

Жадаев А. Г.

Сканирование и распознавание текстов

Самоучитель по работе с ABBYY® FineReader 10



Москва, 2010

УДК 32.973.26-018.2
ББК 004.4
Ж15

Ж15 **Жадаев А. Г.**

Сканирование и распознавание текстов. Самоучитель по работе с ABBYY® FineReader 10. – М.: ДМК Пресс, 2010. – 248 с.: ил.

ISBN 978-5-94074-595-2

Работать с электронными документами во многом удобнее и проще, чем с их бумажными аналогами. Электронный документ можно редактировать, использовать при создании собственных работ, его легко копировать и пересылать по электронной почте. Вместе с тем, многие материалы изначально доступны нам в неотредактированном виде (бумажные или отсканированные документы, цифровые фотографии). Программа ABBYY® FineReader – лучший инструмент для создания электронных копий любых печатных материалов: книг, справочников, журналов, договоров, бланков.

Книга включает описание приемов сканирования и распознавания разных оригиналов – от простых книжных страниц до сложно оформленных документов. А приведенные скриншоты программы позволят читателю быстро освоить интерфейс ABBYY® FineReader и получить практические навыки по работе с программой.

Изложение материала сопровождается практическими примерами. Читатели, которые еще не пробовали самостоятельно переводить печатные материалы в электронный вид, найдут в этой книге простое пошаговое руководство. Для тех же, кто хочет в совершенстве освоить работу с программой, книга откроет многочисленные тонкости настройки для эффективного использования ABBYY® FineReader.

УДК 32.973.26-018.2
ББК 004.4

Все права защищены. Любая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами без письменного разрешения владельцев авторских прав.

Материал, изложенный в данной книге, многократно проверен. Но поскольку вероятность технических ошибок все равно существует, издательство не может гарантировать абсолютную точность и правильность приводимых сведений. В связи с этим издательство не несет ответственности за возможные ошибки, связанные с использованием книги.

ISBN 978-5-94074-595-2

© Жадаев А. Г., 2010

© Оформление, издание, ДМК Пресс, 2010

Содержание

Глава 1

Текст и графика в компьютере	7
Форматы файлов	9
Текстовые форматы	10
Графические файлы	11
Составные (сложные) документы	12
Оптическое распознавание символов (OCR)	14
От чего зависит качество распознавания?	16
Системные требования	17
Резюме	17

Глава 2

Быстрый старт	19
Сканирование в MS Word, MS Excel, PDF	23
Конвертирование изображений и PDF в документ Microsoft Word	28
Вызов сценариев из контекстного меню файла	30
Сканирование и сохранение изображений	31
Резюме	33

Глава 3

Работа в пошаговом режиме	34
Окно программы и настройка рабочего пространства	34
Рабочие окна	35
Окно Страницы	36
Окно Изображение	38
Окно Текст	39
Окно Крупный план	40
Изменение расположения рабочих окон	41
Настройка панелей инструментов	45
Диалоговое окно Опции	51
Документ FineReader	55
Резюме	57

Глава 4

Получение изображений	59
Работа со сканером	60
Параметры сканеров	60
Драйвер и настройки сканера	64
Разрешение	66
Режим сканирования	67
Яркость	67
Параметры страницы	68
Сканирование многостраничных документов	69
Работа с цифровой камерой	69
Параметры цифровых камер	69
Техника съемки	73
Расстояние	73
Освещение	73
Баланс белого	74
Повышение четкости изображения	74
Работа над ошибками	75
Частные случаи	77
Крупноформатные оригиналы	78
Книги	79
Резюме	81

Глава 5

Обработка и анализ изображений	83
Обработка изображения	84
Настройка автоматической обработки	85
Обработка в Редакторе изображений	86
Анализ изображений	92
Области изображения	92
Исправление разбивки на области	94
Свойства области	99
Использование шаблонов областей	101
Анализ таблиц	104
Резюме	107

