

Компьютерная шпаргалка

Ответы на вопросы всегда под рукой

программа

FineReader

распознавание отсканированного
и сфотографированного текста,

а также перевод PDF-файлов



в текстовые



Корнеев А. П., Иванова А. А., Прокди Р. Г.

**Программа FineReader:
распознавание отсканированного
и сфотографированного текста, а также
перевод PDF-файлов
в текстовые**

КОМПЬЮТЕРНАЯ ШПАРГАЛКА



Наука и Техника, Санкт-Петербург, 2010

Корнеев А. П., Иванова А. А., Прокди Р. Г.

Программа FineReader: распознавание отсканированного и сфотографированного текста, а также перевод PDF-файлов в текстовые. Компьютерная шпаргалка. – СПб.: Наука и Техника, 2010. – 80 с.: ил.

Серия «Компьютерная шпаргалка»

В этой книжке вы найдете описание методики распознавания отсканированного и/или сфотографированного текста с помощью самой лучшей и популярной программы FineReader, специально для этого предназначеннной. Кроме того, вы узнаете, как переводить PDF-файлы в текстовые, отсканированные таблицы в таблицы Excel и многое другое. Книжка будет несомненно полезна всем пользователям компьютеров.

Контактные телефоны издательства:

(812) 412-70-25, 412-70-26

(044) 516-38-66

Официальный сайт: www.nit.com.ru

© Прокди Р. Г.

© Наука и Техника (оригинал-макет), 2010

Содержание

3

Глава I. Программа ABBYY FineReader 9.0	
и ее возможности	7
Глава II. Распознавание текста в ABBYY FineReader 9.0	
2.1. Работа с программой.....	11
2.2. Сканирование и распознавание текста с переводом в текстовый документ Microsoft Word	16
Глава III. Настройка распознавания текста, задание областей распознавания, полный контроль	20
Глава IV. Специальные режимы распознавания в FineReader.....	45
Как отсканировать и преобразовать в PDF	45
Распознавание таблиц и сохранение их в формате Excel	46
Сканирование изображений (без текста) в Word	48

Как преобразовать PDF-файл в Word'овский текстовый файл	51
Глава V. Повышение качества распознавания текста	53
Режим обучения при распознавании текста	53
Глава VI. Распознавание текста с цифровых изображений (фотографий)	56
Как правильно фотографировать документы, чтобы их фотографии потом можно было распознать (преобразовать в текстовый документ)	56
Распознавание текста с цифровой фотографии	58
Глава VII. Полезные приемы работы с FineReader'ом. Решение проблем с распознаванием текста	60
Как отключить автоматический запуск распознавания текста после сканирования	60
Как сделать так, чтобы при распознавании отсканированного/ сфотографированного разворота книги FineReader автоматически разбивал разворот на отдельные страницы	62

Как повернуть страницу в окне	62
Как указать язык текста в документе для распознавания.....	64
Как включить или исключить какую-либо область страницы из распознавания	65
Как подправить границы той или иной области	67
Как указать, что данный фрагмент страницы является изображением	67
Использование режима тщательного распознавания	70
Как дать команду на повторное распознавание текста после произведенных перенастроек	71
Как перераспознать отдельную область на странице	72
Что делать, если в распознанном тексте некорректно отображается шрифт или на месте некоторых букв стоят значки «?» или «□»	73
Что делать, если в исходном тексте присутствуют нестандартные (декоративные, математические и т.д.) символы	74
Что делать, если таблица в исходном документе не определена	75
Как вручную задать сетку таблицы, указать разбиение на столбцы и строки	76

Как сохранить документ FineReader с отсканированными страницами для
возможности дальнейшего повторного распознавания 78

Программа ABBYY FineReader 9.0 и ее возможности

Сегодня, в век стремительно развивающихся цифровых технологий, тема электронного распознавания страниц многим может показаться неактуальной. Но на самом деле это далеко не так. Конечно, если бы лет десять назад любому секретарю было известно о существовании такого инструмента, то никакие деньги не остановили бы его приобрести данную программу, даже за собственный счет. Ведь она бы сэкономила не менее трети всего его рабочего времени за год. И это на самом деле так.

И если сегодня большинство документов можно отредактировать простым открытием файла, то это не означает, что данные виды программ потеряли свою сущность. Ведь до сих пор существуют миллионы копий документов, которые не переведены в цифровой формат. Например, вам срочно понадобилось внести изменения в сметный документ, который последний раз редактировался 15 лет назад и до сих пор пылится в архиве

8 организации. Смета состоит из 100 листов + рисунки и чертежи, присутствующие на каждом третьем. А времени дано всего 2 недели, или фирма будет подвергнута серьезным штрафным санкциям. Вы прекрасно понимаете, что вручную это сделать реально только в случае привлечения к данному заданию большого количества людей. Но это невозможно.

Вот примерно для таких ситуаций и нужен FineReader. С помощью него вы легко перелопатите стопы неоцифрованных документов, приведя их в удобный для редактирования в любом текстовом редакторе вид. Современные системы распознавания текстов OCR (Optical character recognition) распознают текст независимо от его шрифта, формата, а также языка. И хотя на рынке присутствует великое множество таких программ, явным и непревзойденным лидером среди них является FineReader от компании ABBYY, которая уже давно занимается проблемой качественного преобразования текста со сканированного изображения в текстовый файл.

На момент написания шпаргалки вышла 9-я версия этой замечательной программы, и существует она в нескольких модификациях: FineReader Home Edition (для домашних пользователей), FineReader Corporate Edition (для ведения документооборота), а также FineReader Professional

Edition, которая сочетает в себе простоту и в то же время всю мощь для распознавания многостраничных документов, при этом имея удобный интерфейс.

Мы рассмотрим именно версию ABBYY FineReader 9.0 Professional Edition, чтобы вы смогли оценить все её непревзойденные возможности. На сайте компании (<http://www.abbyy.ru/Download/>) представлена пробная версия продукта, которую вы сможете скачать и использовать в течение 15 дней или распознать 50 страниц. При этом сохранять или напечатать за раз вы сможете не более 1 страницы документа.

Особенностью 9-й версии является появление технологии адаптивного распознавания документа ADRT. Суть её работы в следующем. Если раньше анализ велся постранично, то сейчас он будет проводиться для всего документа в целом. Результатом стало безошибочное распознавание не только основного текста, но и всех прилагающихся к нему элементов: колонтитулов, сносок, подлиссий, нумераций страниц и т.д. Появилась возможность конвертации результата распознавания в новые версии Microsoft Office 2007, поддерживающие типы документов DOCX и XLSX. Кроме этого, возможно пересыпать их по электронной почте, а также реализо-

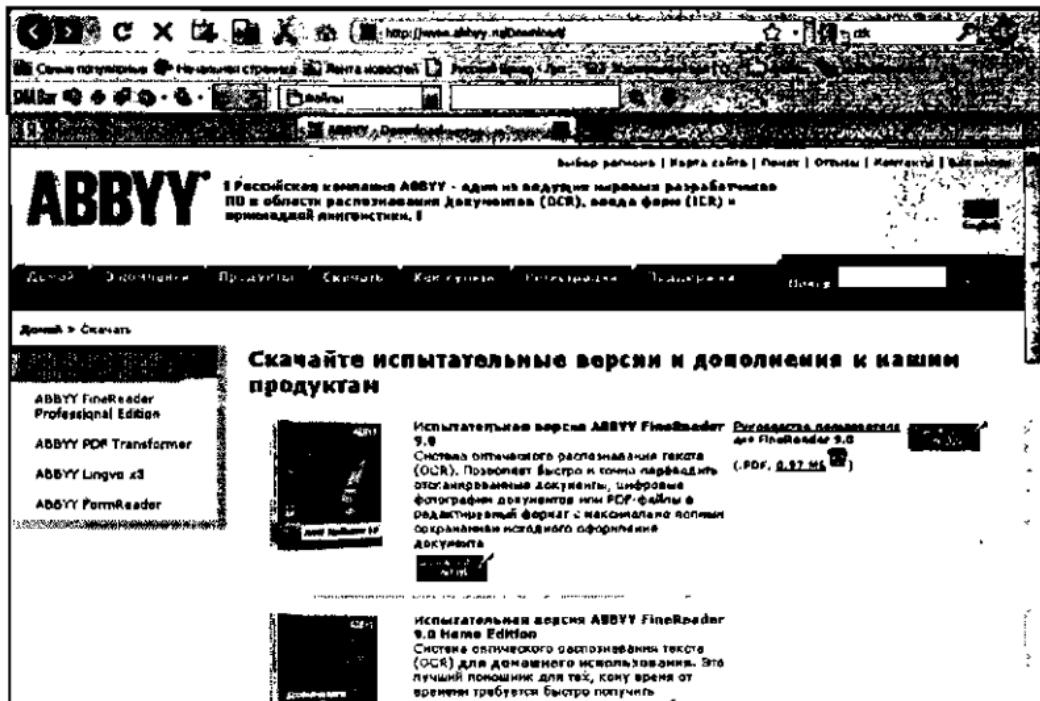


Рис. 1.1. На сайте <http://www.abbyy.ru/Download/> можно скачать программу FineReader

Распознавание текста в ABBYY FineReader 9.0

2.1. Работа с программой

Установка происходит в штатном режиме. Выбираем стандартный режим, который подразумевает полную установку программы. В результате в меню Пуск появится новый пункт **ABBYY FineReader 9.0**, который в свою очередь будет иметь подпункты (рис. 2.1).

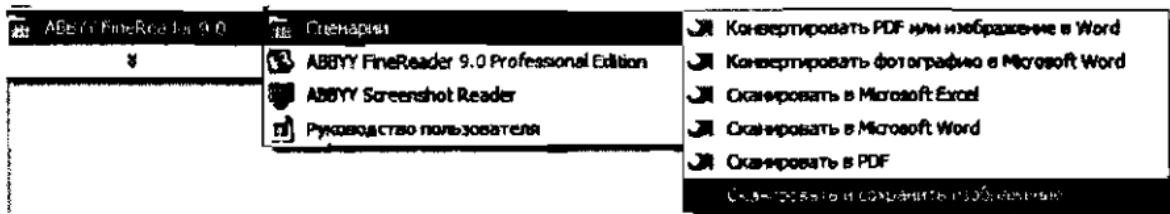


Рис. 2.1. Структура ABBYY FineReader 9.0 в меню Пуск

12

Ярлыки на рабочем столе не предусмотрены при установке программы, поэтому либо устанавливаем их самостоятельно, либо пользуемся пунктами меню.

- **Сценарии** – данный пункт меню подразумевает запуск определенной операции, не открывая главной программы. Вы можете отдельно выполнять следующие действия: *Конвертировать PDF или изображение в Word* (уже готовые файлы изображений или PDF-формата конвертируются в формат Word), *Конвертировать фотографию в Microsoft Word* (то же самое для фотографий), *Сканировать в Microsoft Word*, *Сканировать в Microsoft Excel*, *Сканировать в PDF*, *Сканировать и сохранить изображение* (сканирование всего, что возможно, и сохранение в формате изображений).
- **ABBYY Screenshot Reader** – можно распознать текст, находящийся в любой части экрана. Для этого предусмотрен механизм выполнения скриншотов.
- **ABBYY FineReader 9.0 Professional Edition** – запуск главной программы.

Запускаем программу. Откроется главное окно FineReader. При первоначальном запуске в окне отображается список сценариев, совершаемых программой (подобно меню Пуск) (рис. 2.2).

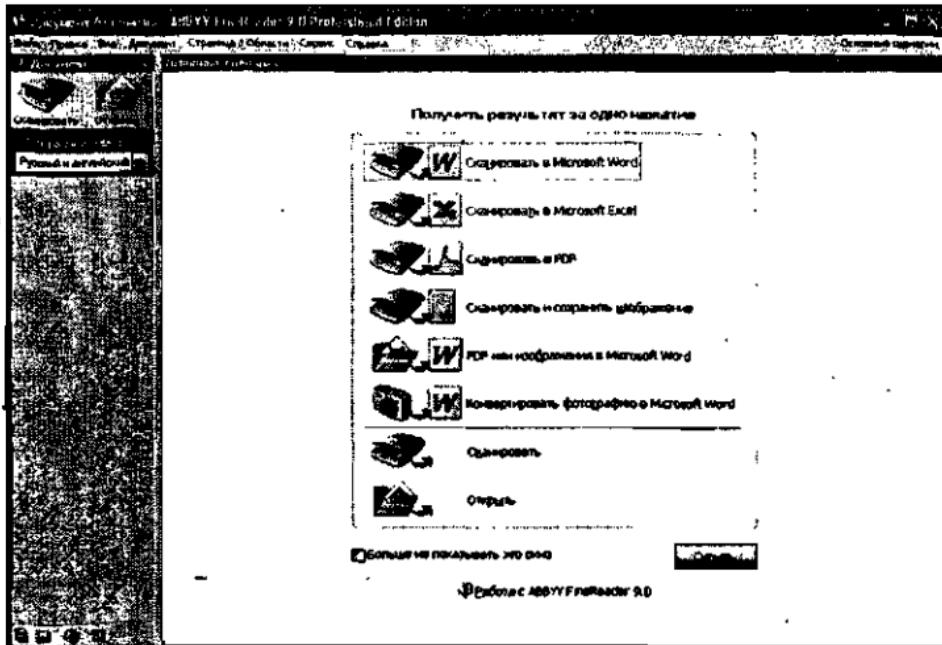


Рис. 2.2. Список сценариев FineReader

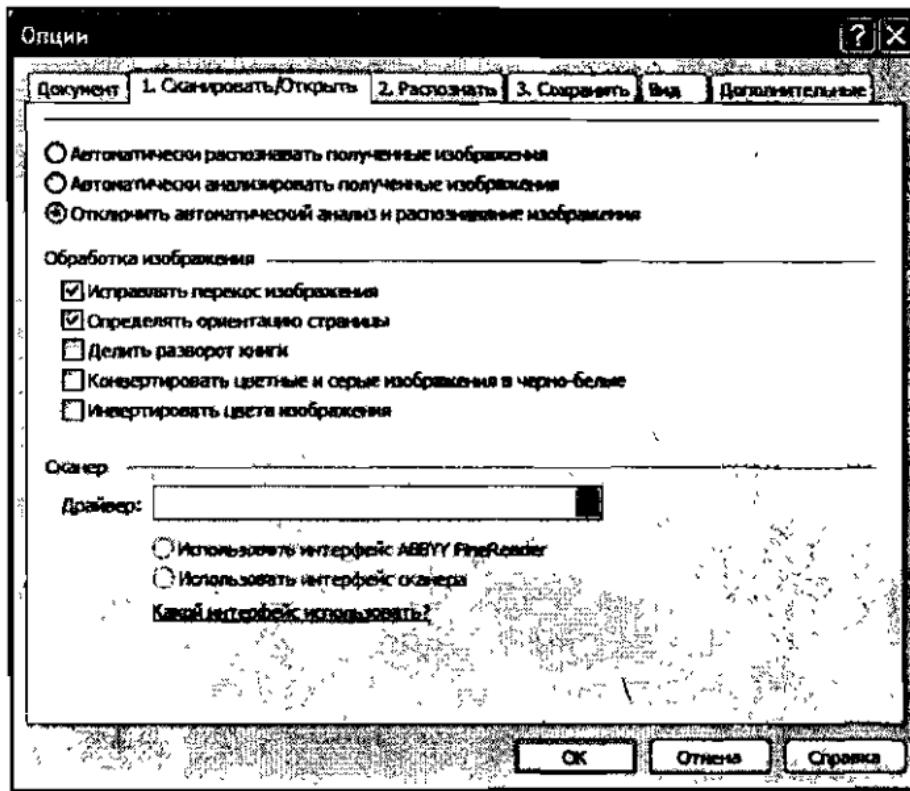


Рис. 2.3. Настройка программы

Установите флажок, чтобы они не отображались при запуске, и щелкните кнопку Скрыть. В следующий раз любую из операций можно вызвать с помощью главного меню программы **Сервис - Сценарии ABBYY FineReader 9.0**.

