

Cognitive PDF/A — технология оцифровки текстовых документов для публикации в Интернете и долговременного архивного хранения *

С. А. Усилин¹, Д. П. Николаев², В. В. Постников³

¹ ООО «Когнитивные технологии»,

Россия, 117312 Москва, пр. 60-летия Октября, 9

² Институт проблем передачи информации им. А. А. Харкевича РАН,

Россия, 127994 Москва, Большой Каретный пер., 19, стр. 1

³ Институт системного анализа Российской академии наук,

Россия, 117312 Москва, пр. 60-летия Октября, 9

В работе рассматриваются вопросы перевода бумажных документов в электронных вид. Предлагается оригинальная технология оцифровки, включающая сжатие, распознавание и упаковку текстовых документов способом, пригодным для долговременного архивного хранения. Технология использует отделение цветных элементов с помощью гистограммы насыщенности, выделение текстовых блоков с опорой на преобразование Хафа и морфологическую фильтрацию, а также методы оптического распознавания.

Введение

Увеличение объемов информации во всем мире и бурное развитие информационных технологий привело к повсеместному использованию электронных документов. Электронные версии обладают целым рядом преимуществ. Во-первых, решается проблема оперативного доступа — поиск необходимой информации, копирование, распечатка занимают секунды. Во-вторых, цифровые документы не разрушаются со временем, не ухудшают своих пользовательских качеств и могут храниться практи-

* Работа выполнена при финансовой поддержке РФФИ (грант №09–07–00473-а).

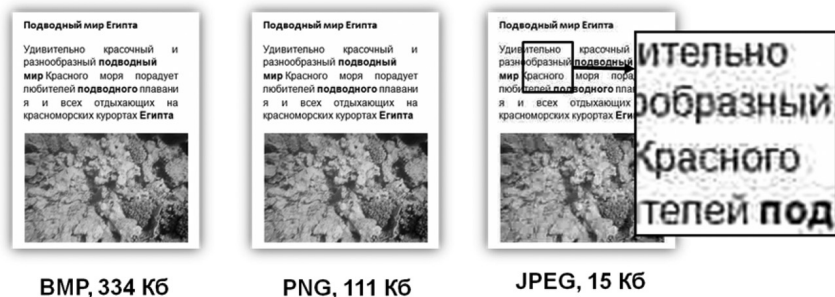


Рис. 1. Применение классических способов сжатия к изображению документа

чески вечно. В-третьих, решается проблема целостности: при правильной организации электронного архива несанкционированное удаление, фальсификация или неумышленная модификация информации невозможны. Помимо перечисленных достоинств, электронные архивы обладают рядом неявных, но удобных с практической точки зрения преимуществ: многопользовательский режим и дистанционный доступ, что абсолютно недоступно при использовании бумажного архива.

Однако большая часть знаний на сегодняшний день до сих пор хранится в бумажном виде. Наиболее популярным способом перевода бумажных документов в электронный вид является сканирование [1, 2]. В результате создаются точные электронные копии (электронные образы) документов со всеми графическими особенностями (рисунками, пометками, подписями, печатями и т. д.). При качественном сканировании получившиеся изображения-образы зачастую оказываются достаточно большого размера. Например, документ формата А4, отсканированный в цветном режиме при разрешении 300 DPI, имеет размер порядка 25 Мб. Использование файлов таких больших размеров неэффективно в электронных архивах, поэтому все больший интерес обретают технологии сжатия получившихся электронных образов. Классические технологии сжатия изображений [3, 4] не применимы, так как в общем случае документы могут содержать как монохромный текст, так и полноцветные графические области. Алгоритмы сжатия изображений без потерь, результативные для монохромных текстов, неэффективны для полноцветной графики, в то время как сжатие с потерями демонстрирует высокие показатели для цветных изображений, однако сильно искажает текстовую информацию (рис. 1). Поэтому обычно для сжатия изображений такого типа используют комбинированный подход.

1. Структурное сжатие изображений документов

Изложим идею структурного сжатия на примере изображения страницы журнала (рис. 2). Классическая страница журнала может содержать фоновый рисунок, один или несколько текстовых блоков, графических элементов (фотографии, схемы, таблицы и пр.) и каких-то пометок. Основная идея структурного сжатия изображений такого рода заключается в выделении структурных блоков, объединении данных блоков в слои (т. е. «расслоение» изображения на текстовые, графические и прочие слои) и сжатие каждого слоя наиболее подходящим образом. Так изображение страницы журнала на рис. 2 расслаивается на четыре слоя: фон, область черного текста, область синего текста и область с фотографией. Для сохранения максимального качества текстовые слои следует сжимать алгоритмами сжатия без потерь (например, CCITT Group 4), в то время как для фотографии вполне допустимо применение методов сжатия с потерями (JPEG).