Глава 6

Распознавание текстов	109
Применение пользовательского эталона	109
Общие правила работы с пользовательскими эталонами	109
Пример обучения и использования эталона	111
Редактирование пользовательских эталонов	117
Распознавание многоязычных документов	120
Пример распознавания двуязычного документа	120
Выбор языка для распознавания документа	123
Создание группы языков	124

Создание пользовательского языка	128
Пример распознавания текста с помощью регулярных выражений	131
Использование словарей	135
Общие правила работы со словарями	135
Пример редактирования и применения пользовательского словаря	139
Резюме	141

Глава 7

Проверка и корректировка распознанного документа 143

Проверка и корректировка документа в программе FineReader	143
Пример корректировки документа в окне Текст	143
Пример проверки и корректировки текста в диалоге Проверка	151
Использование стилей	158
Общие сведения о стилях в FineReader	159
Пример создания и применения пользовательского стиля	159
Окончательная обработка распознанного документа в программе Microsoft Word	162
Пример избавления импортированного в Word документа от стилей FineReader	162
Пример обработки распознанного документа в Word	168
Корректировка таблиц	173
Настройка панели инструментов для работы с таблицами	173
Пример корректировки таблицы	175
Резюме	177

Глава 8

Сохранение распознанного документа 178

Передача документа в приложение	180
Пример передачи распознанного документа в Microsoft Word	180
Пример передачи распознанного документа в Microsoft Excel	181
Пример передачи распознанного документа в Adobe Reader	182
Пример передачи распознанного документа в веб-браузер	183
Сохранение документа в файл	184
Настройка параметров сохранения	184
Примеры сохранения распознанного документа в формате Word	194
Пример сохранения распознанного документа в формате PDF	198
Пример сохранения распознанного документа в формате HTML	202
Пример сохранения распознанного документа в формате TXT	205
Пример сохранения распознанного документа в формате Excel	205
Как сохранить документ FineReader	208
Резюме	210

Глава 9

Сценарии 211

Создание пользовательского сценария	211
Пример: распознавание платежного поручения	212

Другие действия сценариев	222
Менеджер сценариев	225
Копирование сценария	225
Изменение сценария	226
Удаление сценария	226
Экспорт и импорт сценариев	226
Использование сценариев	228
Резюме	230

Глава 10

ABBYY Screenshot Reader	231
Интерфейс и настройки программы	231
Работа с программой	233
Захват изображения	234
Передача в буфер обмена	235
Передача в приложения Microsoft Office	237
Операция Изображение в ABBYY® FineReader	237
Сохранение в файл	238
Примеры использования программы	241
Копирование текста документов с экрана	241
Список файлов	242
Снимки интерфейса	245
Резюме	246

Глава 1

Текст и графика в компьютере

Любая информация в компьютере хранится и обрабатывается в цифровом виде. Хотя в этой книге повсюду упоминаются слова «текст» и «графика», на самом деле компьютер работает исключительно с цифровой информацией. Любая информация в компьютере хранится в виде файлов, а файл – всего лишь последовательность двоичных чисел (единиц и нулей). С помощью двоичных чисел можно закодировать все, что угодно: от команд, которые должен выполнять сам компьютер, до текста, изображений, музыки и фильмов.

Нужно лишь договориться, что в каждом случае будут означать группы двоичных чисел, из которых состоит файл, как их должны воспринимать и обрабатывать компьютерные программы. Алгоритм (правило), в соответствии с которым данные превращаются в цифры и помещаются в файл, называют *форматом файла*. Разумеется, для разного рода информации требуются разные способы кодирования. Более того, одну и ту же информацию можно представить в цифровом виде различными способами, поэтому форматов было изобретено много.

Например, Объединенная группа экспертов в области фотографии (Joint Photographic Experts Group) разработала набор спецификаций, позволяющих значительно сократить размер файла, в котором хранится изображение. Как вы уже догадались по первым буквам английского названия группы, в результате появился формат JPEG. С файлами этого формата работают и компьютерные программы, и цифровые фотоаппараты, и некоторые бытовые устройства.

В этой книге речь идет о файлах, содержащих данные двух видов. Это *текст* и *графика*.

