

Лабораторная работа №1

«Технология автоматического распознавания образов. OCR-системы»

1. Цель работы

Изучить возможности современных OCR-систем и приобрести навыки работы с ними при выполнении автоматического распознавания текста.

2. Подготовка к работе

Изучить основные понятия и принципы технологии автоматического распознавания образов и, в частности, автоматического распознавания (чтения) текста [1,2].

Ознакомиться с материалами по OCR-системам на web-сайтах производителей подобных систем (например, ABBYY и Cognitive Technologies).

Перед выполнением лабораторной работы необходимо установить соответствующее программное обеспечение:

- Ознакомительную версию *ABBYY FineReader* можно загрузить, перейдя по указанной ссылке – ABBYY FineReader 11 PE или запустить файл *ABBYY_FineReader_11_PE_TrialWithoutArabic.exe* в папке «8 - Дополнительные материалы» для установки программного обеспечения.

Системные требования *ABBYY FineReader 11*:

- Процессор с тактовой частотой 1ГГц или выше.
- Операционная система: *Microsoft Windows 7, Microsoft Windows Vista, Microsoft Windows Server 2008 R2, Microsoft Windows Server 2003, Microsoft Windows XP*.
- Объем оперативной памяти: не менее 1024 МБ, дополнительно для каждого ядра процессора 512 МБ.
- Свободное место на диске: 700 МБ для обычной установки и 700 МБ для работы программы.
- *TWAIN*- или *WIA*-совместимый сканер, цифровой фотоаппарат/фотокамера мобильного устройства или факс-модем.
- Видеоплата и монитор с разрешением не менее 1024×768 точек.
- Клавиатура, мышь или другое указательное устройство.
- Для конвертирования *PDF* в форматы документов *Microsoft Word, Excel, PowerPoint* и *Visio* соответствующие приложения *MS Office* должны быть установлены на компьютере.
- Загрузить OCR-систему с открытым кодом *CuneiForm* можно по следующей ссылке – Cognitive OpenOCR CuneiForm или запустить

файл *setup_openocr_cuneiform_rus.exe* в папке «8 - Дополнительные материалы» для установки программного обеспечения.

3. Лабораторное задание

1. Исследовать возможности и особенности *OCR*-систем (например, *ABBYY FineReader*, *CuneiForm* и др.) для выполнения распознавания изображения с текстом и преобразования его в документ *MS Word*, *pdf* и др.
2. Выполнить распознавание подготовленных трех изображений с помощью *OCR*-систем и результаты конвертировать в один из выбранных форматов (*MS Word*, *pdf* и др.) и сохранить в результирующий файл.
3. Сравнить полученные результаты и сделать соответствующие выводы по распознаванию изображений различного разрешения средствами *OCR*.
4. Для распознавания изображения текста низкого качества использовать возможность обучения по шаблону *OCR*-системы *ABBYY FineReader*. Результаты распознавания по шаблону конвертировать и сохранить в результирующий файл.
5. Подготовить отчет для защиты лабораторной работы №1.

4. Требования

- Наличие трех файлов с растровым изображением текста различного разрешения (**низкого:** < 100 dpi, **среднего:** 100–300 dpi, **высокого:** > 300 dpi).
- Для выполнения автоматического распознавания текста необходимо воспользоваться двумя разными *OCR*-системами или различными версиями одной *OCR*-системы.
- Исходные изображения, результирующие файлы с информацией о корректно распознанных и сомнительных символах (относительная величина ошибки в %) и выводы по работе с *OCR*-системами включить в отчет по лабораторной работе №1.

5. Методические указания

Лабораторная работа выполняется с помощью *OCR*-систем (например, *ABBYY FineReader* и *CuneiForm*). Подробное руководство пользователя по системе *ABBYY FineReader 11* расположено в папке «8 - Дополнительные материалы» (файл – *FR11_Guide_Russian.pdf*).

Для загрузки и распознавания подготовленных изображений используются стандартные инструменты и соответствующие пункты главного меню (рис. 1 и рис. 2).