Основное место в алгоритмах структурного сжатия отводится методам расслоения исходного изображения на текстовый и графический слои. На сегодняшний день опубликовано большое количество статей, предлагающих способы отделения текстовых блоков от графики различными методами [5–9]. Несмотря на это, высокую популярность данный подход получил сравнительно недавно. Одним из примеров, полностью реализующим идею структурного сжатия, по праву можно считать формат DjVu [10].

Для сжатия цветных изображений в DjVu применяется специальная технология, разделяющая исходное изображение на три слоя: передний план, задний план и черно-белую (однобитовую) маску. Маска сохраняется с разрешением исходного файла; именно она содержит изображение текста и прочие четкие детали. Разрешение заднего плана, в котором остаются иллюстрации и текстура страницы, понижается для экономии места. Передний план содержит цветовую информацию о деталях, не попавших в задний план; его разрешение понижается еще сильнее. Затем задний

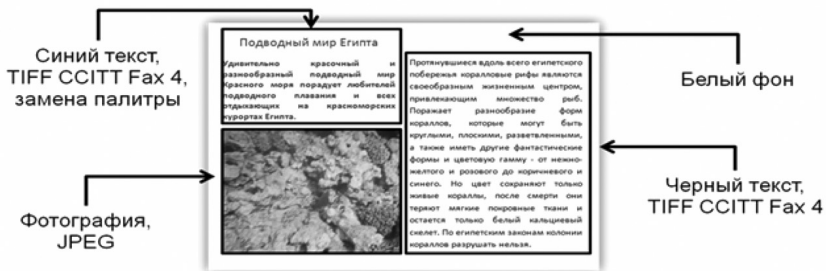


Рис. 2. Структурные блоки изображения документа

и передний планы сжимаются с помощью вейвлет-преобразования (алгоритмом IW44), а маска — алгоритмом JB2.

Несмотря на высокие коэффициенты сжатия изображений документов, DjVu обладает существенным недостатком: на сегодняшний день формат не стандартизован, что затрудняет его использование в качестве средства для создания электронных архивов. К тому же использование одинаковой схемы расслоения для всех типов документов не всегда оправдано, и даже иногда может приводить к значимому искажению документа. Дополнительно стоит отметить, что в формате полностью отсутствуют какие-либо средства обеспечения безопасности и конфиденциальности документов.

2. Технология Cognitive PDF/A

Опишем технологию Cognitive PDF/A, предназначенную для перевода бумажных документов в электронный вид, и процесс оцифровки в соответствии с предлагаемой технологией (рис. 3).

Первым этапом обработки является расслоение исходного изображения. В результате появляются два новых изображения. Первое содержит области исходного изображения, соответствующие текстовой информации (текстовый слой), а второе — графическим элементам (графический слой).

В соответствии с архитектурой алгоритма текстовый слой не должен содержать никаких лишних областей, кроме текстовых блоков. Следовательно, изображение текстового слоя может быть легко распознано без какой-либо предварительной подготовки с помощью внешних OCR-систем.

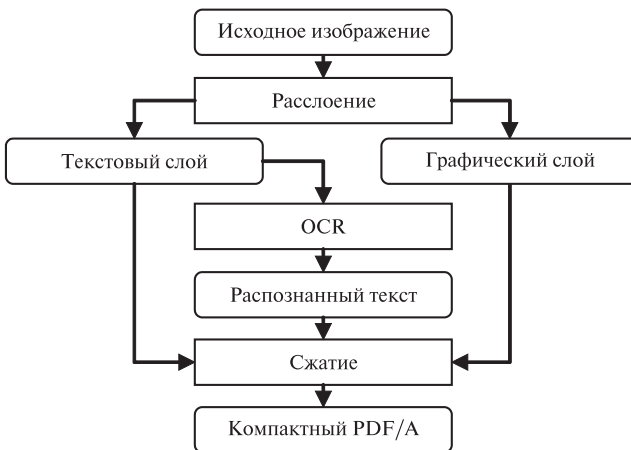


Рис. 3. Схема технологии Cognitive PDF/A

Последним действием является упаковка полученных слоев и распознанного текста в PDF/A. Графический и текстовый слои подвергаются соответствующему сжатию, а распознанный текст упаковывается таким способом, чтобы обеспечить максимальное удобство поиска и копирования информации в документе.