Для компьютера **текст** – последовательность символов. Символами являются буквы разных алфавитов, цифры, знаки препинания. Пробел, разделяющий слова, – тоже символ. Каждому символу соответствует определенный числовой код.

С таким представлением текста компьютеру легко выполнять самые обычные математические операции. Например, он может найти в тексте все символы с определенным кодом, вставить в указанное место коды введенных с клавиатуры символов или заменить определенную последовательность на другую. На этом основана работа всех программ – текстовых редакторов. Простой пример – программа Блокнот (Notepad), поставляемая в составе ОС Windows. В текстовом редакторе можно выделить часть текста, скопировать ее в буфер обмена, а затем вставить в другое место того же документа или в другой документ; найти в тексте определенные последовательности символов и т. д.

Графика, или рисунок, в компьютерном представлении – набор отдельных точек, о каждой из которых известны ее положение и цвет.

ПРИМЕЧАНИЕ

В компьютерах представление цвета зависит от выбранной цветовой схемы.

Точки, из которых состоит изображение, называют пикселями (pixel, от англ. *PICTure'S ELeMent* или *PICTure Cell* – элемент или клетка изображения). Часто их зовут и просто точками (Dots). Таким образом, изображение является «картой точек», по-английски *Bitmap*. Графический файл содержит информацию обо всех точках изображения.

На рис. 1.1 показан пример рисунка «с точки зрения компьютера». Каждая клеточка изображает одну точку (пиксел). Мы видим, что ширина рисунка со-

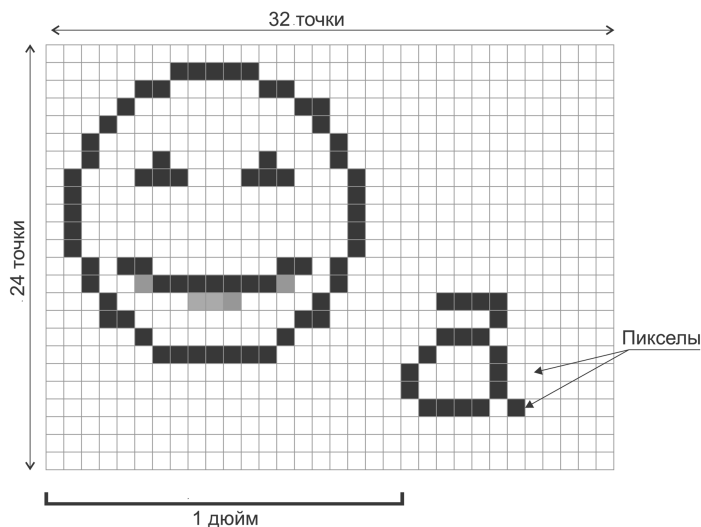


Рис. 1.1 ▾ Пример изображения

ставляет 32 точки, а высота – 24 точки. Белые, или «пустые», точки тоже считаются! В таком случае говорят, что размер этого рисунка – 32г24 точки (пиксела). Вместе с распространением цифровых камер размеры изображения чаще стали выражать общим числом точек. Например, размер этого рисунка – 768 пикселей, или примерно 0,0008 мегапиксела (Мпикс).

На том же примере продемонстрируем еще одну характеристику компьютерного изображения – *разрешение*. Для сравнения под рисунком помещена линейка длиной 1 дюйм (2,54 см). На рисунке ей соответствуют 20 точек. Таким образом, разрешение этого изображения составляет 20 точек на дюйм, или 20 DPI (Dots Per Inch – точек на дюйм).

Возникает вопрос: как же связаны размер изображения и его разрешение? У цифрового изображения есть лишь «истинный размер в пикселах». О разрешении можно говорить, только если мы одновременно уточним, с оригинала какого размера, в сантиметрах или дюймах, было получено (отсканировано, сфотографировано) это изображение, или каков должен быть размер изображения при его выводе на экран. На практике любой графический файл содержит служебную информацию (заголовок), где именно это и указано: размер в пикселах и разрешение в DPI одновременно.