Для тех, кто малознаком с программой, да и вообще для получения более стабильного результата, советуем сразу же поменять стандартную настройку распознавания текста. Для этого щелкните желтый значок, расположенный в левом нижнем углу программы; открывающий настройки программы. Откроется окно **Опции**. Щелкните вкладку **Сканировать/Открыть** и установите переключатель опции **Отключить автоматический анализ и распознавание изображения** (рис. 2.3).

Ну вот, собственно, и все приготовления. Можно приступить к выполнению сценариев. Единственное, что у вас должно быть уже готово, так это установленный и работоспособный сканер. Поэтому это полностью оставим на вашей совести. Хотя можно обойтись и без сканера, например, имея уже отсканированные в другом месте файлы.

2.2. Сканирование и распознавание текста с переводом в текстовый документ Microsoft Word

Для удобного перевода и распознавания отсканированного текста в текстовый Word-документ при запуске FineReader'a в списке сценариев выберите **Сканировать в Microsoft Word**. Очень простая, как и все остальные, операция.

Просто вставьте интересующее вас изображение текста в сканер и запустите соответствующий сценарий. Откроется диалоговое окно выбора сканера, где вы должны будете подтвердить правильность отображаемого в поле сканера либо выбрать нужный из списка (рис. 2.4).

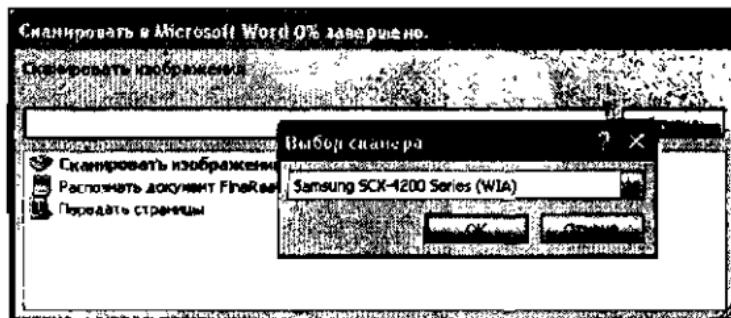


Рис. 2.4. Выбор сканера

Затем откроется окно предварительного просмотра сканируемого изображения (рис. 2.5). Чтобы увидеть его, следует в окне нажать кнопку **Просмотр**. Сканируемая область будет помечена синей рамкой, перемещая границы которой вы сможете указать что нужно, а что не нужно сканировать.

В раскрывающемся списке **Разрешение** вы можете указать то разрешение, которое должно использоваться при сканировании. Чем мельче шрифт, тем разрешение должно быть выше. Для нормального (а для большого и подавно) размера шрифта (не менее 10 пунктов), как правило, используется установленное по умолчанию значение 300 dpi.

В раскрывающемся списке **Режим сканирования** можно выбрать цветность сканирования. Если вы хотите, чтобы сканирование/распознавание было произведено с учетом цвета, то следует установить значение **Цветной**. По умолчанию же установлено значение **Серый**, что наиболее часто подходит для обычного сканирования/распознавания текста и требует меньше ресурсов компьютера. В то же время использование цветного режима может повысить качество распознавания в сложных случаях, так как программа еще сможет анализировать не просто черно-белые оттенки, но и разность в цвете.

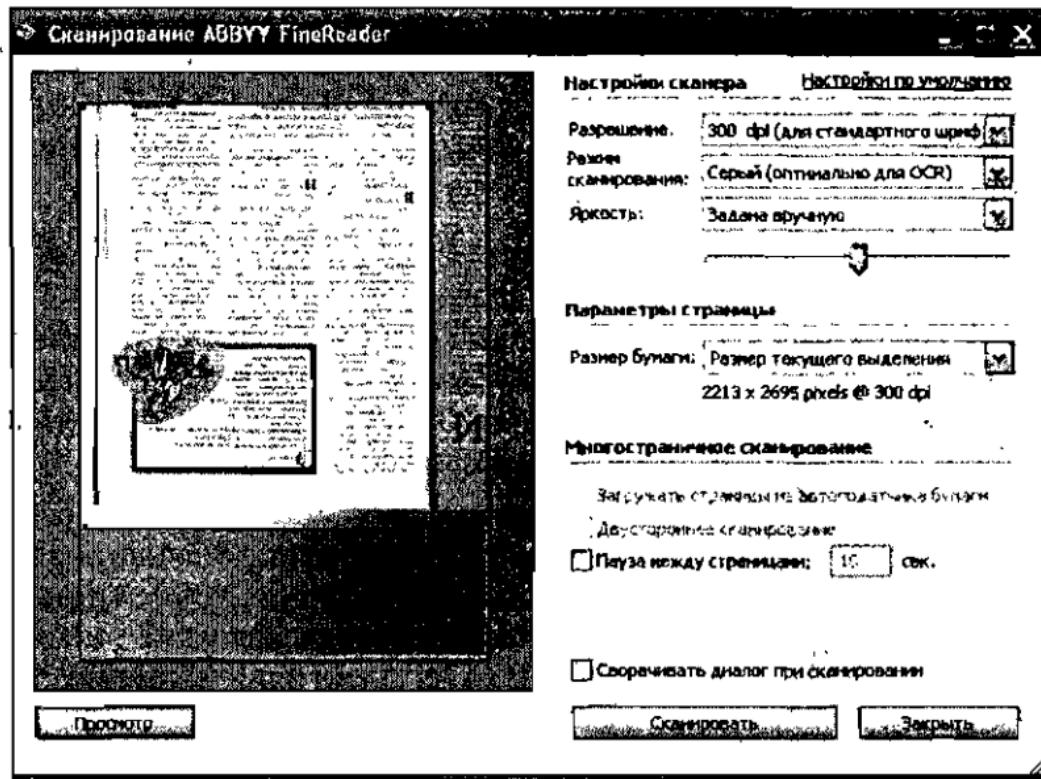


Рис. 2.5. Диалоговое окно «Сканирование ABBYY FineReader»

Ползунком **Яркость** задается степень яркости бумаги. В большинстве случаев оптимальным значением яркости является 50%. Однако если при сканировании буквы получаются «ранеными», излишне светлыми и тонкими, то следует уменьшить яркость. Если же в результате сканирования получились слишком жирные, смазанные буквы со сливающимися краями, то это яркость необходимо увеличить.

Для начала сканирования щелкните кнопку **Сканировать** (рис. 2.5). Таким же образом вы можете отсканировать еще одну или несколько страниц. По окончании сканирования нажмите кнопку **Закрыть** в окне **Сканирование ABBYY FineReader**. Как только документ отсканируется, автоматически запустится процедура распознавания текста, по завершении которой откроется редактор Word. В его окне отобразится уже готовый для дальнейшего редактирования документ. Вот, собственно, и все. Но это если все распозналось корректно, что далеко не всегда бывает. Также необходимо отметить, что по умолчанию рисунки с исходной страницы не переносятся. Только текст. А вполне вероятно вы захотите включить и рисунки в Word-документ.

О том, как настроить параметры распознавания применительно практически ко всем режимам, сказано в следующей главе.

Настройка распознавания текста, задание областей распознавания, полный контроль

Чтобы вручную настроить параметры и области распознавания, следует после сканирования вернуться в окно FineReader, которое будет иметь вид, показанный на рис. 3.1.

Сразу вас предупреждаем, что описываемое нами расположение окон является таковым по умолчанию. В дальнейшем его можно изменять как угодно по своему усмотрению, с помощью пункта **Вид** главного меню программы, или же просто перемещать границы окон с помощью мыши.

Главное окно программы состоит из нескольких подокон (для удобства мы будем называть их окна). В верхней части расположено главное меню. Каждое окно имеет свой порядковый номер, кроме самого нижнего. Там всегда будет отображаться оригинал отсканированного документа. В любой момент вы сможете сравнить определенную область оригинала с

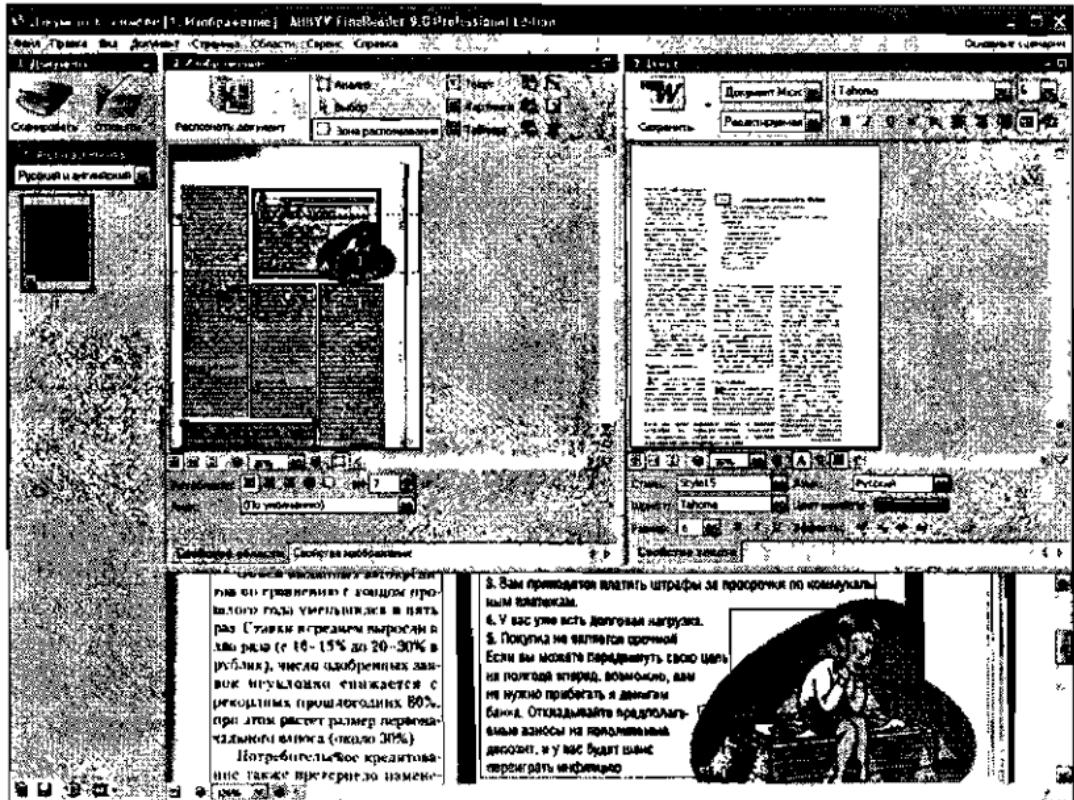


Рис. 3.1. Основное окно программы FineReader

22 той же областью, но уже распознанного текста. Как бы ни был совершен алгоритм оптического распознавания, в любой момент может произойти ошибка. Любая погрешность сканера может неправильно распознать тот или иной символ. Например, английская буква **I** зачастую интерпретируется программой как цифра **1**, буква **O** как цифра **0** и так далее. И вы всегда сможете «на глаз» сравнить корректность распознанной области документа.

Окно номер 1 служит для общего управления документами. Здесь вы можете как отсканировать бумажный оригинал, так и открыть готовый файл уже отсканированного документа, файла изображения либо файла PDF-типа.

Если документ содержит текст иного, чем установленный по умолчанию, языка алфавита, то для корректного его распознавания выберите нужный из списка.

Ниже будут отображаться мини-изображения всех страниц документа. Щелкнув правой кнопкой нужную страницу, вы вызовете контекстное меню, с помощью которого вы сможете выполнять различные операции с ней, например, повернуть её (рис. 3.2).

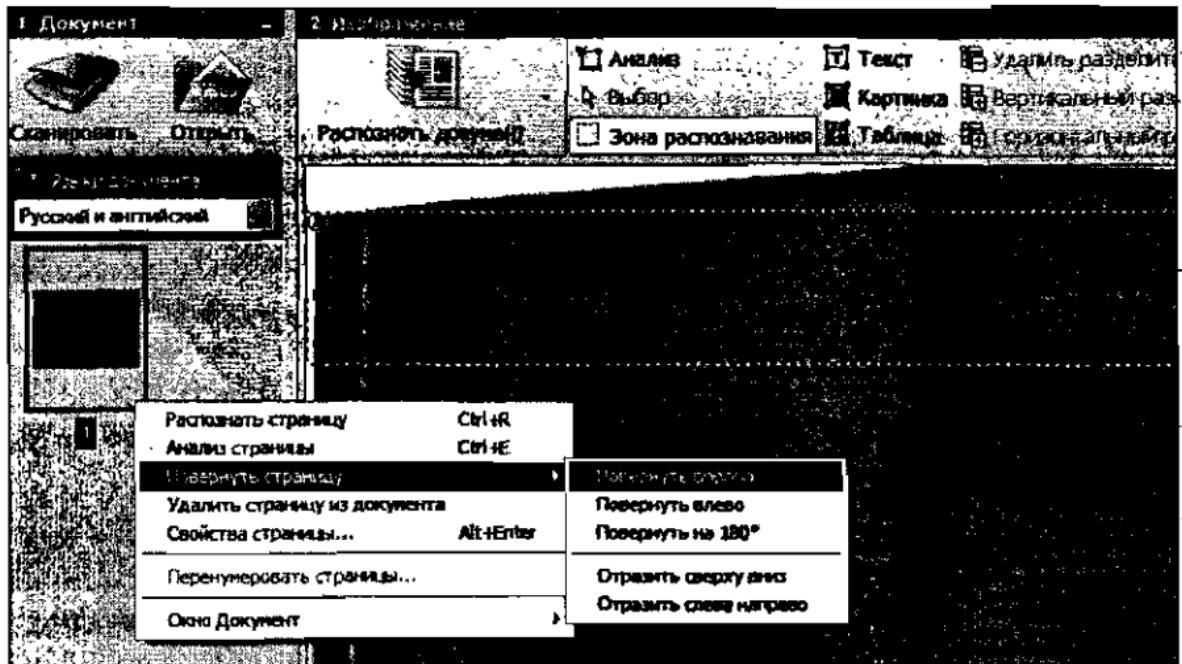


Рис. 3.2. Контекстное меню мини-изображения

24

Второе окно предназначено для распознавания документа. Для выполнения данной операции в автоматическом режиме щелкните кнопку **Распознать документ**. В этом случае программа сама определит области и проведет процедуру распознавания. Но иногда нам требуется распознать лишь определенную часть документа. Или, например, когда в документе вместе с текстом отображен рекламный баннер, также включающий в себя текст. В результате распознавания программа исключит этот текст из баннера, оставив лишь картинку. Но ведь вам этого не нужно. Баннер должен целиком быть сохранен, иначе нарушится его целевое предназначение. Поэтому нужно будет вручную указать область картинки, этим самым дав программе понять, что текст, вошедший в эту область, не нужно распознавать.

Рассмотрим инструменты, присутствующие в данном окне.

Анализ – выполнение автоматического определения областей документа. В результате документ будет разделен на области для дальнейшего распознавания. Смысл областей в том, что текст должен быть отделен от картинок. Текстовые области будут отмечены зеленым цветом, а области изображений – красным (рис. 3.3).