Таким образом, технология Cognitive PDF/A состоит из трех основных частей: расслоение исходного изображения, распознавание текстового слоя с помощью OCR-системы и компактная упаковка получившихся слоев и распознанного текста в PDF/A-файл. Рассмотрим эти части более подробно.

3. Алгоритм расслоения

Разные типы документов обладают различными особенностями. Например, для финансовых документов характерно наличие печатей, подписей и штампов, журнальные статьи могут иметь сложный многоцветный фон, в книги часто включают полноцветные графические элементы. Поэтому технологией Cognitive PDF/A предусматриваются уникальные схемы расслоения для каждого типа документа. Выбор наилучшей схемы осуществляется с помощью алгоритмов предварительной идентификации типа документа [11–13]. Далее в качестве примеров будут рассмотрены схемы расслоения для двух важных типов документов: страница книги и офисный документ.



Рис. 4. Изображение страницы книги



а)



б)

Рис. 5. Изображение страницы книги после бинаризации (а); после применения морфологической операции «размыкание» (б)

Обычно страница книги содержит черный текст на белом фоне и, возможно, графические элементы: рисунки, схемы, графики и пр. (рис. 4). Обычно в книгах области текста и графики не пересекаются. Еще одной ключевой особенностью верстки книг является использование шрифтов близких линейных размеров. Опираясь на эти характерные черты, построим схему расслоения изображения страницы книги.

Шаг 1. Бинаризуем исходное изображение [14, 15], тем самым преобразуем его в монохромный вид (рис. 5 а). Так как изображение в основном содержало черный текст на белом фоне, то процесс бинаризации не должен сильно сказаться на областях, содержащих текстовую информацию.

Шаг 2. С помощью морфологической фильтрации [4, 16] «сольем» слова в единые компоненты связности. Обозначим w и h характерные ширину и высоту символов соответственно. Заметим также, что расстояние между буквами в слове сравнимо с толщиной штриха символа, а расстояние между словами близко ширине символа. Поэтому «склеим» каждое слово

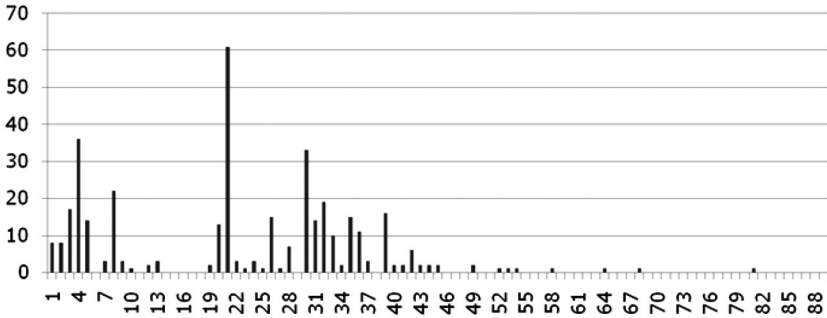
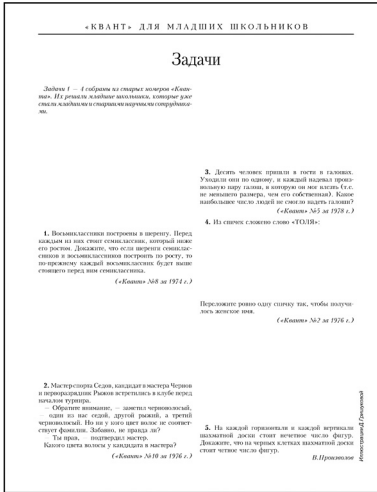
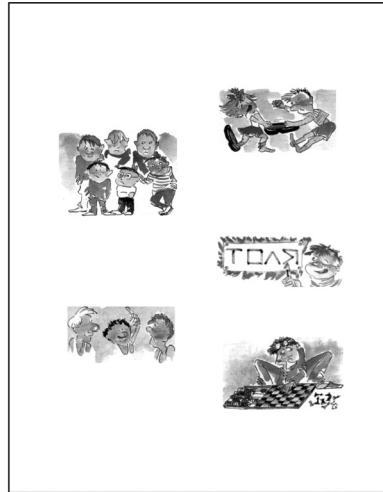


