на классы, предложенная в работе [10]. Основным отличием применяемой здесь классификации является использование новой методики определения степени близости изображений

двух символов при их сравнении. Предложенная методика мало чувствительна к шумам печати и сканирования, так как основана на стабильных характеристиках, которые не зависят от характерных контурных шумов печати и сканирования. Это в значительной степени повышает качество классификации изображений символов, осуществляемой с помощью известного алгоритма «просеивания» [11].

При сравнении двух изображений символов  $S_1$  и  $S_2$  с допустимыми отклонениями  $\Delta H$  (по высоте символа),  $\Delta W$  (по ширине символа) и  $\Delta P$  (по периметру символа) эти изображения накладываются друг на друга с помощью плоско-параллельного переноса так, чтобы их центры тяжести совпадали. Далее подсчитываются две величины:  $R(S_1, S_2)$  – количество «существенных отличий», и  $D(S_1, S_2)$  – количество общих черных точек.

Первая величина – это количество несовпадающих по яркости (белое – черное) точек, которые не являются смежными для совокупности общих черных точек. Таким образом, количество существенных отличий  $R(S_1, S_2)$  игнорирует несовпадения в тех точках, которые лежат на периметрах изображений и, как правило, представляют собою шумы печати и сканирования. Например, при сравнении изображений букв «с» и «е» точки существенных отличий составляют горизонтальный отрезок, который имеется в букве «е» и отсутствует в букве «с».

Вторая величина нужна для обезразмеривания первой, чтобы диапазон возможных значений величины

$$\varepsilon(S_1, S_2) = \frac{R(S_1, S_2)}{D(S_1, S_2)} 100\% \quad (1)$$

для всех пар символов не менялся при изменении размера шрифта и разрешения сканирования.

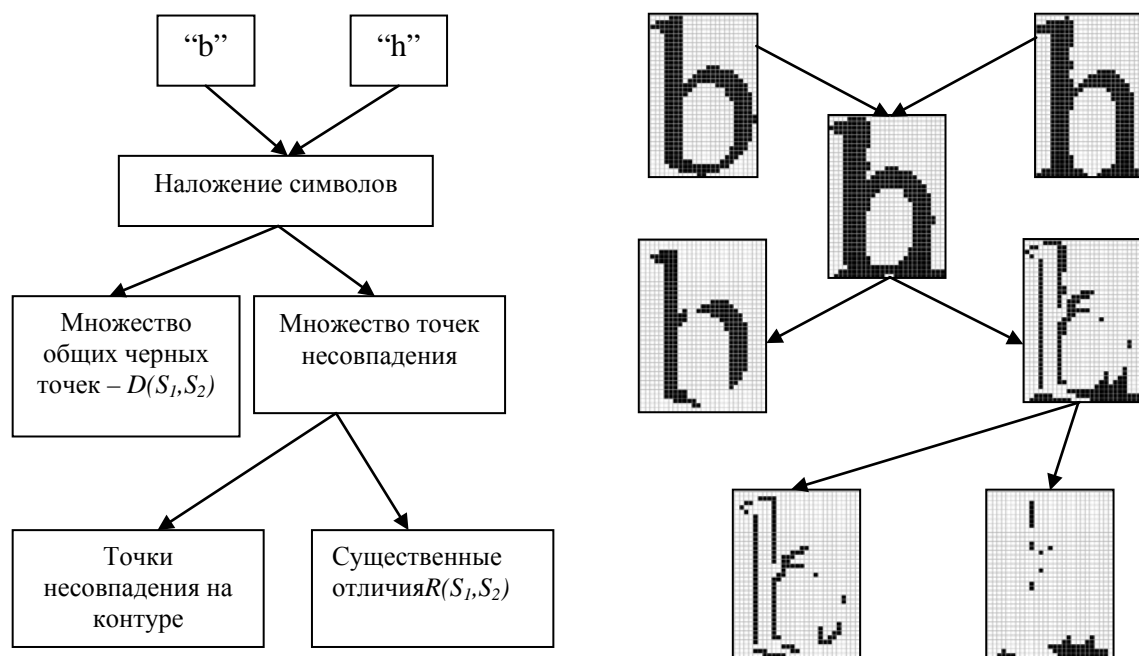


Рис.2. Схема сравнения изображений символов «b» и «h».