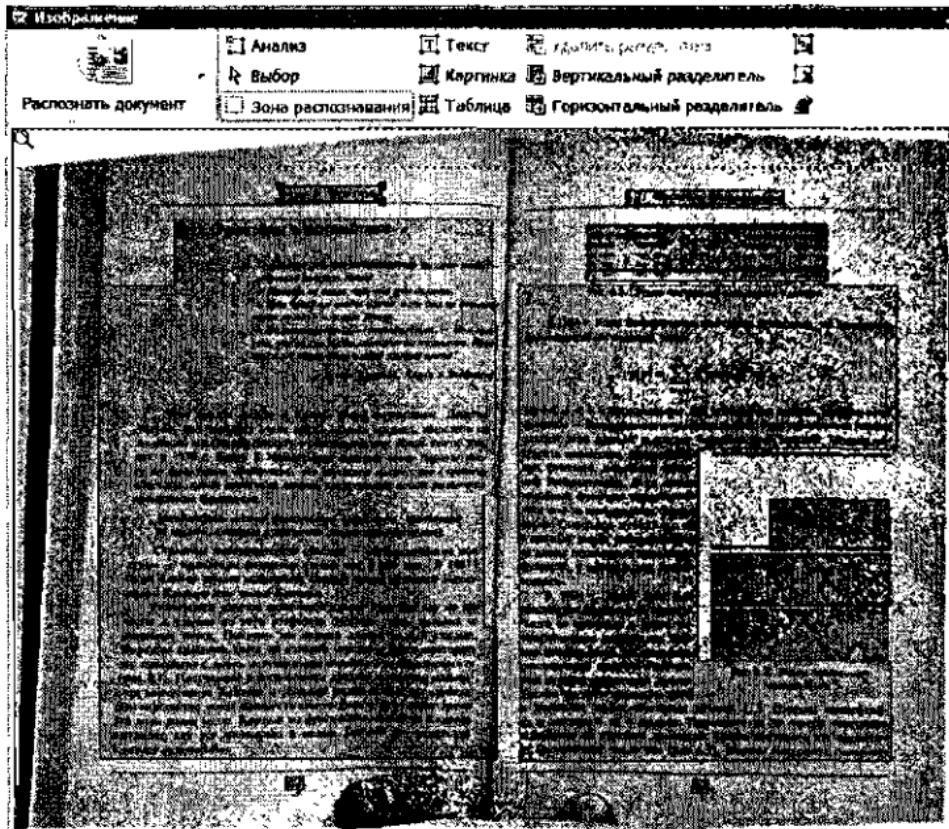


Рис. 3.3. Выполненный анализ страницы

26

Теперь вы можете вручную подправить любую из областей, путем перемещения их полностью, либо перетаскивая границы с помощью мыши. Если какая-то из областей не требуется для перемещения в выходной документ, здесь же её можно удалить.

Любую операцию, выполненную в программе, можно отменить с помощью меню Правка – Отменить или комбинации клавиш **Ctrl+Z**.

Обратите внимание, что каждая область имеет свой порядковый номер. В этой очередности страница будет распознаваться. Если вас чем-то не устраивает порядок, то вы можете перенумеровать области с помощью предназначеннной для этого кнопки.

При желании вы можете не делать автоматического анализа, вручную распределить области. Для этого воспользуйтесь кнопками Текст, Картинка или Таблица соответственно для выделения области для распознавания текста, области для изображений, а также для таблиц. Этим вы значительно повысите качество распознавания документа, так как алгоритм распознавания таблиц немного отличается от текста.

Кстати, размечать и редактировать области вы можете даже в нижней

области окна программы, где находится изображение документа в 27 натуральную величину.

В нашем примере произошел случай некорректного определения области изображения. Программа приняла за текст некоторые символы,

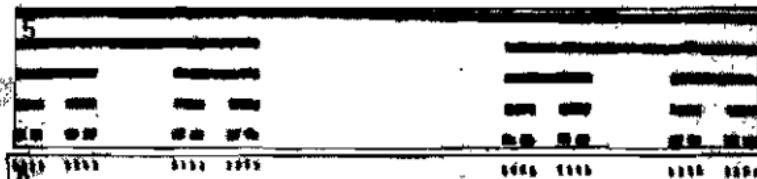


Рис. 6.1. Этапы построения множества Кантора

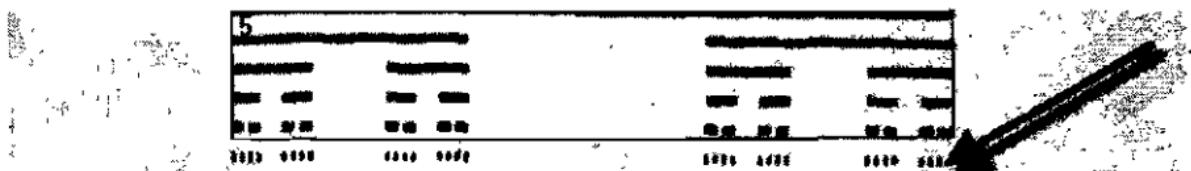
В 1886 г. Карл Вейерштрасс построил функцию, не имеющую производной ни в одной точке

$$y(x) = \sum_{n=0}^{\infty} A^n \cos(B^n \pi x),$$

Рис. 3.4. Неправильное определение области изображения

23 изображенные на странице. Кроме этого, в текстовую область вошла и формула, которая также не будет подлежать дальнейшему редактированию (рис. 3.4). Поэтому нам нужно вручную внести изменения.

Сначала удалите текстовую область, в которую ошибочно попал рисунок. Для этого щелкните кнопку Выбор и затем щелкните ошибочную зеленую область с номером семь. Теперь просто нажмите клавишу **Delete** (рис. 3.5).



6

Рис. 6.1. Этапы построения множества Кантора

В 1886 г. Карл Вейерштрасс построил функцию, не имеющую производной ни в одной точке

$$y(x) = \sum_{n=0}^{\infty} A^n \cos(B^n \pi x),$$

Рис. 3.5. Удаление ненужной области

Стрелкой показано исчезновение некорректной области. Теперь зацепитесь за нижнюю часть области изображения №5 и потяните вниз, чтобы «бездомный» рисунок вошел в неё (рис. 3.6).

Теперь займемся формулой. Рассмотренный способ включения её в соседнюю область изображения нам не подойдет. Здесь нужно создавать



Рис. 3.1. Этапы построения множества Кантора

в 1861 г. Карл Вейерштрасс построил функцию, не имеющую производной ни в одной точке

$$y(x) = \sum_{n=0}^{\infty} A^n \cos(B^n \pi x)$$

Рис. 3.6. Включение рисунка в соседнюю область

30 новую область картинки. Щелкните кнопку Картинка и обведите формулу, тем самым создав новую область. Номер будет присвоен автоматически (в нашем случае – 11) (рис. 3.7).

5 Рис. 6.1. Этапы построения множества Кантора

В 1886 г. Карл Вейерштрасс построил функцию, не имеющую производной ни в одной точке

11

$$y(x) = \sum_{n=0}^{\infty} A^n \cos(B^n \pi x),$$

где $0 < A < 1$, а произведение AB достаточно велико. График этой функции – бесконечно изломанная линия. При увеличении любой участок этой кривой выглядит подобно всей кривой. Можно построить множество разнообразных функций, подобных функции Вейерштрасса. В то время подобные функции представлялись чем-то аномальным. Эрмит писал Стилтьесу в 1883 г.: «Я с дрожью ужаса отворачиваюсь от ваших несчастных проклятых

7



8

Рис. 3.7.
Создание
области
картинки

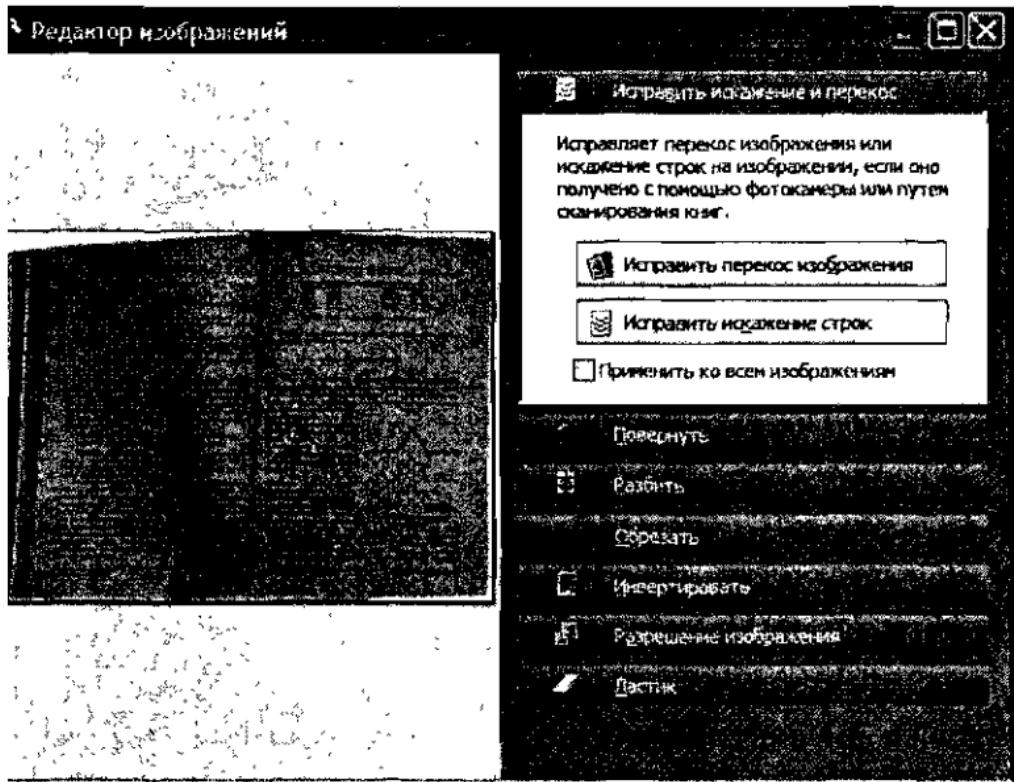


Рис. 3.8. Встроенный редактор изображений

Теперь уже не составит труда исправить ниже еще одну красную область (№7).

Вот, собственно, и все премудрости. Здесь можно еще остановиться на кнопке **Редактировать изображение страницы**. Она служит для перехода в режим редактирования страницы после сканирования (рис. 3.8).

Встроенный редактор, конечно, далек от всяких Фотошопов, однако он и не предназначен для глубокой реанимации изображения, к тому же всегда под рукой. Его вполне должно хватить для подготовки страницы к распознаванию. После сделанных изменений просто щелкните значок закрытия окна (красный крестик).

Теперь все готово для выполнения распознавания страницы, поэтому остается лишь щелкнуть кнопку **Распознать документ**. Начнется процесс распознавания, итогом которого будет появление текста в третьем окне программы (рис. 3.9).

В верхней части окна расположены кнопки управления распознанным текстом, начиная от проверки орфографии с последующим исправлением ошибок и заканчивая конвертацией его в выбранный формат.

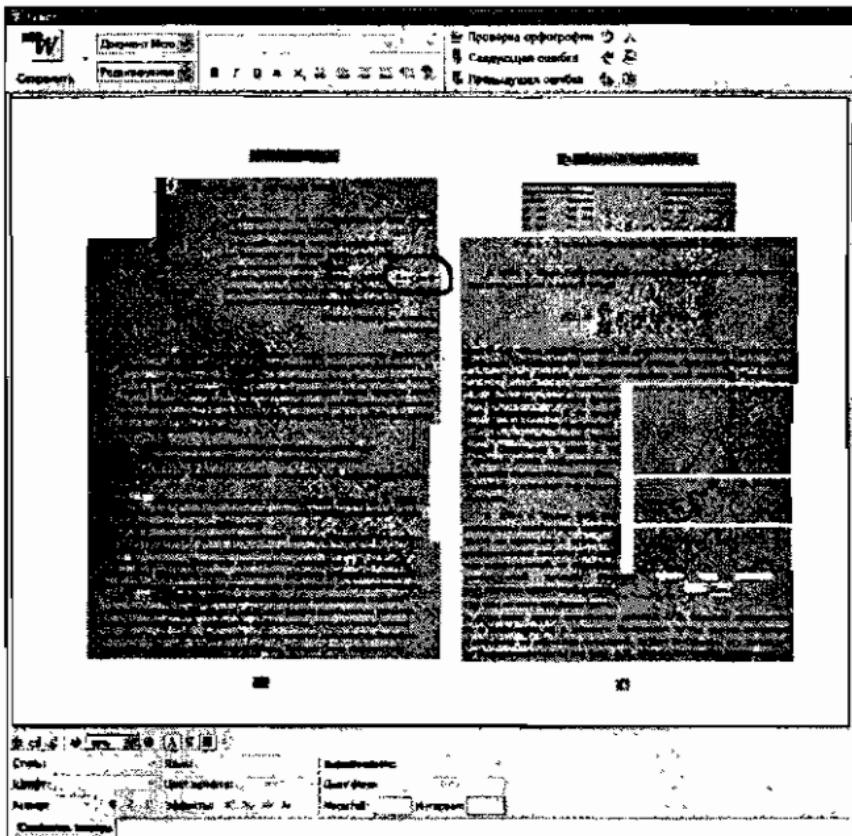


Рис. 3.9. Распознанный документ в окне Текст

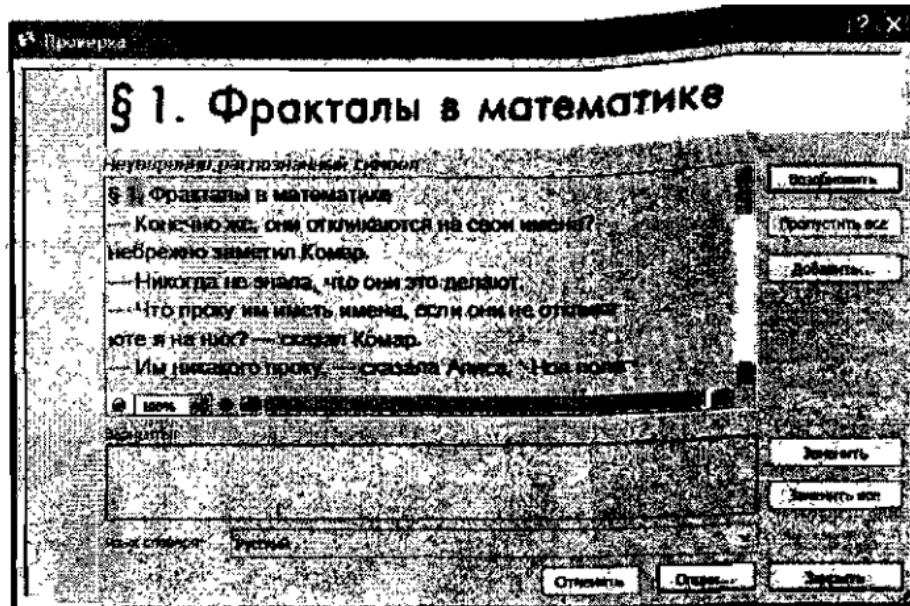


Рис. 3.10. Проверка орфографии

Начнем с орфографии. В случае невозможности идеального сканирования бумажного документа вы можете столкнуться с появлением ошибок в распознанном документе. Если вы когда-нибудь сталкивались со сканированием толстой книги, то знаете, что в области перегиба книги

получаются искривления текста, а то и нечитаемые потемнения (зависит от сканера). Вот и в нашем примере произошел такой же случай (выделено красным цветом на рис. 3.9). В принципе можно все оставить так, как есть, но все-таки лучше исправить ошибки прямо здесь на месте, тем более что программой созданы для этого самые лучшие условия. Зеленым цветом в тексте выделены слова, в которых программа сомневается, что они распознались в соответствии с орфографическими правилами. Красным цветом подчеркнуты слова, которые Fine Reader не знает, то есть их просто нет в базе данных программы.