Рис. 6. Гистограмма высот компонент связности



а)



б)

Рис. 7. Текстовый (а) и графический (б) слои изображения страницы книги

в отдельную компоненту связности, выполнив размыкание с окном $w \times 1$ (рис. 5 б). Прямое вычисление морфологического фильтра имеет линейную сложность от размера окна. Поэтому для ускорения используется алгоритм Ван Херка [17], вычисляющий морфологические операции с прямоугольным примитивом за время, не зависящее от размеров окна.

Шаг 3. Построим гистограмму высот полученных компонент связности (рис. 6). Так как весь текст на странице напечатан примерно одинаковым по размеру шрифтом, то компоненты связности, соответствующие

S-SHINA.RU Расходная накладная № 1875 от 28 августа 2008 09:37:00
 ООО «СИАМ ДИСКОВЫЕ ДИСКИ» Действительна только в случае оплаты в кассе продавца в
 тел. (495) 730-67-34, 223-22-67 день выписки до конца рабочего дня.

Поставщик: Общество с ограниченной ответственностью "Випс Авто"
 Покупатель: Частное лицо (Клиент-"Варшава")
 № машины: Овангеловский мандатер: Савваидисов Руслан Габдувалч Подразделение: Варшавское шоссе, д. 46

№	Артикул	Товар	Количество	Цена	Сумма без скидки	Скидка	Сумма	Номер ГТД	Страна происхождения
1	A3337x	АДДС Бx114.3 Б.Бx15 2738 шт.1.ИИФ Алашу	4 шт	2 332,00	9 328,00	652,90	8 675,04		
2	5y708113	АДДСы 15005 R15 2017 0500уры ШИШТР.Екста	4 шт	3 317,00	13 268,00	928,76	12 339,24		
Итого:				22 596,00	1 581,72	21 014,28			
В том числе НДС:						3 205,57			

Всего наименований 2, на сумму 21 014,28 руб.
Двадцать одна тысяча четырнадцать рублей 26 копеек
 Удостоверяется подлинностью подписями продавца и покупателя, проставлением даты продажи, печати продавца, подписи продавца и покупателя.
 ОБЯЗАТЕЛЬСТВО ПОСМОТРЕТЬ ДИСК И ЗАЩИТУ ПЕРЕД ПОЛУЧЕНИЕМ ОБЪЕКТА ПРОДАЖИ ОБЕСПЕЧИВАЮТСЯ ПОДПИСАТЕЛЯМИ ДОКУМЕНТА, НЕ ИМЕЮЩИМИ
 ОБЯЗАТЕЛЬСТВА ПРОСМОТРЕТЬ ДИСК И ЗАЩИТУ В ОБЪЕМЕ ОБЪЕКТА ПРОДАЖИ В ТЕРМИН 24 ЧАСА С МОМЕНТА ПОЯВЛЕНИЯ ЧАСТНОГО ЛИЦА НА ТЕРРИТОРИИ ПРОДАВЦА И СРОКА УДАЛЕНИЯ ДОКУМЕНТА ОТДАВАЮЩИМ ПРАВОМ ДИСКА. БУДУЩЕЕ ОБЪЕДИНЕНИЕ ПОКУПАТЕЛЯ ЗАКАЗ ПОДЛИСКИ НА ПОСЛЕДНИЙ ДИСКОВЫЙ НАСТРОЙЩИК, ЧТО НЕ ЯВЛЯЕТСЯ ПРЕДСТАВИТЕЛЕМ ОБЪЕКТА ПРОДАЖИ
 ПОДЛИСКИ И ОТКАЗЫВАЮЩИМ ВНЕ ТЕРМИНА, ИСПОЛЬЗОВАНИЕ ЗАЩИТЫ ДОКУМЕНТА И СОСЛАНИЕ НА ЗАКОН РФ О ЗАЩИТЕ ПРАВА ПОКУПАТЕЛЯ, И ТАКЖЕ, ЧТО ЭТО БЫЛ ПЕРЕДАН НА РУКИ ТОВАРУ В ПОЛНОМ
 КОМПЛЕКТЕ, БЕЗ ЗАВЕРШЕНИЯ, В КОМПЬЮТЕРНОЙ СИСТЕМЕ (ТАБЛ). С ПОДПИСАТЕЛЯМИ ПО УТВЕРЖАЮЩИМ И ТЕРМИНАМ, ПРОИЗВОДИТЕЛЕМ ОБЪЕКТА ПРОДАЖИ.