На рис. 2 показаны, какие совокупности точек рассматриваются при сравнении двух трудно различимых букв «b»и«h». Справа показаны их геометрические расположения, а слева – поясняющая схема.

Величина  $\varepsilon$ , выражение (1), является мерой близости изображений двух символов, используемой в алгоритме «просеивания».

Рис. 2 показывает, что даже среди точек существенного отличия имеются одиночные точки и малые группы точек, которые носят случайный характер, вызванный контурными шумами печати и сканирования. Действительные отличия в начертании букв «b»и«h» демонстрируют большие группы точек. Поэтому, вместо определенной выше функции  $R(S_1, S_2)$ , подсчитывающей количество «существенных отличий», в настоящей работе используется ее модификация, которая подсчитывает это число с учетом веса [3]. Весовой коэффициент каждой точки в  $R(S_1, S_2)$  тем больше, чем больше у данной точки таких же смежных точек. Эта функция, будем по-прежнему обозначать ее через  $R(S_1, S_2)$ , дает классификацию с меньшим числом классов, чем первоначальный ее вариант (без учета весов).

Таким образом, предлагаемая метрика  $\varepsilon$ , определяющая степень близости изображений двух символов при классификации алгоритмом «просеивания», мало чувствительна к шумам печати и сканирования. Она основана на стабильных характеристиках  $R(S_1, S_2)$  и  $D(S_1, S_2)$ , которые подавляют (не учитывают) контурные шумы сравниваемых символов при их наложении.

Завершает 1-й этап обработки изображения текста (рис.1) нахождение усредненных изображений в каждом классе. После этого все изображения символов можно отбросить – остается только совокупность «представителей» каждого класса – усредненные изображения. Процедура усреднения состоит в наложении друг на друга всех изображений класса при совмещенных «центрах тяжести» каждого изображения, вычисления среднего значения яркости для каждой точки и округления. На рис. 3 иллюстрируется процесс усреднения для трех классов изображений одного и того же символа «n». Черные точки означают, что среднее значение яркости в них равно 0, серые – что среднее значение меньше или равно  $\frac{1}{2}$ , а светлые – больше  $\frac{1}{2}$ . При округлении черные и серые точки превращаются в черные, а светлые – в белые. Результат показан на рис. 3.

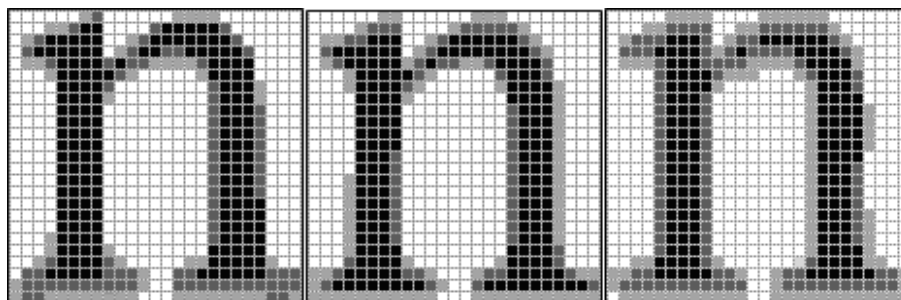


Рис. 3. Изображения различных классов символа “n”.

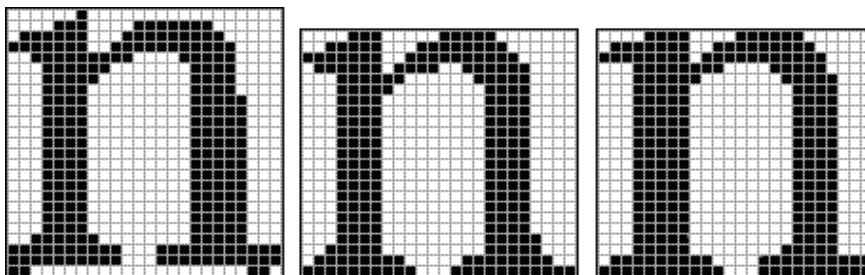


Рис.3. Изображения средних представителей классов изображений символа «n».