Щелкните кнопку **Проверка орфографии**. Программа начнет поиск отклонений от правил. В случае положительного результата откроется окно **Проверка**, в котором будет отображен текст с выделенным словом, в котором программа усматрела отклонение от правил орфографии (рис. 3.10).

Проверьте идентичность подчеркнутого слова с оригиналом в нижней части окна программы. Если слова совпадают, то продолжите проверку, щелкнув кнопку **Пропустить**. Программа ищет следующее неуверенно распознанное слово. Вы снова сверяетесь с оригиналом и продолжаете проверку. Так происходит до тех пор, пока не попадется слово, отличающееся от бумажного экземпляра (рис. 3.11).

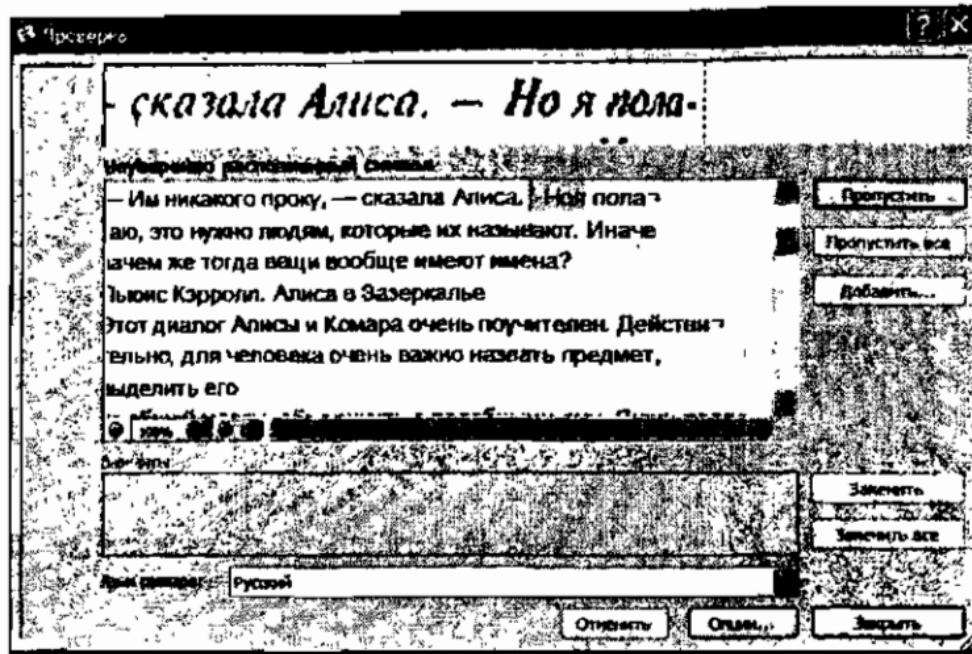


Рис. 3.11. Неправильно распознанное слово

В данном случае вы можете просто вручную ввести правильное слово. Когда в нижней области окна Проверка будут выдаваться варианты слова, из которых, по мнению программы, должен быть правильным. Вы

можете просто щелкнуть его, и слово заменит выделенный синим цветом фрагмент в тексте (рис. 3.12).

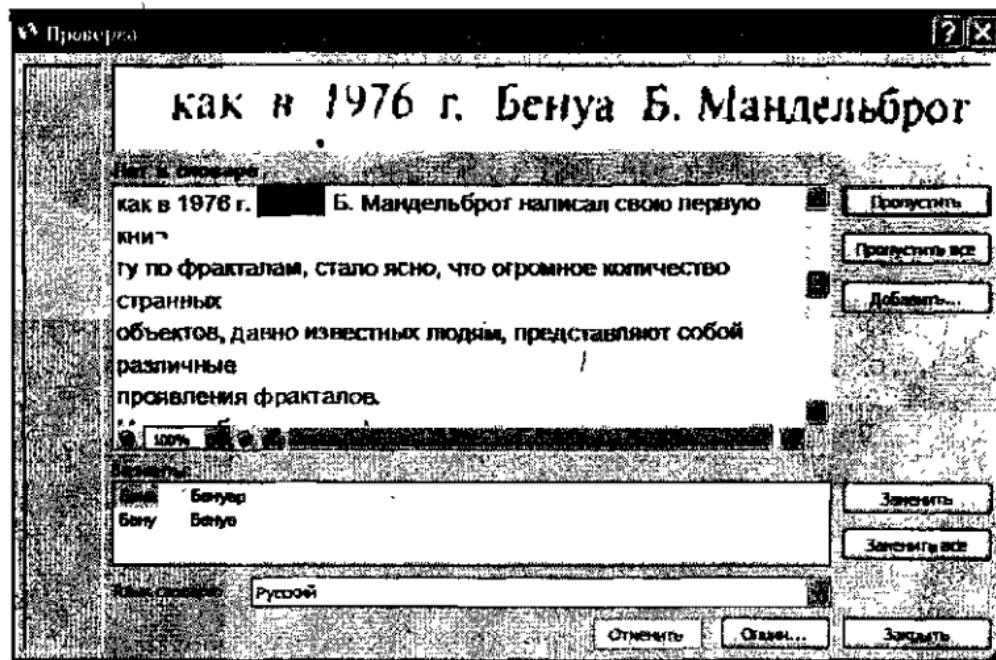


Рис. 3.12. Выбор правильного варианта

733. Но в нашем примере в области **Варианты** нет ни одного правильного, поэтому все равно придется ввести вручную.

После того как вы исправите ошибочный фрагмент, щелкните кнопку **Подтвердить**.

Таким образом вы должны дойти до конца документа. В итоге вы получите совершенно безошибочный (с точки зрения орфографии) вариант текста.

Теперь дело осталось за малым – преобразовать его в нужный формат. Кстати, если вы хотите сохранить текст, например, в PDF-формате, то можно без помощи сторонних приложений выполнить его форматирование (ведь в PDF нет возможности редактирования документа). Для этого в нижней части окна **Текст** имеются самые необходимые инструменты. Конечно, это не Word, но все же и здесь можно привести текст в приличный вид. Для этого просто выделите необходимый фрагмент и с помощью инструментов отредактируйте его на свой вкус.

С помощью раскрывающегося списка выберите нужный формат документа, в который вы желаете преобразовать полученный текст (рис. 3.13).

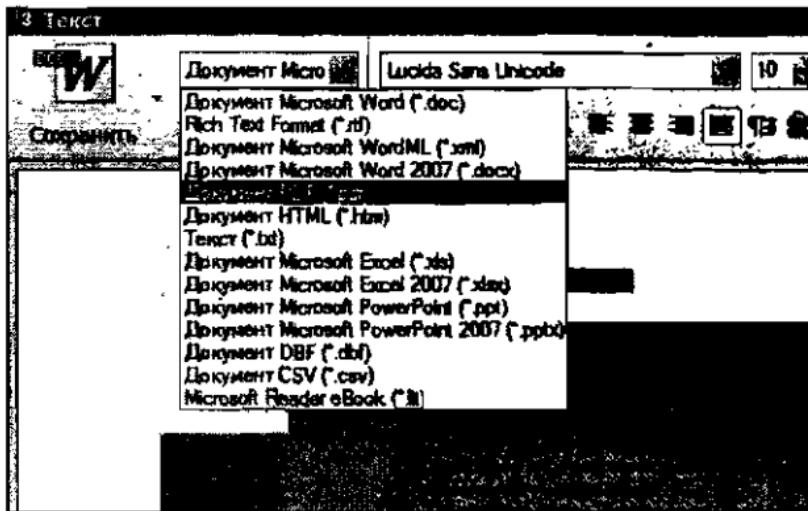


Рис. 3.13. Доступные типы файлов для сохранения

Как видите на рисунке, для сохранения текста доступны практически все самые распространенные типы файлов, поэтому вы никогда не будете иметь проблем.

При выборе того или иного типа выходного файла ниже в раскрывающемся списке можно выбрать дополнительные опции для его создания.

40 Например, если вы будете создавать документ Microsoft Word, то можно будет выбрать **Точную копию**, **Редактируемую копию**, **Форматированный текст** и **Простой текст**.

При выборе точной копии в выходном файле вы получите точное расположение текста и картинок, как оно было на бумажном носителе. Однако редактировать в таком виде не совсем удобно, так как текст и картинки помещены в ограниченные области, размеченные во время анализа страницы перед распознаванием текста. Обычно применяется в случаях, когда актуально именно сохранение оригинального макета (рис. 3.14).

Прямой противоположностью точной копии является редактируемая копия, которая отличается от только что рассмотренной тем, что не имеет ограниченных полей. В таком виде становится удобным редактирование документа.

Но бывают ситуации, когда вам совершенно не нужно сохранение макета страницы. Вас интересует только содержание (текст и картинки), которое в дальнейшем вы будете использовать по своему усмотрению. Для этого существуют опции форматированного и простого текста. Отличаются они только шрифтом текста в выходном файле. Если в форматированном

6.1. Фрактальность в математике

—Конечно же все определяется на самом деле!—
исрелически воскликнул Бенор.
—Но ведь вы же не можете подумать, что я...
—Что придумать? —спросил Бенор.— Но я подумал, что это было бы интересно, и поэтому я попытался... Ну да, я знаю, что это было бы интересно...

Льюис Кэрролл. Алиса в Зазеркалье

Этот диалог Алисы и Кантора очень поучителен! Действительно, для членов этого семейства предметы, выделяющиеся из общего массы, обладают суперпозиционными свойствами математики 1976-го (Бенору Б. Кантором) раскрыты тема «переворачивающие фракталы», стала ясна, что первоначальное множество становится «бессмыслицей»: двумя наименее логичными представлениями собственной реальности!

Конфузорные функции или перекруты, скручивающие, изогнувшиеся и смешанные подобиямоку?

Стремящиеся к математической объективности — математические «студии» — виноваты в возникновении в XIX веке... В 1863 г. Готфр. Кантор (именно: математик, автором труда называемого математиками Кантором или письмом Кантора) у

расмотрел некоторое единство единиц. Рассматривая различные виды изображения градиентных структур, он вывел концепцию тона. Получив двойную длину 1/2 единиц, выраженную седьмую треть, канторидного отрезка и будучи повторять эту процедуру, он сошел с пути чистых математиков до физики! (рис. 6.11). Далее он вывел математическое значение множества, связанных с кантором: Кантор и обладает рядом необычайств. Это неслыханное «единство» множества, это неслыханное время «степени времени» нулю. Каждый из фракталов изображает Кантора-антагонист, как и множество в целом. Говорят, что множества Кантора-антагонисты?



Рис. 6.1. Этапы построения множества Кантора

1

В 1866 г. Карл Вейерштрасс (история физики: «изящная производная») показал, что

где $0 < A < 1$, и приводит к 48 достаточно членам. График этой функции — бесконечное множество линий. При увеличении любой участок этой кривой поддается описанию в Аксиомах Г

стремящимся к равнобедренным фракталам, подобным фракталам Вейерштрасса. В это время подобные функции представлялись как аналитические. Эрнст Кантор (математик) в 1883 г. в «Словаре математических терминов» включил в себя «функции Кантора», у которых нет производных.

В 1884 г. Харди (академик Кембриджского университета) ввел в курс математического анализа матрицы квадратной, квадратичной, квадратурной, кривой. Каждая из которых единичной длины, удовлетворяет условию третьей и дополнительной единицы длины 1/3. Следует обратить внимание на концепцию «Фракталы». Применив эту концепцию в Канторе, можно получить изображение Кантора-антагониста. Будет повторять эту процедуру бесконечное число раз (рис. 6.2). Многопостроитель Кантора?



тексте шрифт останется неотличимым от оригинала, то во втором случае вы получите текст с шрифтом, который установлен по умолчанию в настройках программы (рис. 3.15).

Глава 3. Фракталы

§ 1. «Фракталь в математике»

— Конечно, он интересен не потому, что я избрался захотел Кантор? — Айзенштадт, спросил ее Альберт
— Что-то учительшина? — сказал один из учеников — сказал Кантор.
— Неизвестно, — сказал Альберт. — Но я избрался это, чтобы говорить
Львову Карлову, Альбера Зандерса!
Этот диалог Альбера очень получился, — добавил Гельфанд, — я люблю математику, я кинь пасхальное, если в 1976 г. Петру Б. Майдельбрехт объектов, даже не имеющих размеров, представляют собой фракталы.
Мой королевский профессор, старший, состоящий из четырех категорий.
Странные математические объекты — математические конструкции, которые строятся из множества компонентов Кантора или Льва Кантора.
Рассмотрим их на единичной единице. Равнодействующая части отрезка длиной π из $1/3$ каждого. Вырежем среднюю третью из каждой бесконечности (рис. 6.1). Получим бесконечное множество конечных, имеет конечность бесконечную, а также самое первое отрезок и говорит, что множество Кантора скончалось!

807

§ 1. «Фракталы в математике»



Рис. 6.1. Открытие настройки Кантора

В 1884 г. Карл Вейберштрасс построил функцию, известную в

математике как «функция Кантора».

Если $f(x) = 1$, то производная AB достаточно велика. График этой функции — бесконечный спираль, выходит из одиночной кривой. Можно построить множество разнообразных

Глава 3. Фракталы

§ 1. «Фракталы в математике»

— Конечно же, она откладывает на свою книгу? побежал заметил Кантор?

— Никогда не читала, что они это делают?

— Чуть-чуть имена, есть они постепенно, когда я начну? — сказал Кантор?

— Ничемного прошу, — сказала Альба. — Но я подожду, это нужно людям, которые начнут

Львову Карлову, Альбера Зандерса!

Этот диалог Альбера очень получился. Добавьте тему для человека очень в

интересных ему Львову Карлову, если в 1976 г. Балуу Б. Майдельбрехт национализированы объекты, давшие название людям, представляют собой различные приемы

Майдельбрехт-из-за фракталов и структур, состоящих из частей, которые в каком-то

Странные математические объекты — математические конструкции — привлекли ваши

интересство, которое теперь называют множествами Кантора или вымысли Кантора?

Рассмотрим их на единичной единице. Равнодействующая части отрезка из $1/3$ отрезка отрезка длиной π из $1/3$ каждого. Вырежем среднюю третью из каждого бесконечности (рис. 6.1). Получим бесконечное множество конечных, что скажет множеству конечность бесконечную, в то же самое время что и настройки, а также самое первое отрезок и говорит, что множество Кантора скончалось.

807

§ 1. «Фракталы в математике»

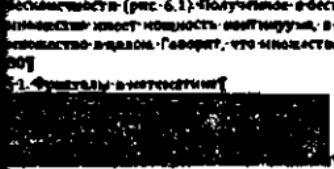


Рис. 6.1. Открытие настройки Кантора

В 1884 г. Карл Вейберштрасс построил функцию, не имеющую производной ни в одной-

точке, кроме конечной.

Если $0 < A < 1$, то производная AB достаточно велика. График этой функции — бесконечный спираль, выходит из одиночной кривой. Можно построить множество разнообразных

Рис. 3.15. Выходной файл Word с форматированным (слева) и простым (справа) текстом

Таким образом, вы можете сохранять текст в различных форматах, используя нужные опции. Приведем примеры некоторых выходных форматов нашего текста (рис. 3.16 ... 3.18).