Заказ оформил: (Марьяна Мелева)
 Отпустил: Получил: *Мелева Р*

Вашему вниманию! При покупке в нашем магазине комплекта шин или дисков Вы получаете 10% на услуги Шиномонтажа и Автосервиса!
 Не забывайте спрашивать купоны на скидку у менеджеров.

ОПЛАЧЕНО
 "___" _____ 2008

Рис. 8. Пример изображения финансового документа

словам, образуют на гистограмме один или несколько четко выраженных максимумов. Поэтому, анализируя гистограмму, можно вычислить характерный размер шрифта h_{font} , которым набран текст на странице, и, соответственно, выделить область на изображении, соответствующую текстовой информацией (области, соответствующие компонентам связности с высотой порядка h_{font}).

Зная область расположения текста на исходном документе, построим маску расслоения, после чего применим ее для получения графического и текстовых слоев (рис. 7).

Поскольку для выделения текстовых блоков используются быстрые алгоритмы морфологической фильтрации с прямоугольным окном, то очень важно, чтобы текстовые блоки были выровнены относительно осей изображения. Поэтому перед морфологией выполняется «выравнивание» изображения с опорой на быстрое преобразование Хафа [18].

Для цветных изображений финансовых документов (счетов-фактур, квитанций, договоров и пр.) не характерны перечисленные выше особенности изображения страницы книги, так как графические элементы (печати, подписи, рукописные пометки) часто накладываются на текстовые блоки (рис. 8). Следовательно, использовать для расслоения вышеописанный алгоритм неразумно. Построим схему расслоения, опираясь на цветовые характеристики изображения. Цветовая насыщенность [4] черного текста и белого фона близка к нулю, в то время как для синих печатей и подпи-

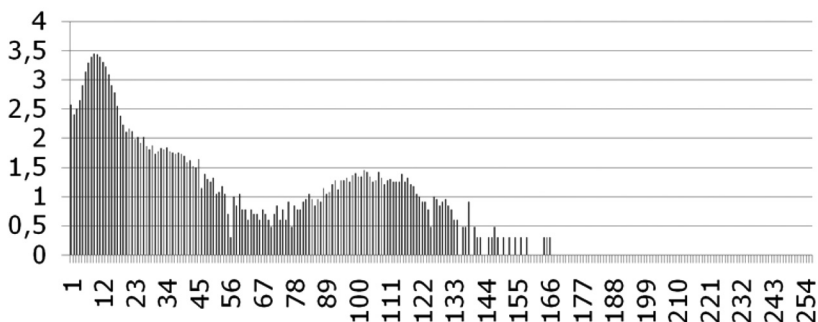


Рис. 9. Гистограмма цветовой насыщенности

сей это значение велико. Принимая во внимание это свойство, построим следующую схему расслоения.

Шаг 1. Построим гистограмму цветовой насыщенности (рис. 9), т. е. зависимость $y = \log N_x$, где N_x — количество пикселей изображения, насыщенность которых равна x .

Шаг 2. Заметим, что на гистограмме четко выделяются два класса: первый сформирован пикселями с малыми значениями цветовой насыщенности, второй — с большими значениями. Пиксели из первого класса составляют области изображения, соответствующие фону и черному тексту, из второго — графическую часть изображения. Найдем порог разделения двух классов t^* методом Отсу [15].

Шаг 3. Если величина порога t^* слишком мала (меньше, чем некоторое минимальное значение t_{\min}), то исходное изображение содержит только черный текст и не содержит каких-либо цветных графических элементов. В противном случае расслоим исходное изображение следующим образом: пиксель исходного изображения (x, y) принадлежит текстовому слою (рис. 10 а), если значение его цветовой насыщенности меньше порогового $s(x, y) < t^*$; иначе — пиксель (x, y) принадлежит графическому слою (рис. 10 б).