В полученном на предыдущем этапе семействе усредненных изображений по-прежнему присутствуют разные изображения одних и тех же символов (рис. 4). Ситуацию можно улучшить, применив еще одну классификацию, призванную отождествить изображения одного и того же символа.

Для этой классификации, как и для предыдущей, используется нечувствительное к контурным шумам классифицирующее правило  $\varepsilon$  (1), но алгоритм классификации выбирается другим. Дело в том, что теперь изображения одного символа очень близки друг к другу – случайная компонента изображения в значительной мере подавлена. Поэтому существенно меньше опасность спутать изображения двух разных символов. Это позволяет применить алгоритм «наращивания областей».

Алгоритм «наращивания областей» состоит в том, что на первом шаге, начиная с произвольно выбранного элемента классифицируемого множества, к его классу присоединяются все достаточно близкие элементы. На втором шаге к вновь присоединенным элементам добавляются все элементы, близкие к ним. Процесс «наращивания» повторяется до тех пор, пока на каком-то шаге не окажется новых элементов, которые можно было бы присоединить. Далее все элементы «выращенного» класса исключаются из классифицируемого множества и «выращивается» следующий класс. Алгоритм заканчивает работу, когда в классифицируемом множестве не остается ни одного элемента.

После получения новых классов, изображения в каждом из них усредняются [12], и получившийся набор изображений представляет собой «графический словарь». На рис. 4 показано теперь единственное изображение символа «n», вошедшее в графический словарь.

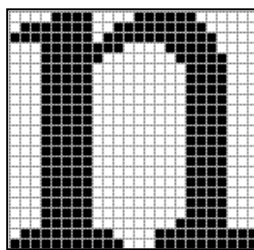


Рис. 4. Представитель класса изображения символа «n», после повторной классификации, который заносится в «графический словарь».

В таблице 1 для различных разрешений приводятся число классов после основной и повторной классификаций. Для сравнения приведено число классов после классификации алгоритмом JB2 (формат DjVu в режиме Bitonal).

Таблица 1.

Разрешение изображения текста (dpi)	Количество классов в исходном	Количество классов после основной	Количество классов после второй	Количество классов после классификации

	изображении	классификации $\varepsilon_{opt} = 6\%$	классификации $\varepsilon_{opt} = 6\%$	алгоритмом JB2
600 dpi	3558	197	72	314
500 dpi	3557	137	72	259
400 dpi	3557	130	71	199
300 dpi	3545	122	95	235
200 dpi	3890	237	148	451

Данные, приведенные в таблице, демонстрируют достаточно высокую эффективность как первой, так и повторной классификаций и несомненные преимущества предлагаемого алгоритма перед алгоритмом JB2 .

Второй этап – сжатие «графического словаря» изображений символов, полученного на первом этапе (рис. 1).

Все изображения символов «графического словаря» требуется разложить на вертикальные элементы строки (ВЭС), как показано на рис. 6. Как видно из рис. 6, все ВЭС формирующие изображение символа “е” имеют одинаковый размер по высоте, что является необходимым условием для их автоматической классификации.

При помощи горизонтальной гистограммы строки определяется высота строки как расстояние между линией верхних выносных элементов (ЛВВЭ) и линией нижних выносных элементов (ЛНВЭ). Так же определяются линия строчных знаков (ЛСЗ) и базовая линия (БЛ), рис. 5.