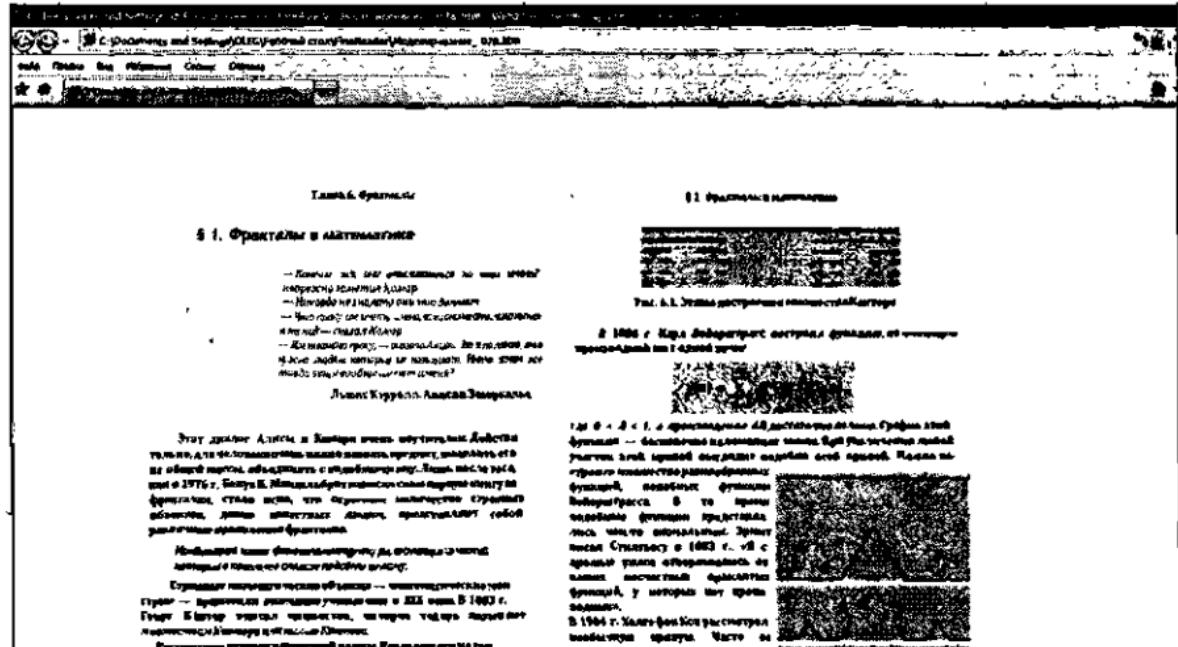


Рис. 3.16. HTML-страница

Глава 6. Фракталы

§ 1. Фракталы в математике

§ 1. Фракталы в математике

— Конечно же они определяют не своим цветом?
нарвало Камор.

— Наверно не знала, что они это делают.

— Что скажут ими цветом, если они не охлаждаются в море? — сказал Камор.

— Им придется прокур, — сказала Алиса. — Но я помоложе, привычней к инженерной математике. Кто же заслужил бы еще краевого шинца?

Лицес: Королева Алиса в Зазеркалье

Этот диалог Алисы и Камира очень заучительный. Действительно, для человека очень важно видеть предмет, выделяющийся из общей массы, обединять с подобными ему. Лицес после того, как в 1976 г. Бенуа Б. Мандельброт написал свою первую книгу по фракталам, стало ясно, что огромное количество странных объектов, даже непонятных людям, представляют собой различные проявления фракталов.

Мандельброт называл фракталами структуры, состоящие из частей, которые в конечном смысле подобны целому.

Строительные математические объекты — «математические конструки» — привлекли внимание ученых еще в XIX веке. В 1883 г. Георг Кантор описал множество, способное покрывать множеством Кантора.

Рассмотрим отрезок единичной длины. Разделим его на три части и удалим из него среднюю третью, оставив ее концевые точки. Получим два отрезка длиной по $1/3$ каждой. Выполним среднюю третью из каждого отрезка и будем повторять эту процедуру с новыми получившимися отрезками до бесконечности (рис. 6.1). Полученные с бесконечным числом множеств называются множествами Кантора и обладают рядом любопытных свойств. Это множество имеет плотность континуума, а также такие пренебрежимо малые измерения, что

так $0 < A < 1$, в произведение AB достаточно велико, функции — бесконечно маленькие лягушки. Рук ушия участок этой кривой выходит за пределы всей кривизны спектра множества разрывобанных функций, подобных функциям Мандельброта. В то время подобные функции представлялись чисто эпиграфическими. Эрик Стенлессу в МИИТ гг. «Н с доказал удалось отыскать в них настоящими предметами фракталами, у которых нет пренебрежимо малых измерений».

В 1904 г. Хелле фон Кох разработал необычную кривую. Часть ее входит в курс математического анализа как пример непрерывной, но нондифференцируемой кривой. Рассмотрим отрезок единичной длины. Удалим из него среднюю третью и добавим две отрезки длиной $1/3$. Отрезки

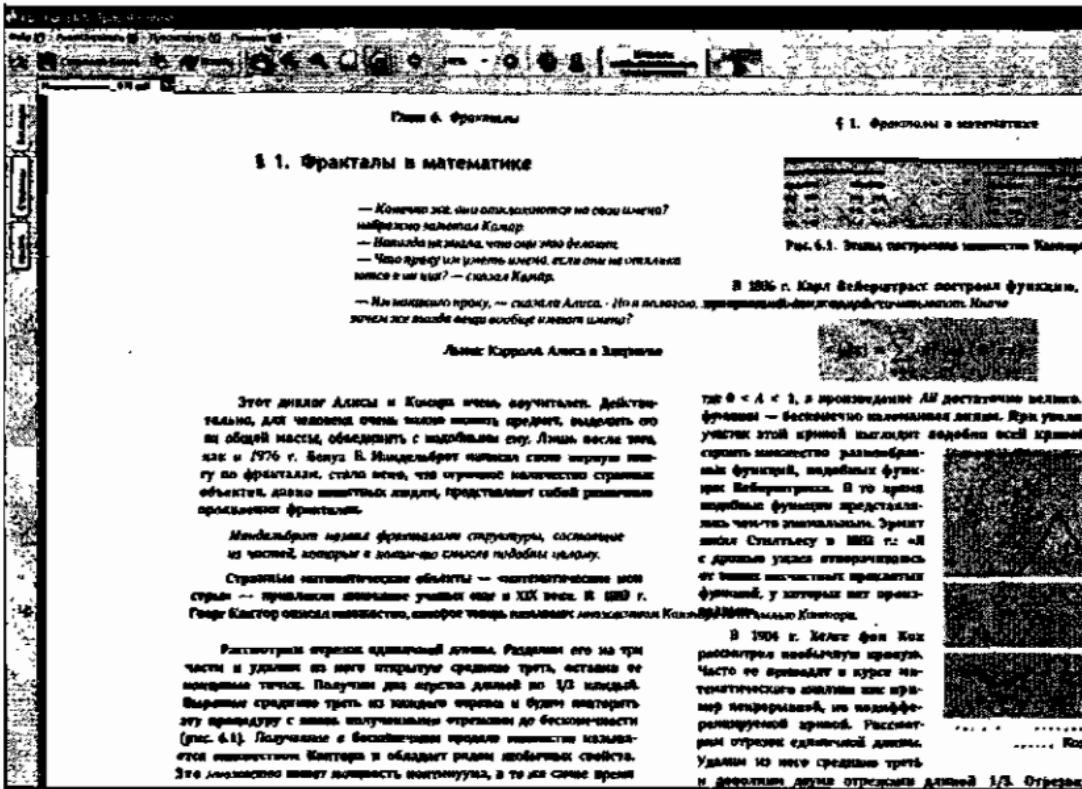


Рис. 3.17. PDF-документ

Специальные режимы распознавания в FineReader

Как отсканировать и преобразовать в PDF

Если вы желаете отсканировать микст (текст с картинками), то очень удобно преобразовать его в формат PDF. Запустите соответствующий сценарий, предварительно поместив лист бумаги в сканер. Далее снова выделите нужную область для сканирования. Вот теперь может потребоваться поменять опции сканирования в случае, если вам необходимо получить цветной документ. Дело в том, что по умолчанию программой установлен режим сканирования в градации серого (для ускорения процесса). Поэтому поменяйте данную опцию, выбрав в раскрывающемся списке пункт Цветной. При выборе величины разрешения следует руководствоваться общим принципом. Для простого распознавания текста достаточно 300 dpi (меньше не рекомендуется), а вот для того, чтобы в результате преобразования получить качественное изображение, нужно установить как

46 · минимум 600 dpi.

Распознавание таблиц и сохранение их в формате Excel

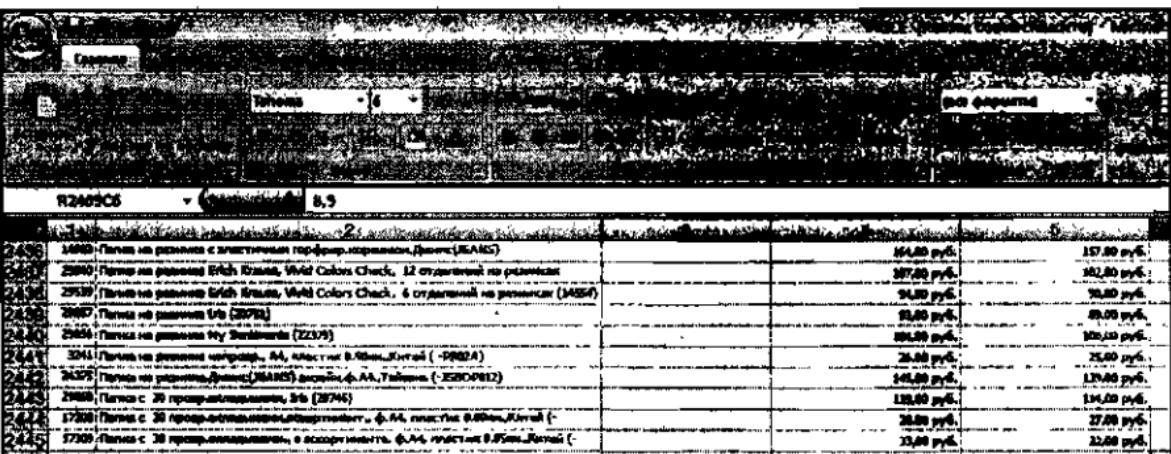
Если вам понадобится распознать таблицу, то нужно поступить немного другим образом. После сканирования таблицы она отобразится в окне **Изображение** (рис. 4.1).



Рис. 4.1. Отсканированная таблица

Как вы можете заметить по картинке, в данном случае выделенная область окрашена в синий цвет. Это обозначает, что анализируемым объектом является таблица. И если вам понадобится вручную определить нужную область, то для этого нужно применить кнопку Таблица, очертив ею границы таблицы.

После процедуры распознавания текста проверка орфографии проводится в штатном режиме и ничем не отличается от проверки обычного



The screenshot shows a Microsoft Word document with a table selected. The table has 13 columns and 14 rows. The first column contains product codes like 'R240505' and '240505'. The second column contains descriptions of the products. The third column contains quantities ('1', '2', etc.). The fourth column contains unit prices ('157,00 руб.', '162,00 руб.', etc.). The fifth column contains total prices ('157,00 руб.', '162,00 руб.', etc.). The table is styled with alternating row colors (white and light gray). The entire table is highlighted with a thick blue border.

R240505	240505	8,9		
14001	Папка на резинке с эластичным поролон.корешком, фане (БАБС)	1	157,00 руб.	157,00 руб.
240505	Папка на резинке Etch blouse, Metal Colors Check, 12 страницей на резинке	2	162,00 руб.	324,00 руб.
25037	Папка на резинке Etch blouse, Metal Colors Check, 6 страницей на резинке (БАБС)	1	162,00 руб.	162,00 руб.
24057	Папка на резинке Etch (БАБС)	1	162,00 руб.	162,00 руб.
24054	Папка на резинке Ну Эмблема (223/9)	1	162,00 руб.	162,00 руб.
24141	Папка на резинке матовая, А4, цветы 0,4мм/бумага (-29821)	1	20,00 руб.	20,00 руб.
24025	Папка на резинке, бумага (БАБС) акрил. ф. А4, Тайвань (-25009812)	1	140,00 руб.	140,00 руб.
24060	Папка с 20 прозрачными листами, ф. А4, пластик 0,4мм, Китай (-28746)	1	120,00 руб.	120,00 руб.
27208	Папка с 30 прозрачными листами, ф. А4, пластик 0,4мм, Китай (-28745)	1	20,00 руб.	20,00 руб.
27209	Папка с 30 прозрачными листами, в ассортимент, ф. А4, пластик 0,4мм, Китай (-28745)	1	22,00 руб.	22,00 руб.

Рис. 4.2. Сохранение таблицы в формате Microsoft Excel

48 текста.

При выборе выходного файла для таблицы вы выбираете любой формат, но естественнее было бы, конечно, сохранить её в формате Microsoft Excel, что мы и проделаем. В результате получаем свежеиспеченный XLS-файл (рис. 4.2).

Сканирование изображений (без текста) в Word

Если вы готовите документ Word, содержащий большое количество картинок, то с помощью FineReader вы сможете легко осуществить такую задачу. Сначала вы просто сканируете нужные вам изображения, при этом сохраняя файл в одном из множества предлагаемых форматов. Для этого воспользуйтесь сценарием **Сканировать и сохранить изображение**. Когда все картинки будут оцифрованы, можно приступить к выполнению сценария.

Выберите пункт **Сервис – Сценарии ABBYY FineReader - Конвертировать фотографию в Microsoft Word**. Откроется окно, в котором нужно

выбрать файлы изображений, которые вы хотите включить в выходной документ (рис. 4.3).

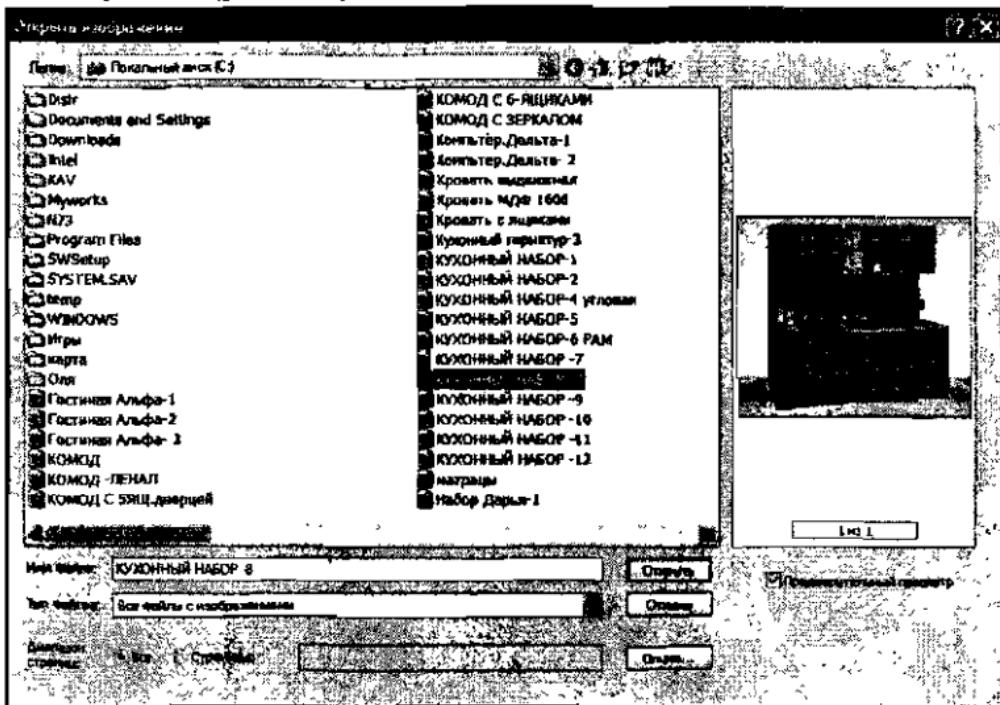


Рис. 4.3. Выбор файлов для конвертирования в Word

50

Как видно по рисунку, любое изображение можно на месте просмотреть. Остается лишь щелкнуть кнопку **Открыть**, и вы сразу же получите новый документ Word с уже включёнными в него изображениями. Не правда, как удобно! (рис. 4.4).