4. Распознавание текстового слоя

В результате расслоения исходное изображение разделяется на два новых: изображение, содержащее только текстовую информацию, и изображение, содержащее только графические элементы. Изображение, соответ-

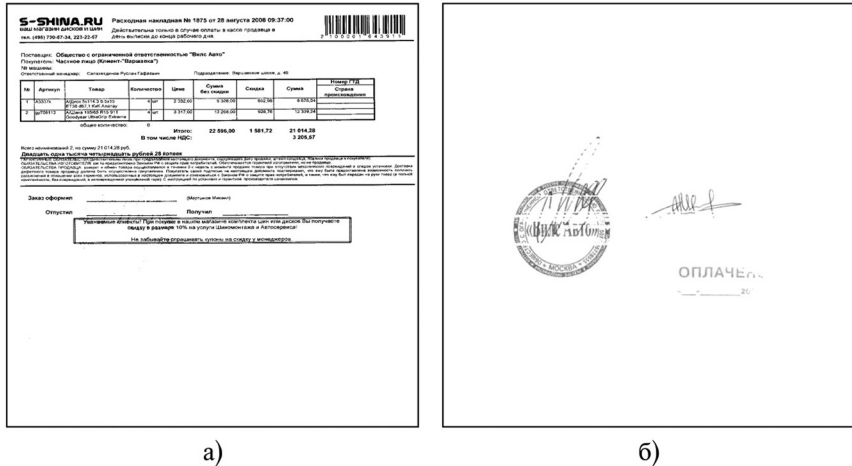


Рис. 10. Текстовый (а) и графический (б) слои изображения финансового документа

вующее текстовому слою, может быть легко распознано без какой-либо предварительной подготовки с помощью внешних OCR-систем.

В программной реализации технологии Cognitive PDF/A в качестве OCR-модуля выбрана система оптического распознавания текстов OCR CuneiForm [19], предназначенная для распознавания любых полиграфических шрифтов, получаемых с принтеров, за исключением декоративных и рукописных.

Система снабжена многими современными алгоритмами и технологиями, повышающими качество распознавания, например адаптивное распознавание и алгоритмы распознавания, использующие нейронные сети. Адаптивное распознавание [20] — метод, основанный на комбинации двух видов алгоритмов распознавания печатных символов: шрифтового (multi-font) и шрифтонезависимого (omnifont). Система генерирует внутренний шрифт для каждого вводимого документа, основываясь на хорошо пропечатанных символах, т. е. используется динамическая настройка (адаптация) на конкретные входные символы. Таким образом, метод совмещает универсальность и технологичность безшрифтового подхода и высокую точность распознавания шрифтового, что позволяет кардинальным образом повысить качество распознавания. Алгоритмы, использующие нейронные сети для распознавания символов, строятся следующим образом [21]. Поступающее на распознавание изображение символа (растр) приводится к некоторому стандартному размеру (нормализуется). Значения яркости в узлах нормализованного растра используются в качестве входных параметров нейронной сети. Число выходных параметров

нейронной сети равняется числу распознаваемых символов. Результатом распознавания является символ, которому соответствует наибольшее из значений выходного вектора нейронной сети. В систему встроены специальные алгоритмы для распознавания текста с матричного принтера, плохих ксерокопий факсов.

Первоначально система CuneiForm была разработана компанией Cognitive Technologies как коммерческий продукт и поставлялась с некоторыми моделями сканеров. Однако в 2008 году было анонсировано открытие исходных текстов программы в соответствии с лицензией BSD, что позволяет использовать данное программное обеспечение в сторонних продуктах.

5. Сжатие и упаковка в формат PDF/A

Полученные в результате расслоения текстовый и графический слои, а также распознанный текст сохраняются в формате PDF/A. Данный формат является стандартом ISO 19005–1:2005, базируется на описании стандарта PDF версии 1.4 от Adobe Systems Inc. и предназначен специально для долгосрочного архивного хранения электронных документов [22]. Несмотря на то что PDF/A является подмножеством формата PDF, существует ряд различий, обусловленных требованиями, предъявляемыми к PDF/A как к формату долгосрочного хранения электронных документов. Так, например, обязательным для PDF/A является:

- внедрение всех используемых шрифтов, в том числе шрифтов из списка «стандартных для PDF»;
- внедрение цветового профиля (если PDF/A-файл содержит изображение) — файла, в котором содержится информация о том, как выводное устройство (монитор, принтер и пр.) должно передать цвет. Важным является тот факт, что включаемый цветовой профиль должен быть аппаратно-независимым;
- обязательное наличие метаданных с указанием версии используемого формата, заголовка документа, списка авторов, краткого описания, даты создания и последней модификации файла документа, а также ключевых слова для осуществления поиска. Спецификацией PDF/A также оговорен формат представления метаданных — Adobe Extensible Metadata Platform (XMP) [23].