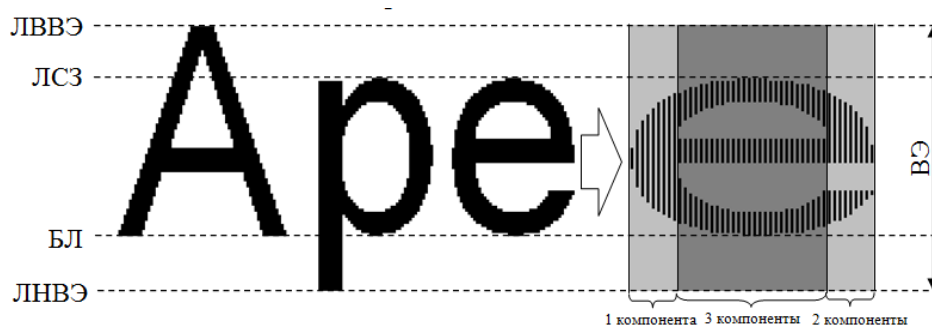


Рис. 5. Разложение изображения символа “е” на вертикальные элементы строки по количеству связанных компонент.

В данном случае, ВЭС составляющие изображение символа “е” будут распределены в группы ВЭС с одной, двумя и тремя связными компонентами (рис. 6). Разделение ВЭС на группы по количеству компонент необходимо для того, чтобы при их классификации исключить возможность сравнения ВЭС с различным количеством компонент и дальнейшим их возможным объединением в один класс.

Основным моментом алгоритма «просеивания» является выбор метода, который определяет понятие достаточной близости двух элементов (ВЭС). Этот метод определяется тем важным обстоятельством, что шумы печати и сканирования носят контурный, а не структурный характер, то есть искажения символов происходят только на их границе.

Возможные варианты несоответствия связных компонент сравниваемых ВЭС представлены на рис. 6. Допуск несовпадения связных компонент ( $\Delta = \pm 1$  точка), обусловлен тем фактом, что искажения, вносимые контурными шумами сканирования, в каждой связной компоненте могут присутствовать только на её концах - сверху и снизу.

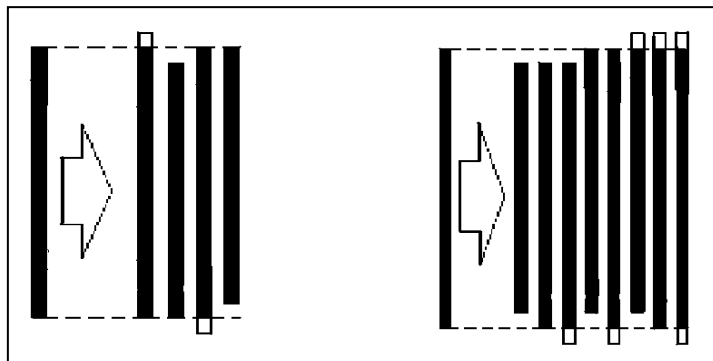


Рис. 6. Возможные варианты близких связанных компонент ВЭС.

В качестве центров первого приближения для алгоритма автоматической классификации используются векторы, являющиеся средними в полученных методом «просеивания» классах, то есть

$$e_j = \frac{1}{M_j} \sum_{x \in X_j} x, \quad (2)$$

где  $M_j$  – число элементов в классе  $X_j$ ,  $j=1,2,\dots,m$ , и  $m=m(\Delta)$  – число получившихся классов.

Предварительное применение метода просеивания значительно улучшает сходимость алгоритма  $k$ -средних и, следовательно, сокращает вычислительное время. Это объясняется тем, что уже на первом шаге центры нулевого приближения аппроксимируют члены своего класса с точностью, не меньшей значения параметра  $\Delta$ . Еще одним преимуществом предложенной предобработки является то, что метод просеивания позволяет автоматически определить необходимое число классов.

Для кодирования “графического словаря” символов необходимо составить карту размещения ВЭС, которая определяет размещение векторов  $S_1, \dots, S_k$ , и для каждого вектора  $S_j$  указать его представителя, в качестве которого используется его центр  $e_j$ .

Далее карта размещения ВЭС и сам словарь ВЭС кодируется стандартными методами сжатия без потерь в бинарном виде.

Фрагмент полученного словаря ВЭС для изображения текста с разрешением 300dpi представлен на рис. 8 (серые полосы – искусственно введенные разделители для отдельных классов ВЭС).