Для работы в программе FineReader изображение обычно получают со сканера или цифровой камеры. При этом в настройках сканера задают разрешение получаемого изображения, например 150 или 300 DPI. Размер изображения будет зависеть от размера сканируемого оригинала. При сканировании с разрешением 300 DPI стандартной книжной страницы получится изображение размером примерно 2400×1600 точек, или, по «фотографической» терминологии, около 4 Мпикс.

Цифровой фотоаппарат выдает снимки определенного размера. Часто цифровые камеры позволяют регулировать размер снимка. Максимально возможный размер является одной из самых важных ее характеристик. В обиходе такой параметр называют «разрешением», хотя на самом деле это именно *размер* выходного изображения в точках.

Нужно различать *размер изображения* и *размер файла*. Первый выражается в количестве точек (пикселей) и характеризует размер самого изображения. Размер (объем) файла показывает, сколько места занимает файл на диске, и измеряется в байтах (килобайтах, мегабайтах). Размер файла зависит не только от размера изображения, но и от формата файла, а также использованных алгоритмов сжатия: если при сохранении изображения применяется сжатие информации, то файл получается меньше.

Форматы файлов

Сложилось так, что разработчики программ регулярно создавали и предлагали новые форматы файлов. В настоящее время существуют сотни различных форматов, каждый из которых обладает своими особенностями, достоинствами и, возможно, недостатками по сравнению с другими. Тип файла (документа) – бо-

лее общая характеристика содержимого файла, например «текст», «изображение», «звук», «архивы» и т. д. Так, к графическим форматам, в которых компьютер работает с изображениями, относятся форматы **BMP**, **TIFF**, **GIF**, **JPEG** и многие другие.

Любая программа способна открывать, обрабатывать и сохранять файлы только определенных типов и форматов. Некоторые программы поддерживают лишь один-два формата, но чаще приложения работают с целым рядом форматов. Рассмотрим некоторые типы и форматы файлов (документов), которые могут нам встретиться при работе с программой FineReader. Заметим, что от версии к версии набор поддерживаемых форматов, как графических, так и текстовых, может изменяться.

Текстовые форматы

Самым старым и простым типом файлов является текстовый. По сути, текстовый файл представляет собой только последовательность символов (букв, цифр, знаков препинания и математических символов), отображаемых на экране или выводимых на печать. Помимо них, в текстовом файле могут содержаться «символ перевода строки», «символ разрыва строки», «символ табуляции» и «символ конца страницы», которые не отображаются на экране и не распечатываются на бумаге. Это так называемые «непечатаемые символы».

Текстовый файл не содержит информации ни о размере или начертании шрифта, ни о выравнивании строк или отступе первой строки абзаца. То, как отображается текст на экране, зависит только от настроек программы, в которой открыт этот файл.

Текстовые файлы имеют расширение **ТХТ**. Работу с такими файлами должны поддерживать все программы, которые способны так или иначе обрабатывать текстовую информацию. На основе текстового файла было создано еще несколько форматов для более узких применений.

Файлы формата **CSV** (*Comma Separated Values* – значения, разделенные запятыми) являются такими же текстовыми файлами, как файлы формата **txt**, но специально предназначены для хранения данных в виде списков или таблиц. Каждая строка считается строкой таблицы, а запятые разделяют текст на части, которые должны быть помещены в отдельные ячейки. Когда вы открываете такой файл обычным текстовым редактором, например Блокнотом, то видите обычный текст. Если же открыть файл **CSV** программой для работы с электронными таблицами, например Microsoft Excel, то такая программа воспримет все запятые как разделители ячеек и автоматически разместит данные в ячейки таблицы. Иногда для разделения значений используется точка с запятой. При открытии или сохранении файла программы предлагают уточнить, какой символ следует считать разделителем.

В формате **HTML** (*HyperText Markup Language* – язык разметки гипертекста) создаются файлы, которые служат основой веб-страниц и предназначены для просмотра с помощью программ-браузеров. Эти файлы имеют расширение **HTML** или **HTM**. По существу, это тоже текстовые файлы, но в тексте есть спе-