Для увеличения коэффициента сжатия графический и текстовый слои сжимаются разными способами. В силу специфики содержимого графический слой приводится к разрешению 100 DPI и кодируется алгоритмом JPEG. Текстовый слой заключает в себе основную информацию документа, следовательно, текстовый слой сохраняется в исходном разрешении, а для кодирования используются алгоритм сжатия без потерь CCITT Group 4.

Заключение

Для определения эффективности технологии было выбрано несколько изображений различного типа, которые были сжаты алгоритмами JPEG (уровень компрессии, при котором сохраняется читабельность), DjVu и Cognitive PDF/A.

По результатам сравнения видно (см. табл. 1), что описанная технология по степени сжатия на порядок обходит JPEG, однако существенно проигрывает DjVu. Такую значительную разницу в размере можно объяснить тем, что файл в формате PDF/A помимо полезной информации (собственно изображений и распознанного текста) содержит также вспомогательные данные, необходимые для долгосрочного хранения. К тому же формат DjVu использует более совершенный алгоритм сжатия с потерями (IW44), основанный на вейвлет-преобразованиях, который не поддерживается настоящей версией формата PDF/A. Однако несмотря на меньший размер файлов DjVu качество сжатия офисных документов довольно низкое, что особенно проявляется в районе печатей и подписей.

Изложенный в работе алгоритм оцифровки текстовых документов реализован программно и внедрен в систему Cognitive Forms [2] в качестве средства архивирования документов. Данная программная подсистема была промышленно внедрена и используется для массового ввода документов. В проекте, реализованном для ОАО «Магнитогорский металлургический комбинат», разработанный алгоритм сжатия применяется для обработки бухгалтерских документов (счета-фактуры, накладные, акты, договора). В проекте, реализованном для ООО «РусФинанс», алгоритмы сжатия используются для обработки изображений документов, необходимых для выдачи кредитов (заявления на выдачу кредита, свидетельства ИНН, ксерокопии паспорта и водительского удостоверения и пр.).

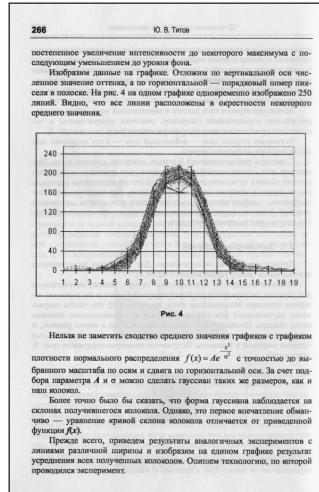
Таблица 1

Сравнение различных технологий сжатия

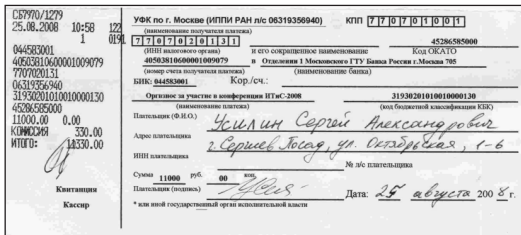
Изображение	JPEG	DjVu	Cognitive PDF/A
Рис. 11 а	718 KB	50 KB	94 KB
Рис. 11 б	373 KB	29 KB	51 KB
Рис. 11 в	642 KB	23 KB	52 KB
Рис. 11 г	120 KB	20 KB	40 KB
Рис. 11 д	292 KB	17 KB	29 KB



а)



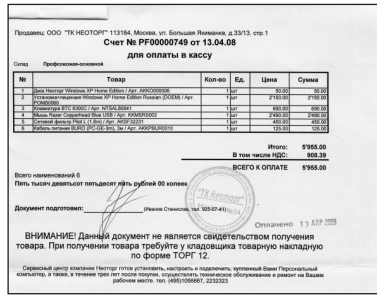
б)



в)



г)



д)

Рис. 11. Примеры изображений различного типа: страница книги с цветными иллюстрациями (а); страница книги с графикой (б); квитанция (в); финансовый документ с насыщенной печатью (г); финансовый документ с тусклой печатью и штампом (д)