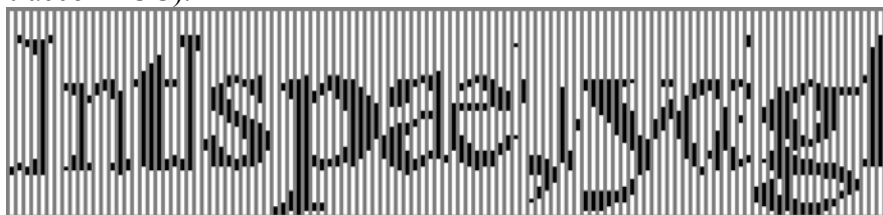


Рис. 7. Фрагмент словаря ВЭС после автоматической классификации.

По сравнению с алгоритмом JB2 (DjVu), использование I-го и II-го этапа обработки в предложенном алгоритме позволяет уменьшить объемы выходных данных на 29% при 200 dpi, 25% при 300 dpi, 24% при 400 dpi, 24,5% при 500dpi и 25,5% при 600 dpi соответственно.

В таблице 2 приведены значения коэффициентов сжатия всего изображения текста после I-го, I-го и II-го этапов обработки и степень сжатия алгоритма JB2 формата DjVu, а также выигрыш в процентном соотношении.

Таблица 2

Разрешение изображения текста (dpi)	200	300	400	500	600
Коэффициент сжатия изображения текста (I этап)	62,38	135	250,48	353,54	436,7
Коэффициент сжатия изображения текста (I и II этап)	74,3	166,18	267,18	361,76	445,34
Преимущество I-го и II-го этапов в сравнении с I-м этапом в (%)	<b>16%</b>	<b>19%</b>	<b>6%</b>	<b>3%</b>	<b>2%</b>
Коэффициент сжатия изображения текста JB2 (DjVu)	52,63	124,16	202,41	272,91	330,73
Преимущества предложенного метода (I и II этап) в сравнении с JB2 в (%)	<b>29%</b>	<b>25%</b>	<b>24%</b>	<b>24,5%</b>	<b>25,5%</b>

На рис. 9 приведены графики отображающие значения коэффициентов сжатия данных при I-м этапе обработки, последовательном применении I-го и II-го этапов обработки и при сжатии изображения текста алгоритмом JB2 формата DjVu. Преимущество предложенного метода составляет около 25% для всех разрешений изображения (табл.2).

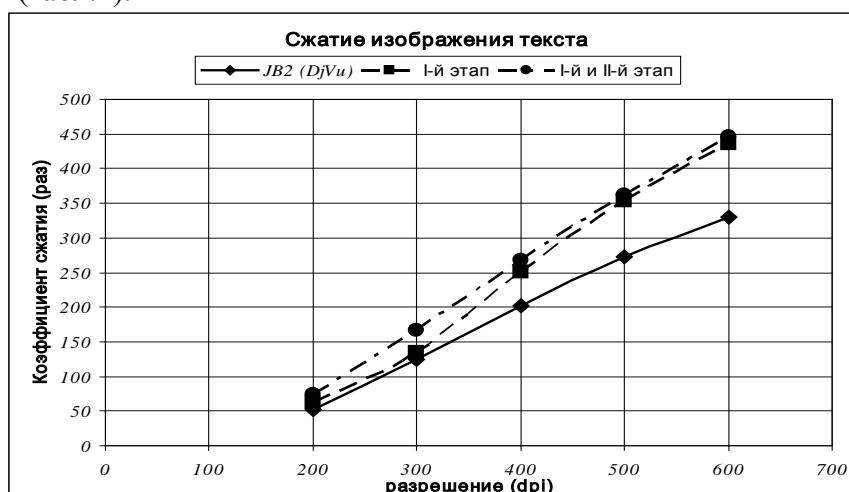


Рис. 8. Сравнительная характеристика результатов сжатия изображения текста для различных разрешений.

Качество классификации первого этапа (рис. 1) предлагаемого метода значительно выше, чем у алгоритма JB2. Количество классов, получающееся в результате предложенной классификации, более чем в два раза меньше при всех разрешениях сканирования. Это является основной качественной характеристикой предложенного метода и дает широкие возможности повышения информативности этого алгоритма в инженерных реализациях.