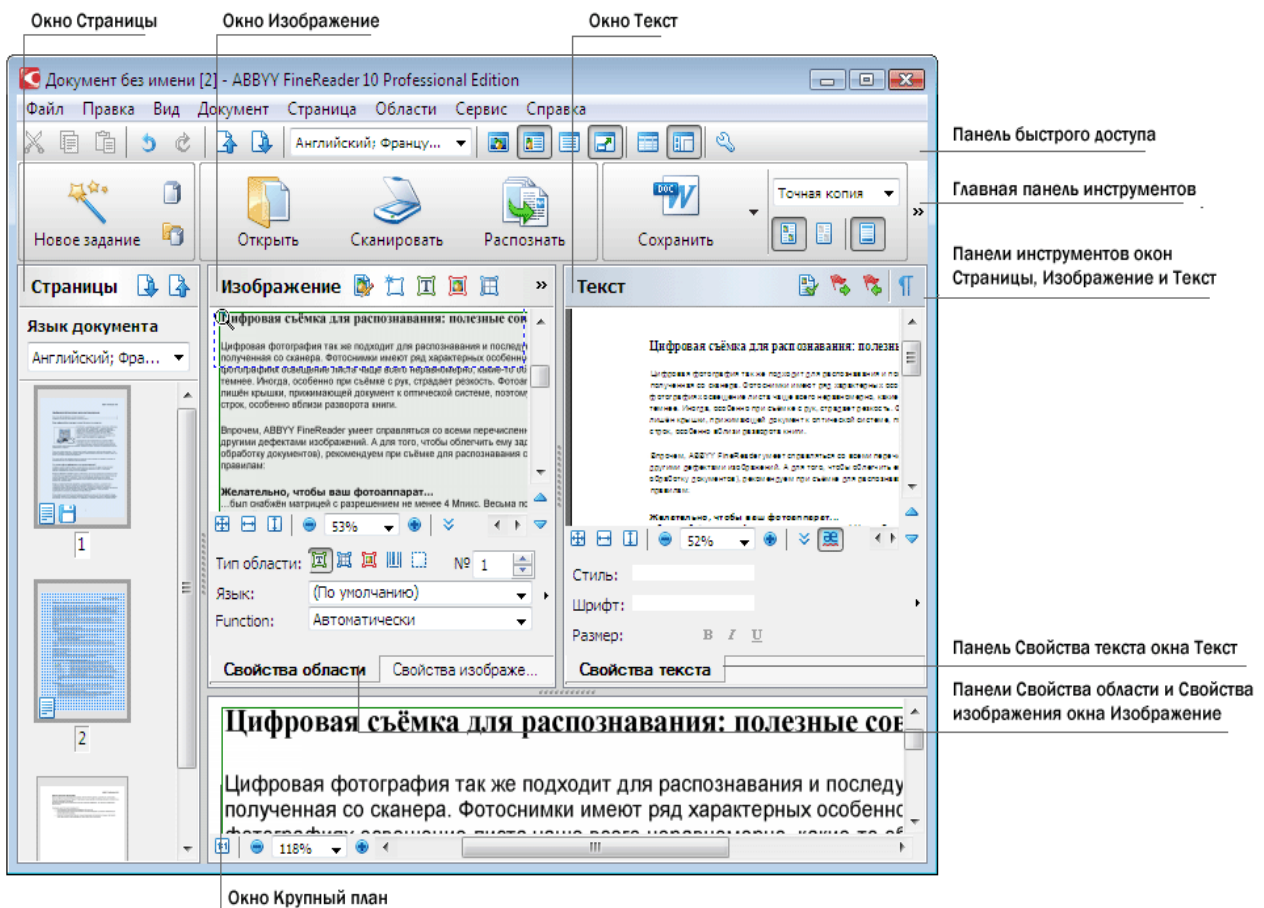


Рис. 1. Главное окно программы *ABBYY FineReader*

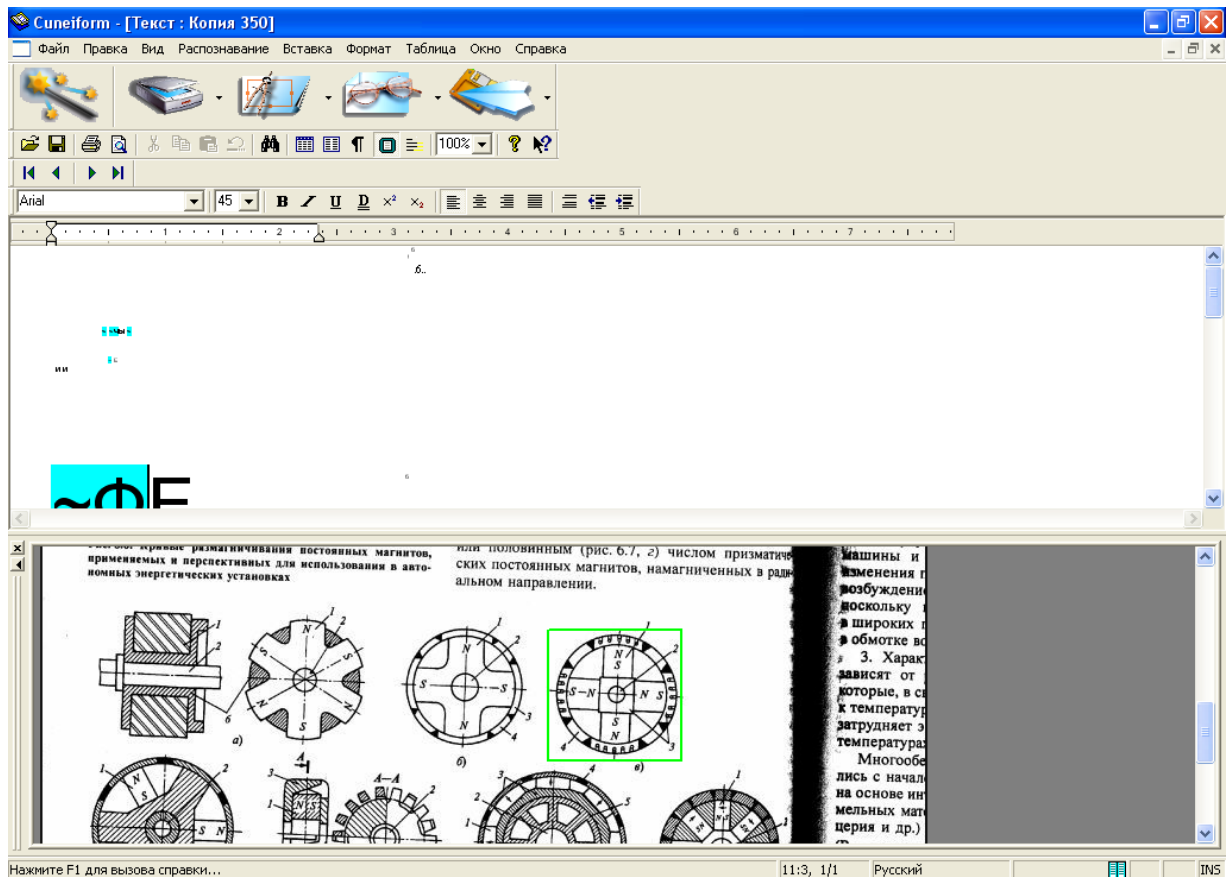


Рис. 2. Главное окно программы *CuneiForm*

Процедура распознавания с обучением в *ABBYY FineReader* предполагает предварительное **создание и обучение эталона**.

Рассмотрим эту процедуру по шагам:

1. Необходимо открыть диалог **Опции** (меню **Сервис>Опции...**) на закладке **Распознать**.
2. В группе **Обучение** установить переключатель в положение **Распознавание с обучением**.
3. Нажать кнопку **Эталоны...**
4. В открывшемся диалоге **Редактор эталонов** нажать кнопку **Новый...**
5. В появившемся диалоге **Создать эталон** ввести имя нового эталона и нажать **ОК**.
6. Нажать кнопку **Закреть** в диалоге **Редактор эталонов**, а затем кнопку **ОК** в диалоге **Опции**.
7. В окне **Изображение** нажать кнопку **Распознать**. Если в процессе распознавания встретится неизвестный символ, откроется диалог **Ручное обучение эталона** с изображением этого символа (рис. 3).
8. Необходимо обучить эталон **символам** или **лигатурам**. Лигатуры – это сочетания двух или трех символов, которые из-за особенностей их начертания невозможно разделить при обучении и которые поэтому сразу обучаются как комбинации символов. Обучение лигатурам происходит аналогично обучению отдельным символам.