Теперь вы можете перемещать её в любое место с помощью инструментов Word.

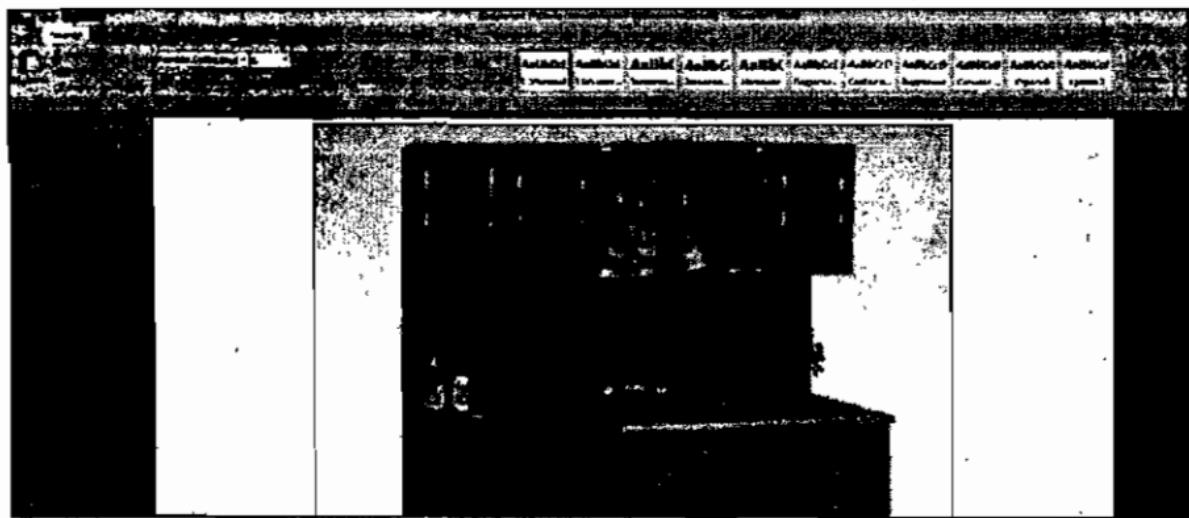


Рис. 4.4. Документ Word с конвертированной фотографией

Как преобразовать PDF-файл в Word'овский текстовый файл

51

Кто когда-нибудь пользовался PDF-файлами, знает, что без специальных версий Acrobat Reader невозможно редактировать данный документ. Хотя первоначально он и был придуман чисто для удобного чтения документов и другого. Но с помощью FineReader вы можете конвертировать PDF-файл в DOC-файл. Кстати, как FineReader, так и Acrobat Reader являются продуктами одной компании ABBYY.

Для конвертации PDF-файла запустите сценарий **PDF или изображения в Microsoft Word**. В открывшемся окне выберите нужный файл и щелкните кнопку **Открыть**. Кстати, если документ имеет большое количество страниц, вы можете выбрать только нужный диапазон.

Как и в предыдущем примере, вы сразу же получите новый документ Word с включенными туда страницами PDF-файла. Причем теперь вы сможете запросто редактировать его. Если будут обнаружены ошибки распознавания, можно вернуться в окно FineReader'a, установить в нем нужные параметры и задать вручную границы областей распознавания, а затем заново произвести распознавание.

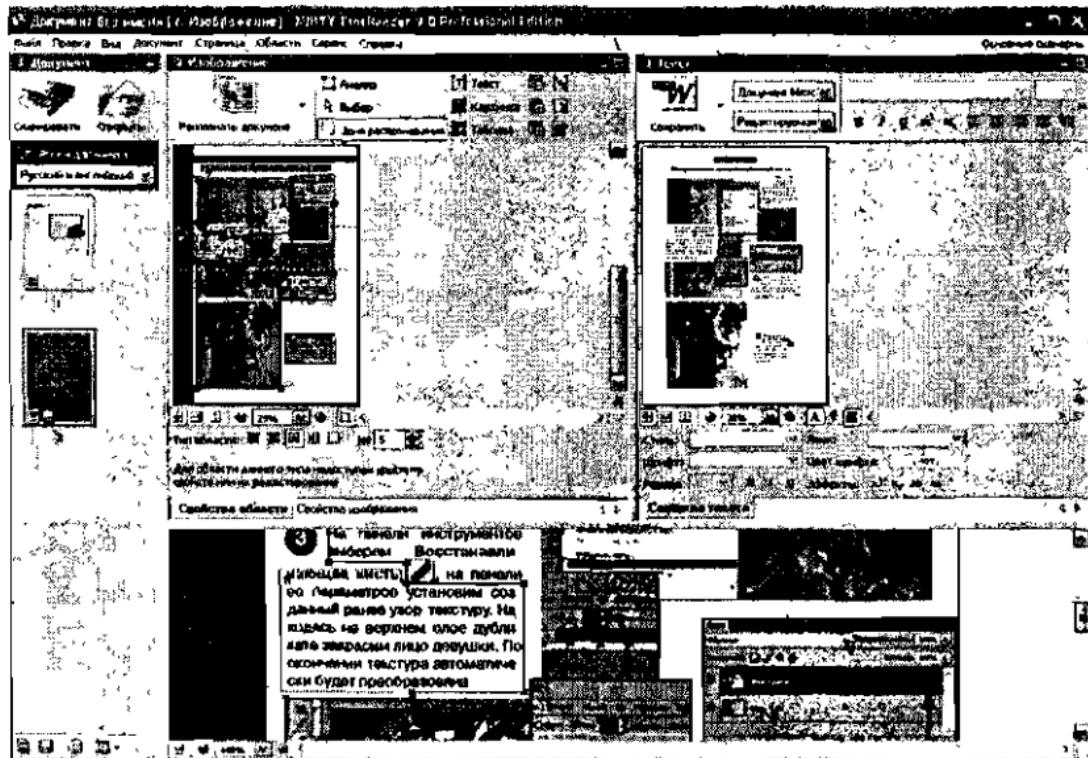


Рис. 4.5. Преобразование PDF-файла в текстовый файл

Повышение качества распознавания текста

Режим обучения при распознавании текста

Если вы решили распознать документ, содержащий множество специфических терминов, символов и их сочетаний, то программа не опознает большую часть текста. Поэтому вам лучше использовать режим обучения во время распознавания. При этом программа не будет принимать самостоятельных решений, а по каждому непонятному символу попросит у вас помощи.

Для включения данного режима перейдите **Сервис – Опции**. В открывшемся окне перейдите во вкладку **Распознать**. Установите переключатель опции **Распознавание с обучением**. Далее щелкните кнопку **Эталоны**, затем – **Новый** (рис. 5.1).

Ведите название эталона и щелкните кнопку **OK**. Закройте оба окна и

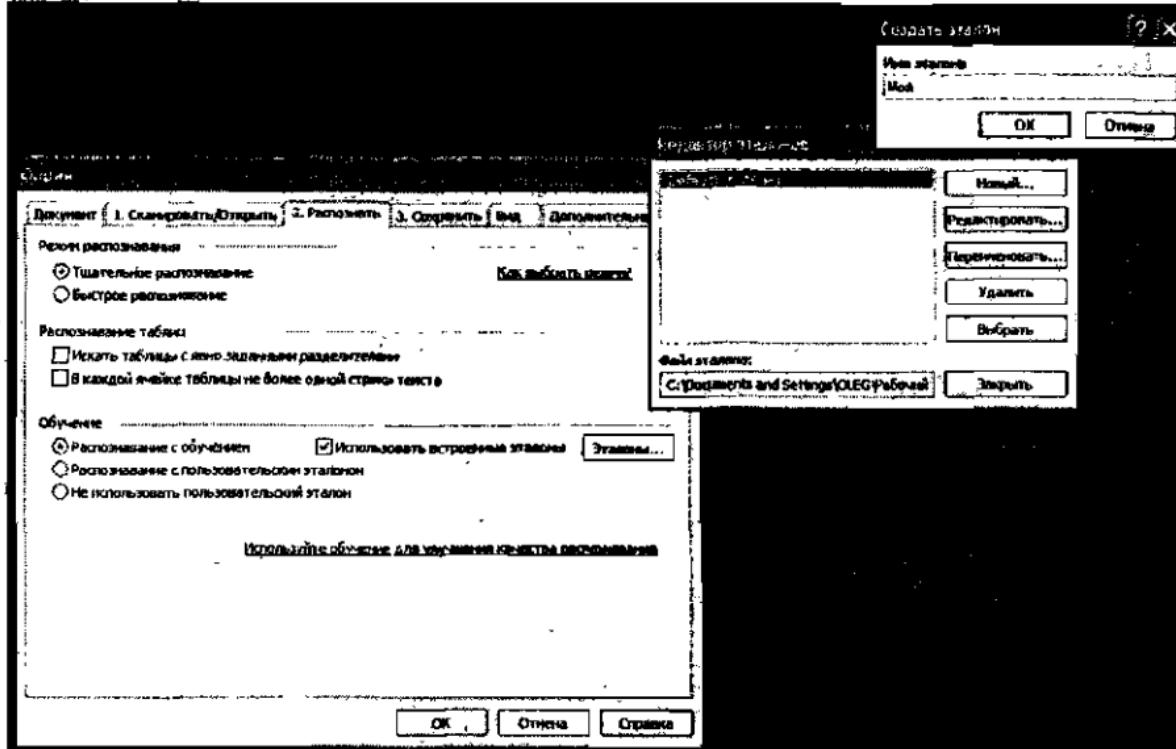


Рис. 5.1. Создание нового эталона

запустите сценарий распознавания. При первом же встретившемся неуверенно распознанном символе программа откроет окно Ручное обучение эталона, где вы должны «развеять сомнения» программы и ввести правильное значение символа, выделенного зеленой рамкой (рис. 5.2).

В примере на рисунке нужно ввести буквы **г**е и щелкнуть кнопку **Обучить**. При необходимости можно задать также формат символа (жирность, наклон, подчеркивание и др.).



Рис. 5.2. Введение правильного символа

Распознавание текста с цифровых изображений (фотографий)

Как правильно фотографировать документы, чтобы их фотографии потом можно было распознать (преобразовать в текстовый документ)

При фотографировании текстовых документов с целью дальнейшего распознания фотографий с преобразованием их в текстовые документы необходимо соблюдать определенный набор правил и условий:

Наиболее очевидное правило состоит в том, чтобы следить за тем, чтобы фотографируемая страницы целиком умещалась в кадре.

Необходимо добиться того, чтобы на поверхность фотографируемого документа падал равномерный текст без каких-либо теней и бликов. Лучше всего использовать дневное освещение (ну или освещение сразу несколькими лампами, чтобы избежать теней).

При фотографировании объектив фотоаппарата следует располагать параллельно поверхности документа, примерно по центру документа. Также следует по максимуму избавиться от неровностей бумаги. Например, если страница загибается ближе к корешку, ее необходимо прижать. Распознавание с геометрически неровных поверхностей может быть произведено с ошибками.

Что касается технических характеристик фотоаппарата, который можно использовать для рассматриваемых целей, то по минимуму они должны соответствовать следующим условиям:

- Размер матрицы 2 млн. пикселов.
- Переменная дистанция фокусировки. Не рекомендуется использовать фотоаппараты с фиксированным фокусом (как правило, такие камеры установлены в сотовых телефонах и КПК).

Рекомендуемыми же значениями характеристик фотоаппаратов являются следующие:

- Размер матрицы 5 млн. пикселов.
- Возможность отключения фотовспышки.
- Возможность установки диафрагмы вручную – т.е. наличие режима

58

приоритета диафрагмы или ручного режима.

- Режим ручной фокусировки.
- Объектив с оптической стабилизацией изображения, при его отсутствии рекомендуется использовать штатив.
- Оптический зум.

В ходе съемки фотоаппарат желательно размещать на расстоянии 50-60 см. Если имеется хороший оптический зум, то хорошего результата можно достичь, подальше расположив фотоаппарата увеличив зумом и использовав вспышку. Вспышку на близком расстоянии использовать нельзя, так как она будет давать блики и пересветы. А вот издалека она дает более-менее равномерное освещение.

Распознавание текста с цифровой фотографии

Чтобы распознать текст с имеющейся цифровой фотографии следует либо при старте FineReader'a выбрать сценарий **PDF или изображения в Microsoft Word**, либо уже, находясь в рабочем окне Finereader'a в строке меню выбрать **Сервис → Сценарии ABBYY Finereader → PDF или изображения в Microsoft Word** (рис. 6.1).

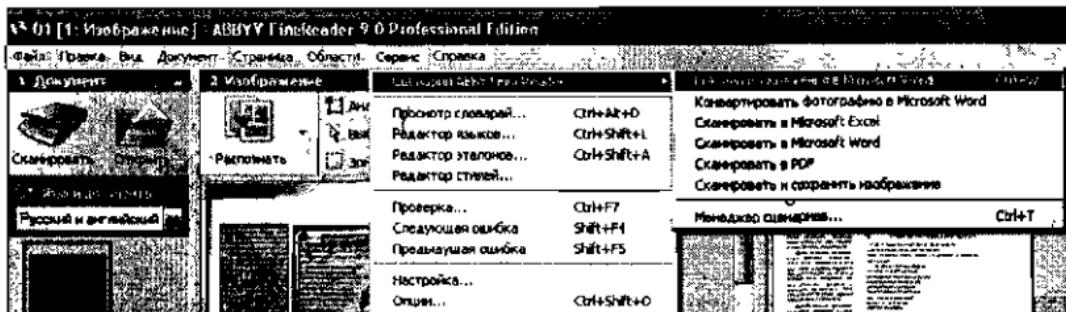


Рис. 6.1. Выбор нужного сценария

Далее вам будет предложено указать файл изображения. Как только вы его выберите, будет произведено автоматическое распознавание и сразу будет создан новый текстовый Word-документ. Если в автоматически получившемся результате вас что-то не устраивает, то вы вольны настроить вручную параметры и произвести повторно распознавание. Для этого следует вернуться в рабочее окно FineReader'a и произвести необходимые изменения, как было описано в гл.3.