Используя дополнительный этап обработки изображений текста, который состоит в разбиении “графического словаря” символов на совокупность ВЭС (рис. 1) и их классификации удалось сформировать новый словарь ВЭС и карту их размещения, которые заменяют исходный “графический словарь”, имея существенно меньший объем в совокупности.

Предложенная двухэтапная обработка изображения текста позволила увеличить степень компрессии при различных разрешениях исходного изображения: на 16% для 200dpi; на 19% для 300dpi; на 6% для 400dpi; на 3% для 500dpi; на 2% при 600dpi, по сравнению с I-м этапом обработки. Качество “графического словаря” символов при визуальном оценивании, в результате второго этапа обработки, не изменилось, что в свою очередь не влияет на восприятие изображения всего текста [13].



Сравнение с лучшим в настоящее время специальным алгоритмом для сжатия изображений текста – JB2, входящим в формат DjVu, показало, что предлагаемый 2-х этапный алгоритм сжатия изображения текста имеет преимущество в степени сжатия данных в среднем на 25% и на 9 % лучше результатов работы [10] при наиболее часто используемых на практике разрешениях изображения текста (200-600dpi).

Полученный результат является знаковой характеристикой метода и открывает широкие возможности повышения информативности представления данных в графических форматах.

### Литература.

- 1) *Д.Сэломон*. Сжатие данных, изображений и звука. – Москва: Техносфера, 2004. – 368с.
- 2) Компьютерное зрение. Современный подход / *Д. Форсайт, Д. Понс*; Пер. с англ. – М.: Вильямс, 2004. – 928 с.
- 3) Цифровое кодирование графики. Тематический выпуск. ТИИЭР. – М.: Мир, 1980. – Т. 68, № 7. – 214с.
- 4) *Gupta Maya R., Stroilov A.* Segmenting for wavelet compression // Data Compression Conference (DCC). – USA, Utah, Snowbird, 2005. – Режим доступа: <http://www.cs.brandeis.edu/>.
- 5) *Гонсалес Р., Вудс Р., Эддингс С.* Цифровая обработка изображений в среде MATLAB. – Москва: Техносфера, 2006.- 616с.
- 6) *Иванов В.Г., Любарский М.Г., Ломоносов Ю.В.* Сокращение содержательной избыточности изображений на основе классификации объектов и фона // Проблемы управления и информатики. – Киев, 2007. – № 3. – С. 93-102.
- 7) *Иванов В.Г., Ломоносов Ю.В., Любарский М.Г.* Сжатие изображений на основе автоматической и нечеткой классификации фрагментов // Проблемы управления и информатики. – К. – 2009. – №1 – С.52-63.
- 8) Technical Papers from AT&T Labs: Электронный ресурс.– Режим доступа: <http://djvuzone.org/techpapers/index.html>.
- 9) <http://www.djvu.org/> - портал DjVu-сообщества
- 10) *Иванов В.Г., Любарский М.Г., Ломоносов Ю.В.* Сжатие изображения текста на основе выделения символов и их классификации // Проблемы управления и информатики. – Киев, 2010. – № 6. – С. 111-122.
- 11) Прикладная статистика: Классификация и снижение размерности: [Справочник] / *С.А.Айвазян, В.М. Бухштабер, И.С. Енюков и др.*; Под ред. С.А. Айвазяна. – М.: Финансы и статистика, 1989. – 607с.
- 12) Обработка изображений на ЭВМ/ *Е.А. Бутаков, В.И. Островский, И.Л. Фадеев*.- М.: Радио и связь, 1987.-240с.:ил.
- 13) *Иванов В.Г., Ломоносов Ю.В., Любарский М.Г.* Сжатие изображения текста на основе формирования и классификации вертикальных элементов строки в графическом словаре символьных данных // Проблемы управления и информатики. К. – 2011. – №5 – С. 98-109.

## КЛАССИФИКАЦИЯ АМОРТИЗАТОРОВ РАДИОЭЛЕКТРОННЫХ СРЕДСТВ НА ОСНОВЕ ФАСЕТНОЙ СТРУКТУРЫ

Лысенко А.В.

Пенза, ФГБОУ ВПО «Пензенский государственный университет»