Литература

1. Арлазаров В. Л., Славин О. А. Алгоритмы распознавания и технологии ввода текстов в ЭВМ // Информационные технологии и вычислительные системы. 1996. Т. 6. № 1. С. 48–54.
2. Арлазаров В. В., Постников В. В., Шоломов Д. Л. Cognitive Forms — система массового ввода структурированных документов // Управление информационными потоками. М.: URSS, 2002. С. 37–49.
3. Ватолин Д., Ратушняк А., Смирнов М., Юкин В. Методы сжатия данных. Алгоритмы сжатия изображений. М.: Диалог-МИФИ, 2002. 99 с.
4. Гонсалес Р., Вудс Р. Цифровая обработка изображений. М.: Техносфера, 2005. 1072 с.
5. Fletcher L. A., Kasturi R. A Robust Algorithm for Text String Separation from Mixed Text/Graphics Images // IEEE Transactions on Pattern Analysis and Machine Intelligence. 1988. Vol. 10. № 6. С. 910–918.
6. Hasan Y., Karam L. Morphological text extraction from images // IEEE Transactions on Image Processing. 2000. Vol. 9. № 11. С. 1978–1983.
7. Hirata N. S. T., Barbera J., Terada R. Text segmentation by automatically designed morphological operators // Computer Graphics and Image Processing. 2000. С. 284–291.
8. Yuan Q., Tan C. L. Text extraction from gray scale document images using edge information // IEEE Document Analysis and Recognition Processing. 2001. С. 302–306.
9. Cao R. Separation of overlapping text from graphics // IEEE Document Analysis and Recognition Processing. 2001. С. 44–48.
10. Lizardtech DjVu Reference. <http://www.lizardtech.com>
11. Postnikov V. V. Flexible Forms Identification // Proceedings of the 5th German-Russian Workshop on Pattern Recognition and Image Understanding (GRWS98). Hamburg, 1999.
12. Усилин С. А., Николаев Д. П., Постников В. В. Быстрый алгоритм совмещения изображений документов в произвольной геометрической модели // Труды конференции «Информационные технологии и системы». 2008. С. 471–477.
13. Безматерных П. В., Николаев Д. П., Постников В. В. Метод идентификации типа документа по структуре проекций его изображения на координатные оси // Труды конференции «Информационные технологии и системы». 2008. С. 498–501.
14. Nikolaev D. P. Segmentation-based binarization method for color document images // Proceedings of 6th Open Russian-German Workshop on Pattern Recognition and Image Understanding. 2003. С. 190–193.
15. Trier O. D., Taxt T. Evaluation of binarization methods for document images // IEEE Transactions on Pattern Analysis and Machine Intelligence. 1995. Vol. 17. № 3. С. 312–315.
16. Куроптев А. В., Николаев Д. П., Постников В. В., Усилин С. А. Выделение графических примитивов и текстовых блоков на изображениях документов с помощью морфологических операций // Труды 51-й научной конференции МФТИ. Современные проблемы фундаментальных и прикладных наук. Часть 9. Инновации и высокие технологии. М.: МФТИ, 2008. С. 29–31.

17. *Herk M. van.* A Fast Algorithm for Local Minimum and Maximum Filters on Rectangular and Octagonal Kernels // Pattern Recognition Letters. 1992. С. 517–521.
18. *Nikolaev D. P., Karpenko S. M., Nikolaev I. P., Nikolayev P. P.* Hough Transform: Underestimated Tool in the Computer Vision Field // 22st Europ. Conf. Model. Simulat. ECMS 2008. Nicosia, Cyprus, 2008. P. 238–243.
19. OCR CuneiForm. <http://www.cuneiform.ru>
20. *Арлазаров В. Л., Астахов А. Д., Троянкер В. В., Котович Н. В.* Адаптивное распознавание символов // Интеллектуальные технологии ввода и обработки информации. М.: URSS, 1998. С. 39–56.
21. *Мисюрёв А. В.* Использование искусственных нейронных сетей для распознавания рукопечатных символов // Интеллектуальные технологии ввода и обработки информации. 1998. С. 122–127.
22. ISO 19005–1:2005. Document management — Electronic document file format for long-term preservation. Part 1: Use of PDF 1.4 (PDF/A-1).
23. Adobe Systems Incorporated. Extensible Metadata Platform (XMP) Specification. <http://www.adobe.com>