Полезные приемы работы с FineReader'ом. Решение проблем с распознаванием текста

Как отключить автоматический запуск распознавания текста после сканирования

При выборе того или иного сценария распознавание будет автоматически запускаться сразу после того, как будет завершено сканирование. Однако это далеко не всегда бывает удобно. Чтобы отключить автоматический запуск распознавания текста следует в окне Finereader'a, в строке меню выбрать **Сервис → Опции**. Далее в появившемся окне **Опции** перейдите на вкладку **Сканировать/Открыть**. Здесь вы должны будете сместить переключатель **Автоматически распознавать полученные изображения** либо в положение **Автоматически анализировать полученные изображения**, либо в положение **Отключить автоматический анализ и распознавание изображений**. Рекомендуется устанавливать в положение **Автоматически анализировать полученные изображения**, так как в этом случае программа

по крайне мере будет пытаться автоматически определить где какие области с текстом находятся, отделить их от, например, графических областей. Если же выбрать последний вариант (**Отключить автоматический анализ и распознавание изображений**), то программа вообще ничего сама по своей инициативе делать не будет. Для всего вы должны будете сами нажать соответствующую кнопку в окне FineReader'a.

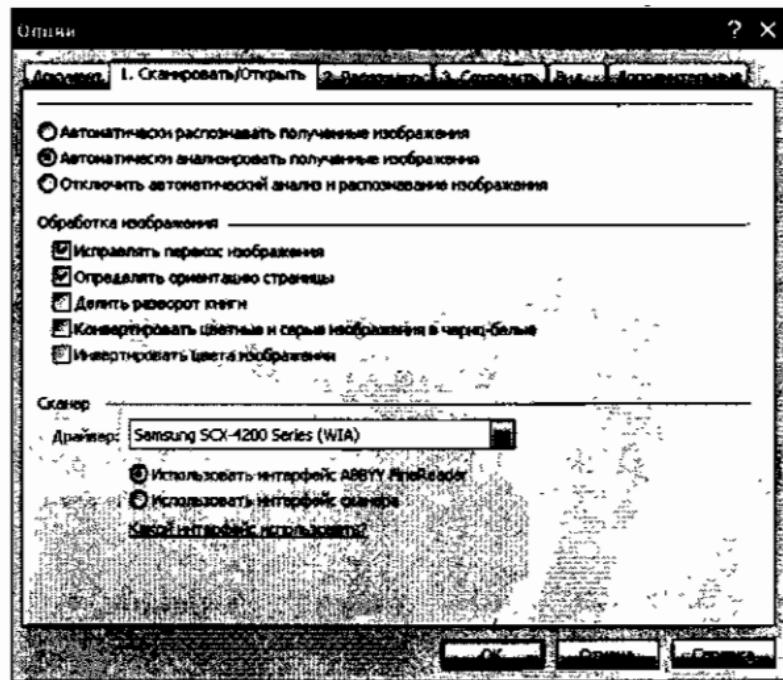


Рис. 7.1. Выбор режима автоматического распознавания

Как сделать так, чтобы при распознавании отсканированного/сфотографированного разворота книги FineReader автоматически разбивал разворот на отдельные страницы

Довольно часто сканируются книги. Так вот для ускорения работы гораздо эффективнее сканировать книги не по отдельной странице, а разворотами (то есть по две страницы, образующими разворот). Получается ровно в 2 раза быстрее. Но при этом необходимо настроить Finereader таким образом, чтобы программа автоматически разделяла разворот на отдельные страницы. В противном случае каждый разворот будет восприниматься как одна страница.

Чтобы произвести интересующую нас настройку необходимо в окне Finereader'a, в строке меню выбрать **Сервис → Опции**. Далее в появившемся окне **Опции** перейдите на вкладку **Сканировать/Открыть**. Здесь вы должны будете установить флажок **Делить разворот книги** и нажать **OK**.

Как повернуть страницу в окне

Во-первых, в настройках программы должен быть включен режим автоматического определения ориентации страницы, когда вы можете

отсканировать страницу вверх ногами, а программа автоматически повернет ее в нужном направлении на нужный угол. Если что, для включения данного режима необходимо в строке меню выбрать Сервис → Опции, далее в появившемся окне Опции перейти на вкладку Сканировать/Открыть и установить флагок Определять ориентацию страницы.

Во-вторых, вы в любой момент времени вы можете принудительно повернуть в окне FineReader'a текущую страницу, выбрав в строке меню Страница → Повернуть страницу и указав затем куда и как необходимо выполнить поворот (рис. 7.2).

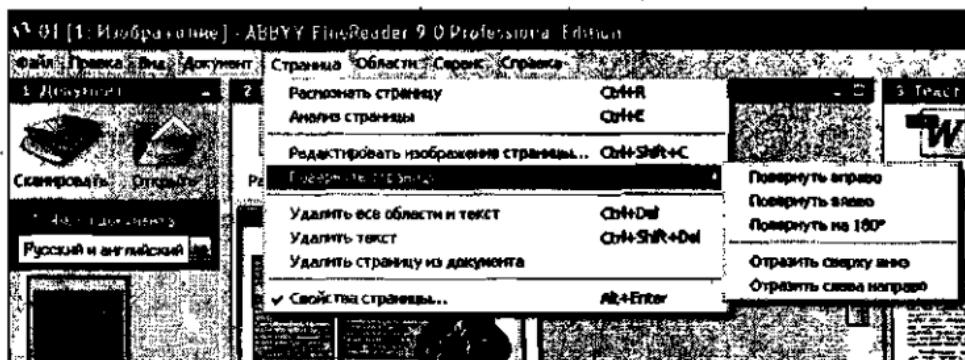


Рис. 7.2. Выбор варианта поворота

Как указать язык текста в документе для распознавания

В базовую комплектацию FineReader'a, как правило, включается два языка – русский и английский, в то время как в версии Professional языков много: итальянский, французский, немецкий и др. Так вот при распознавании иностранного текста или русскоязычного текста с вкраплениями иностранных слов, бывает полезно указать язык распознавания (если используемая версия программы поддерживает много языков).

Причем Finereader позволяет устанавливать пару языков, чтобы можно было эффективно распознавать смешанный текст.

Чтобы выбрать язык распознавания следует в окне

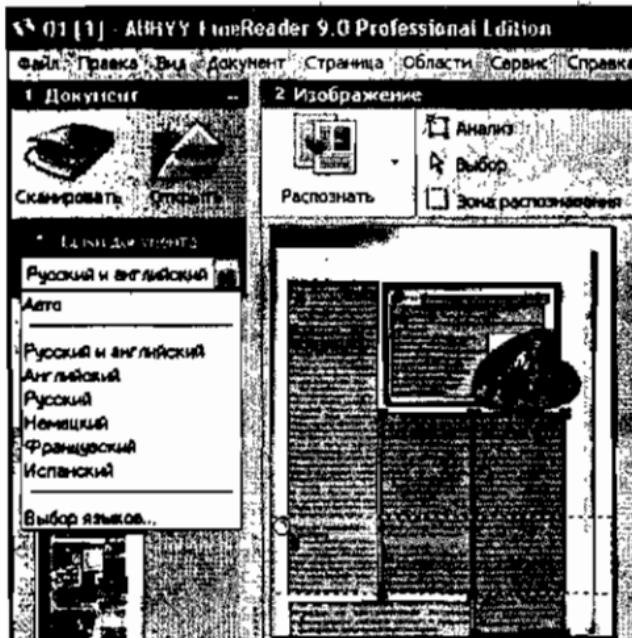


Рис. 7.3. Установка нужных языков

FineReader'a, в области Документ воспользоваться раскрывающимся списком Языки документа. Именно в нем и производится выбор.

65

Как включить или исключить какую-либо область страницы из распознавания

Вполне вероятна ситуация, когда вам требуется распознать не весь текст со страницы, а лишь отдельные части. После сканирования FineReader автоматически определяет/задает области сканирования, обводя их рамками. Причем каждой области присваивается определенный порядковый номер. Рамками зеленого цвета обозначаются текстовые области (которые затем должны будут переведены в текст), а рамками



Рис. 7.4. Удаление области

66 красного цвета – области с графическими изображениями, которые распознавать не надо, но которые должны быть включены в итоговый документ в виде вставленных изображений. Чтобы исключить какую-либо область из распознавания следует щелкнуть по ней правой кнопкой мыши и в появившемся контекстном меню выбрать команду **Удалить область**. После этого рамка с данной области будет снята, а ее содержимое никак не будет восприниматься программой и не будет участвовать в конечном документе (результате распознавания).

Чтобы включить какую-либо область в распознавание необходимо щелкнуть мышкой по:

- Кнопке  **Текст** – если вы хотите включить текстовую область
- Кнопке  **Картина** – если вы намерены включить в конечный документ картинку (без распознавания, в виде графического изображения)
- Кнопке  **Таблица** – если вы хотите включить область с таблицей.

Далее вам остается лишь очертить рамку вокруг той области на странице, которую вы хотите включить в распознавание и в конечный документ в выбранном качестве (текста, картинки или таблицы).

Подкорректировать границы той или иной области можно просто перетаскивая мышкой ее границы.

Как указать, что данный фрагмент страницы является изображением

Иногда программа Finereader не совсем корректно справляется с определением того, где текстовая область для распознавания, а где – графическое изображение, которое присутствует на отсканированной странице и которое распознавать не надо. Например, вы отсканировали листовку с рекламой магазина. На листовке размещена фотография магазина. Так вот по умолчанию FineReader при распознавании вывеску магазина на фотографии также переведет в текст, разрезав фотографию на несколько частей. А этого, скорее всего, вам совершенно не нужно: пусть фотографии магазина будет как фотография, в виде графического изображения. Возможна и обратная ситуация, когда Finereader посчитает часть текста содержимым графического изображения и не станет его распознавать.

Чтобы переопределить ту или иную область для распознавания как

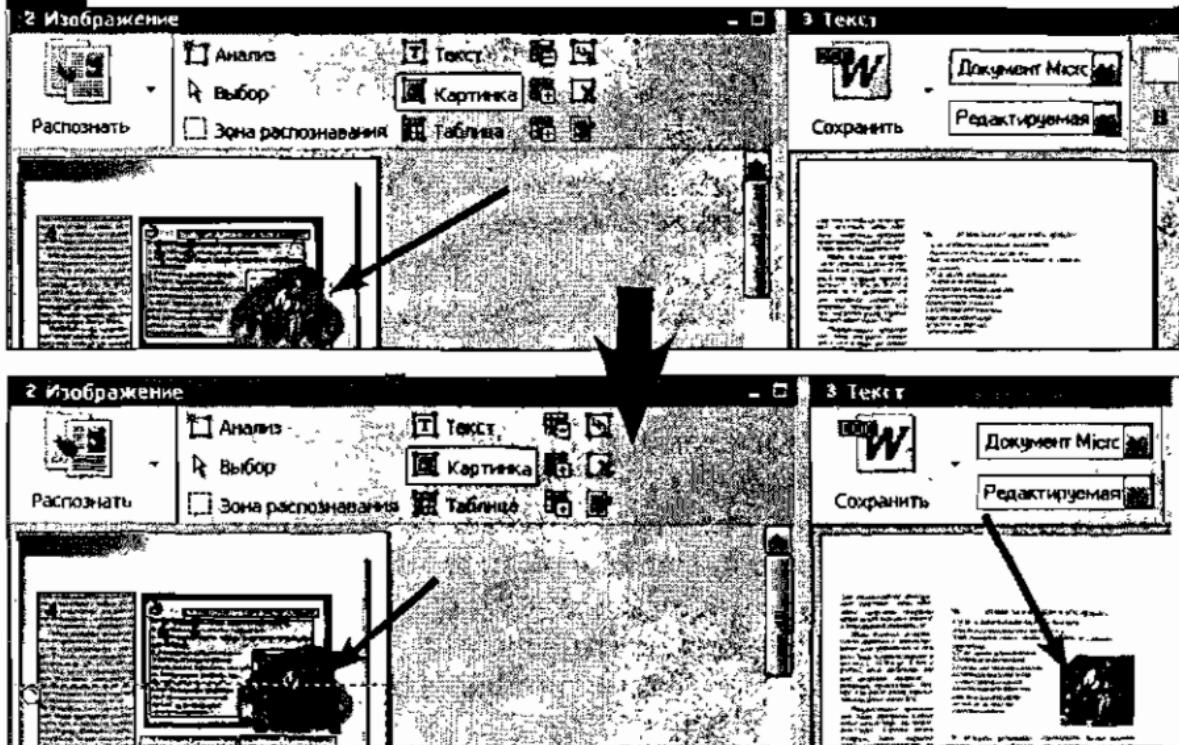


Рис. 7.5. Выделение области с изображением

текст или как изображение, следует для начала снять выделение областей с тех районов страницы, где разбиение на области произведено ошибочно. Затем, нажав кнопку **Текст** следует обвести области, подлежащие распознаванию как текст, а нажав кнопку **Картинка** обвести области, подлежащие распознаванию как изображения. Кроме того существует возможность изменения типа уже существующей области: например, можно текстовую область переделать в область изображения и наоборот. Для этого необходимо щелкнуть по области правой кнопкой мыши и в появившемся контекстном меню выбрать **Изменить тип области**, а затем указать на какой именно тип следует изменить (рис. 7.6).

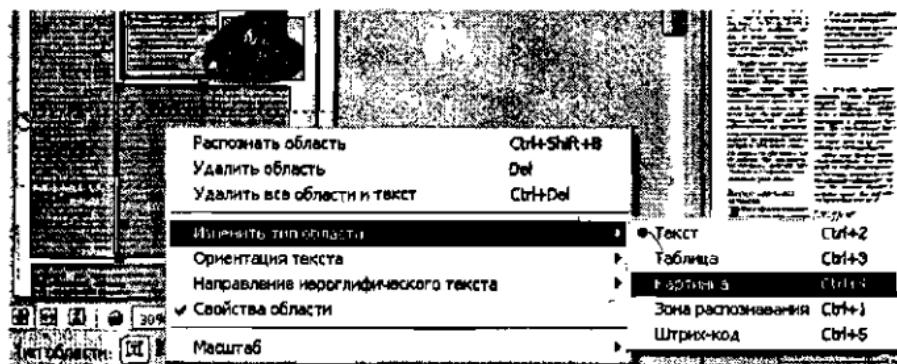


Рис. 7.6. Изменение типа области

Использование режима тщательного распознавания

По умолчанию в FineReader'е используется режим быстрого распознавания, которые отличается наибольшей скоростью и подходит для большинства случаев. Однако бывают сложные случаи, когда программе сложно распознавать текст и в быстром режиме она может не справляться.

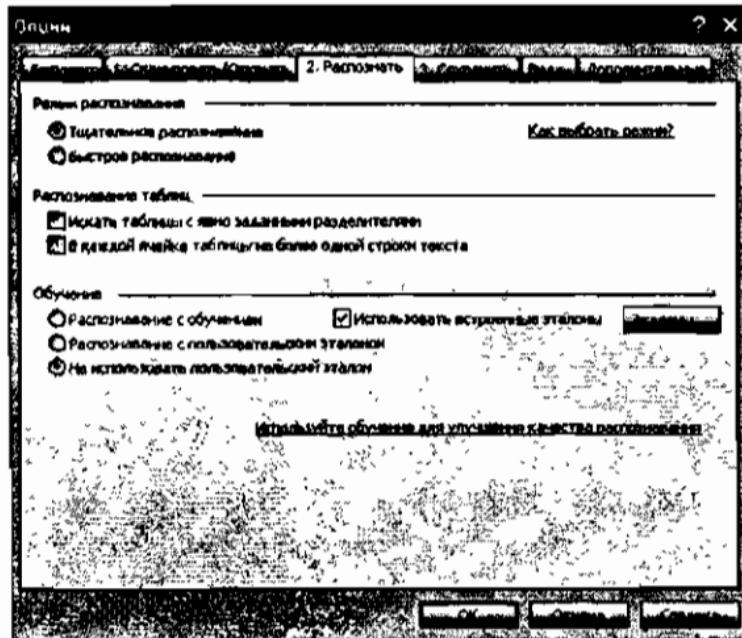


Рис. 7.7. Включение тщательного распознавания

К сложным случаям относятся, например, текст на цветном фоне, документы с таблицами (особенно трудно приходится FinReader'у когда линии сетки не прорисованы и используются цветные ячейки) и т.д.