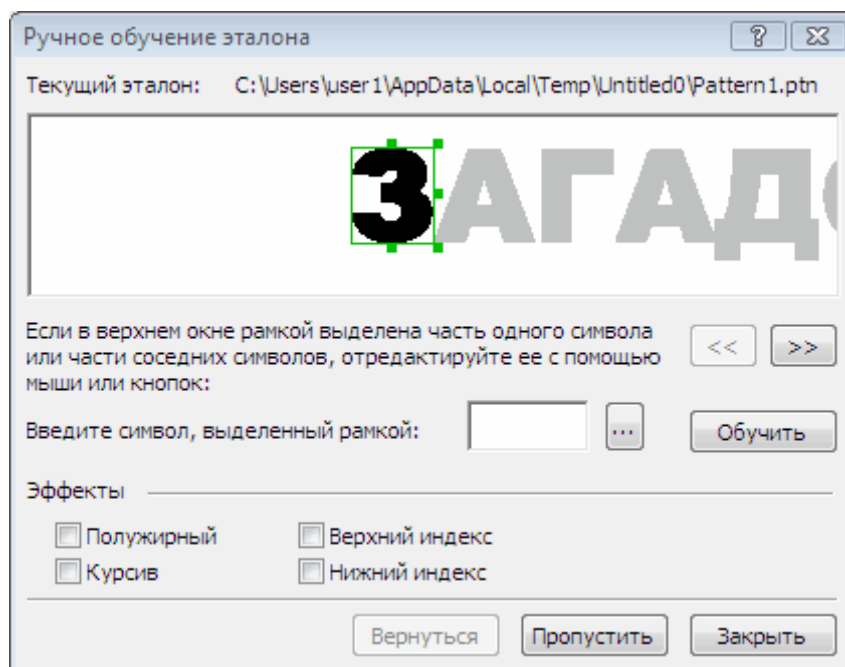


Рис. 3. Диалог Ручное обучение эталона

В процессе обучения можно вернуться к редактированию предыдущего символа нажатием кнопки **Вернуться**, которая действует в пределах одного слова. В этом случае последняя обученная пара «изображение – символ» будет удалена из эталона.

Отметим, что обучение возможно только для символов, входящих в алфавит языка. Если требуется обучить программу символам, которые нельзя ввести с клавиатуры, то для их обозначения можно использовать комбинации из двух символов, или можно скопировать требуемый символ из диалога **Вставка символа**.

В одном эталоне может содержаться до 1000 новых символов. Созданный эталон можно использовать только для распознавания текстов, использующих тот же шрифт, размер и отсканированных с тем же разрешением, что и исходный документ, на котором данный эталон обучался. Сохранить созданный эталон для работы с другими документами *ABBYY FineReader* можно сохранив настройки документа *ABBYY FineReader* в файл набора опций (*.fbt). В дальнейшем этот пользовательский эталон может быть отредактирован через диалог **Редактор эталонов**, а при необходимости отключен.

Для отключения пользовательского эталона достаточно на закладке **Распознать** диалога **Опции** (меню **Сервис>Опции...**) установить переключатель в положение **Не использовать пользовательский эталон**.

Литература

1. Башмаков А.И., Башмаков И.А. Интеллектуальные информационные технологии: Учеб. пособие. –М.: Изд-во МГТУ им. Н.Э. Баумана, 2005. – 304 с.
2. Варшавский П.Р., Куриленко И.Е., Михайлов И.С. Программное обеспечение интеллектуальных систем: учебное пособие / – М.: Издательский дом МЭИ, 2011. – 64 с.

Контрольные вопросы

1. Что такое *OCR*?
2. Области применения *OCR*-систем?
3. На каких трех принципах базируются все *OCR*-системы?
4. Основные виды классификаторов, применяемые в *OCR*-системах?
5. Перечислите основные особенности *OCR*-систем *ABBYY FineReader* и *CuneiForm*.