Вот в подобных случаях и

рекомендуется использование так называемого режима тщательного распознавания. Который занимает существенно больше времени, но отличается лучшей результативностью и меньшим количеством ошибок.

Чтобы включить данный режим следует в строке меню выбрать **Сервис → Опции**. Далее в появившемся окне **Опции** перейдите на вкладку **Распознавание** и переключатель **Режим распознавания** установить в положение **Тщательное распознавание** (рис. 7.7).

Как дать команду на повторное распознавание текста после произведенных перенастроек

Вы в любой момент можете запустить повторное распознавание текущей страницы или всего документа, воспользовавшись кнопкой **Распознать** в области **Изображение**. Рядом с кнопкой, справа, находится направленная вниз стрелочка, щелкнув по которой вы сможете выбрать, что

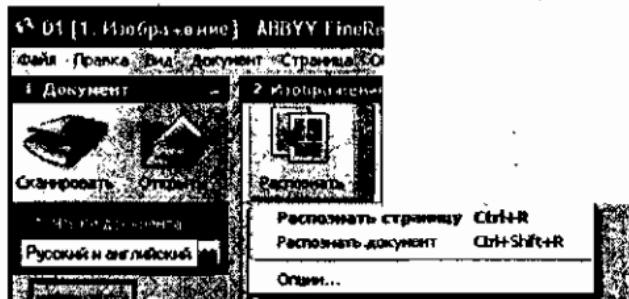


Рис. 7.8. Повторное распознавание

72 именно следует распознавать: только текущую страницу или весь документ целиком (рис. 7.8).

Кроме того можно воспользоваться горячими клавишами: «Ctrl» + «R» – для распознавания страницы, «Ctrl» + «Shift» + «R» – для распознавания всего документа.

Как перераспознать отдельную область на странице

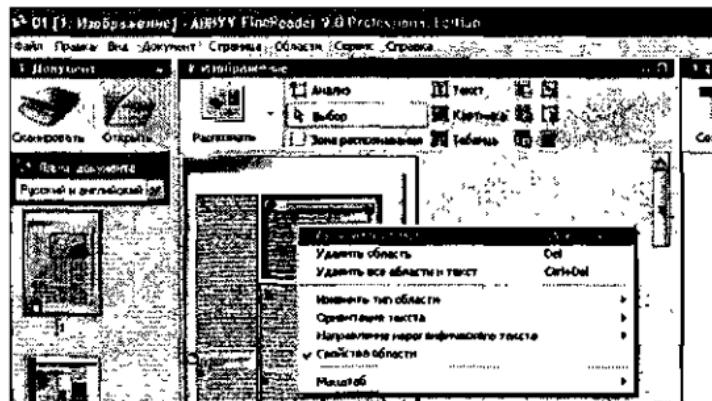


Рис. 7.9. Распознать область

В FineReader'е вы можете повторно распознать не только текущую страницу или документ целиком, но отдельную область на странице (если в других, допустим, вас все устраивает). Для этого необходимо щелкнуть правой кнопкой мыши по требуемой области и в появившемся контекстном

меню выбрать команду **Распознать область**.

Что делать, если в распознанном тексте некорректно отображается шрифт или на месте некоторых букв стоят значки «?» или «□»

Причина, по которой в конечном документе, а также в окне **Текст** (где предварительно отображается распознанный текст) вместо некоторых букв стоят значки «?» или «□», состоит в том, что в уже распознанном документе используется шрифт, содержащий не все символы языка документа. Соответственно чтобы избавиться от подобной проблемы необходимо поменять шрифт.

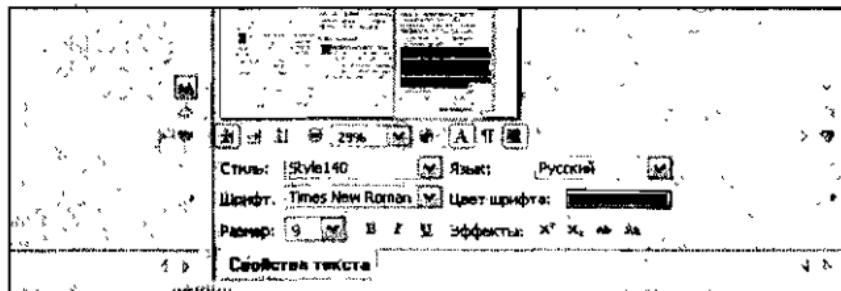


Рис. 7.10. Задаем шрифт

Сделать это можно несколькими способами в зависимости от ситуации, в которой вы находитесь. Самый простой способ состоит в следующем:

- 74**, В окне FineReader'a, в области Текст (№3) выделите фрагмент текста, в котором некорректно отображается шрифт.
2. Под предварительным видом распознанного текста в области Текст должна располагаться панель Свойства, на которой вы и сможете настроить шрифт, в частности в раскрывающемся списке Шрифт выбрать сам шрифт (рис. 7.10). Если панели Свойства по каким-либо причинам у вас на экране нет, то щелкните по абзацу правой кнопкой мыши и в контекстном меню выберите пункт Свойства. Панель должна появиться
3. На панели Свойства текста в списке шрифтов выберите шрифт.

Что делать, если в исходном тексте присутствуют нестандартные (декоративные, математические и т.д.) символы

Если в вашем исходном документе присутствуют какие-либо нестандартные символы, то для повышения качества распознавания рекомендуется в ходе этого процесса использовать режим **Распознавание с обучением**. При этом вы создадите несколько эталонов букв/спецсимволов, встречающихся в тексте и все они в дальнейшем будут распознаваться нужным об-

Что делать, если таблица в исходном документе не определена

Довольно распространенной ошибкой распознавания является неверное определение FieReDer'ом таблиц, когда программа воспринимает ее как набор строк текста. Особенно это часто происходит, когда таблица не имеет явно очерченных границ, либо когда границы прорисованы не все.

В таких случаях рекомендуется вручную задать область таблицы. Для этого необходимо:

1. В области **Изображение (№2)** щелкнуть по кнопке **Таблица** и обвести область страницы, которая должна распознаваться как таблица.
2. Щелкните по области правой кнопкой мыши и в появившемся контекстном меню выберите команду **Распознать область**, чтобы перераспознать таблицу.

Как вручную задать сетку таблицы, указать разбиение на столбцы и строки

Распространенной ошибкой распознания таблиц является некорректное понимание FineReader'ом разбиения таблиц на строки и столбцы. Особенно часто это наблюдается в тех случаях, когда в таблице неявно отрисованы все границы. Кроме того, иногда бывает полезно отсканированный текст оформить в виде таблицы, даже если изначально он не был в табличном виде.

В FineReader'e имеются средства, позволяющие вручную отрисовывать сетку таблицы для распознаваемого документа, указывая разбиение на строки и столбцы. Процедура всего этого такова:

1. Для начала необходимо выделить всю область таблицы, как было сказано в предыдущем пункте шпаргалки.
2. Далее, чтобы отрисовать в таблице вертикальный разделитель (разделяющий столбцы таблицы) следует щелкнуть по кнопке  и в области таблицы указать, где именно должен он находиться.
3. Аналогично для отрисовки горизонтального разделителя следует воспользоваться кнопкой .

4. Чтобы удалить какой-либо из имеющихся разделителей воспользуйтесь кнопкой  , выбрав ее, а потом щелкнув мышкой по ненужному разделителю.

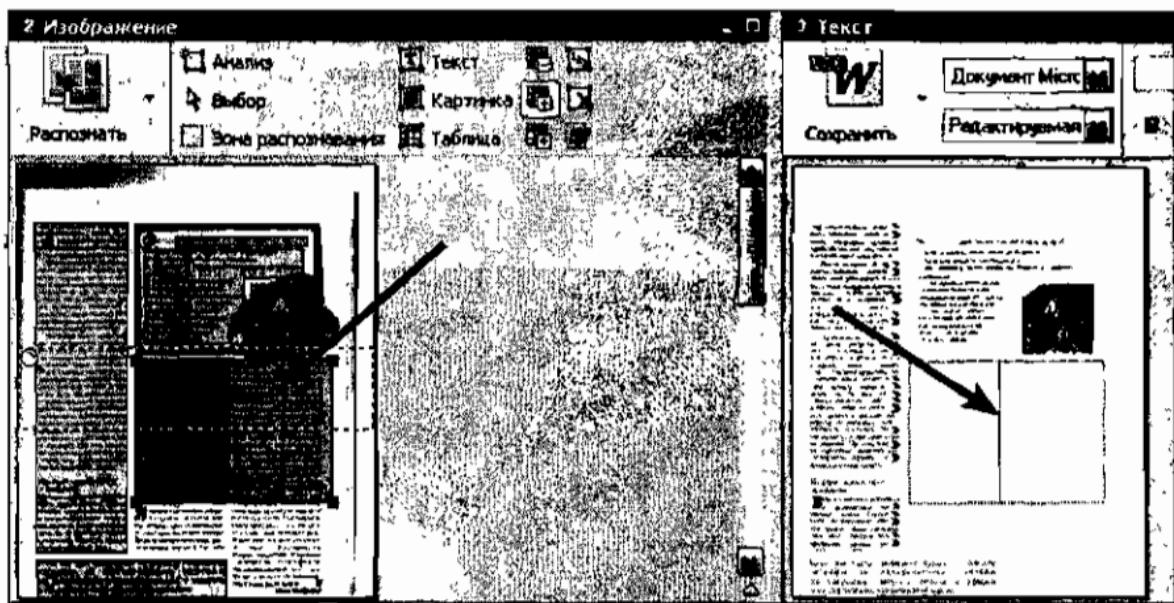


Рис. 7.11. Простановка разделителей в таблице

Как сохранить документ FineReader с отсканированными страницами для возможности дальнейшего повторного распознавания

Чтобы оставить для себя возможность повторного распознавания отсканированных/сфотографированных страниц в дальнейшем рекомендуется сохранять документы FineReader для дальнейшего возможного использования. Тем более, что

сделать это достаточно просто: в строке меню выберите **File → Сохранить документ FineReader**. После этого вам лишь потребуется указать имя файла (рис. 7.13).

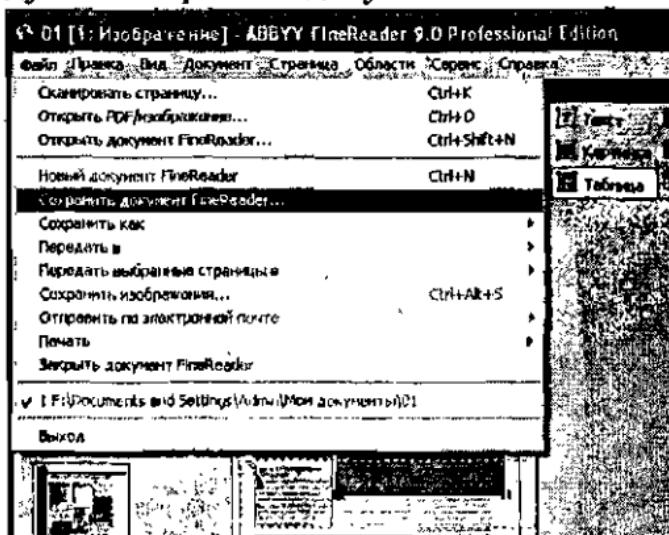


Рис. 7.12. Сохранение документа FineReader

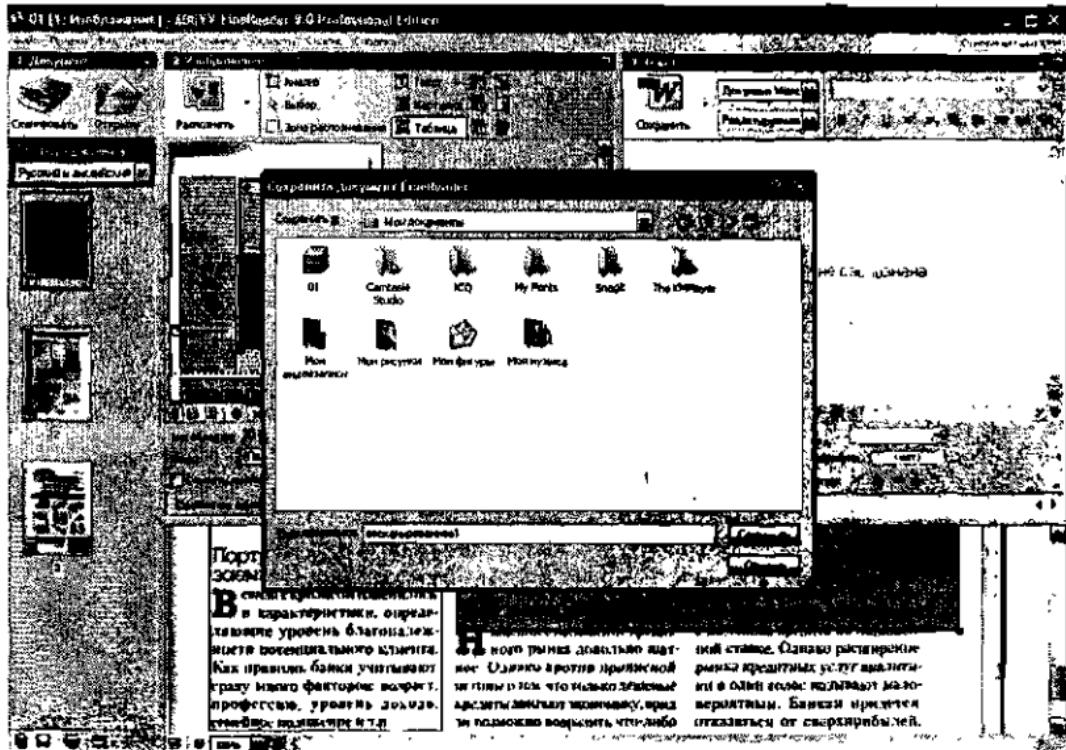


Рис. 7.13. Задание имени файла

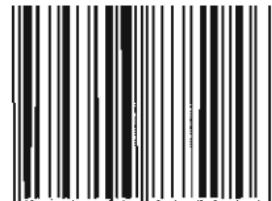
ООО «Наука и Техника»
Лицензия 1000350 от 23 декабря 1999 года.
198097, г. Санкт-Петербург, ул. Маршала Говорова, д. 29.
Подписано в печать 21.07.2009. Формат 60×88/32.
Бумага газетная. Печать офсетная. Объем 2,5 п. л.
Тираж 5000 экз. Заказ № 1019.

Отпечатано с готовых диапозитивов
в ООО «Гипнография Правда 1906».
195299, Санкт-Петербург, Киринская ул., 2.

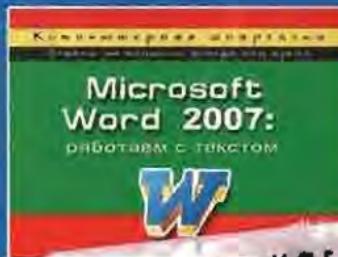
Компьютерные шпаргалки
от издательства Наука и Техника



ISBN 978-5-94387-594-6



9 785943 875946



om.ru

www

Ответы на вопросы всегда под рукой

НиТ
Наука и Техника