

# Содержание:

1. Введение.
2. Ввод и представление медицинских данных.
3. Методы ввода форм:
  - Традиционный метод ввода форм (вручную)
  - Современный метод ввода форм
  - Автоматизированный ввод данных
4. Что такое формы и что используют во избежание ошибок при их заполнении.
5. Системы оптического распознавания символов.

## Введение

При обращении за медицинской помощью, за справкой в медицинские учреждения многие сталкиваются с требованием дать согласие на сбор и обработку своих медицинских персональных данных. Часто это мотивируется тем что создается электронная карточка пациента и весь учет переводится в электронную форму.



# Ввод и представление медицинских данных.

В связи с большим объемом данных, требующих ручного ввода, возникает вопрос об автоматизации этого процесса. Наиболее оптимальным методом такой автоматизации является оптическое распознавание образов из готовых форм. При современном уровне развития технологий сканирования и распознавания скорость ввода можно довести до 100 000 листов за сутки при качестве автоматического распознавания 95-98%.

Автоматическое распознавание данных позволяет избежать таких ошибок, как неправильный набор данных оператором, набор данных в неправильных полях формы, расхождение данных по формату. Довольно простым и прозрачным становится масштабирование (увеличение числа рабочих мест) системы. Одновременно и качество ввода данных перестает зависеть от человеческого фактора, а скорость ввода легко прогнозируется и планируется.

**Ввод форм** – перевод данных, содержащихся в информационных полях заполненных форм, в электронный вид – включает в себя:

- получение (захват) данных из формы;
- оцифровка и сохранение изображения исходной формы.

Как правило, процесс считается завершённым, когда все заполненные поля формы обработаны, а все данные введены, проверены и импортированы в формат используемой базы данных. При этом обычно требуется не только обеспечить высокое качество данных, но и минимизировать трудозатраты.



# Методы ввода форм

- **Традиционный метод ввода форм (вручную)** силами одного сотрудника позволяет обеспечить небольшой объём обработки – до 20 форм в день. При больших объёмах ручной ввод чреват ненадежностью результатов, огромными трудозатратами, организационными издержками.
- **Современный метод ввода форм** подразумевает использование средств автоматизации. Так, например, широко распространенная в России система ABBYY FormReader решает задачу автоматического ввода данных с бумажных форм достаточно эффективно. Именно с ее помощью специалисты Департамента здравоохранения г. Москвы в сжатые сроки обработали результаты диспансеризации почти двух миллионов детей и подростков и в настоящее время имеют в распоряжении полную и постоянно обновляемую базу данных.

Автоматизированный ввод данных с бумажных форм с применением этой технологии состоит из следующих этапов:

- пачку заполненных форм сканируют при помощи скоростного сканера (обычно применяют аппараты с производительностью не менее 10 страниц в минуту);
- подавляющее большинство символов распознаётся системой автоматизированного ввода данных;
- символы, относительно которых сложилось несколько гипотез, автоматически передаются для проверки оператору системы ввода,
- подтверждённую информацию экспортируют в базу данных.







Потоковое сканирование направлений на клиничко-диагностические исследования (институт им. Склифосовского)

Как показывает практика, автоматизация ввода позволяет повысить скорость обработки в 10 раз, что позволит снизить число сотрудников, вовлечённых в обработку тоже в 10 раз. Все операции, кроме укладки пачки форм в приёмный лоток сканера и проверки сомнительных символов, выполняются автоматически. При этом качество данных оказывается на несколько порядков выше. Причины этого очевидны: влияние человеческого фактора сведено к нулю. Основной объём работы выполняется компьютером, который не устает и никогда не допускает опечаток .

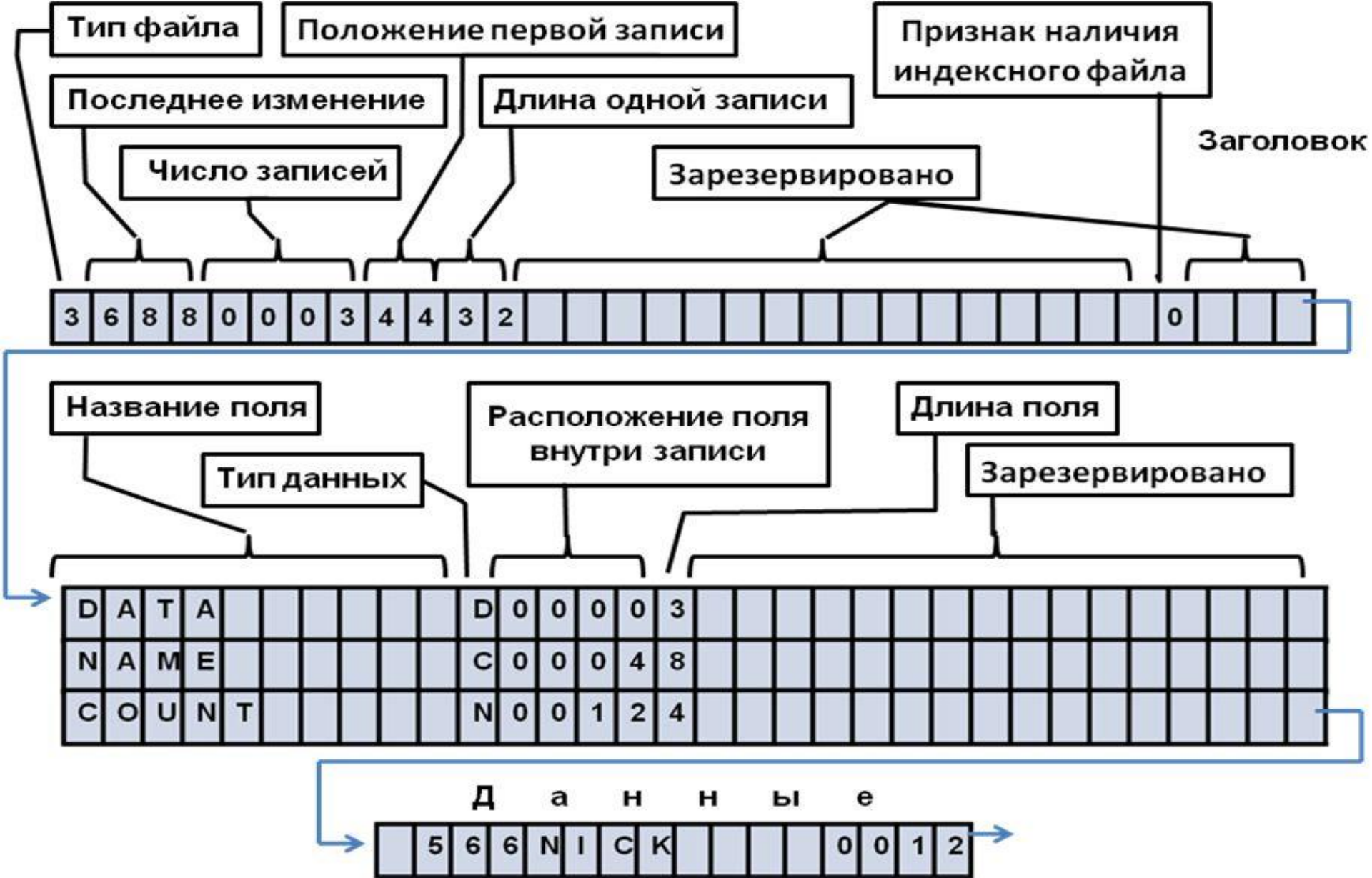




Производительность труда одного оператора, использующего ABBYY FormReader 6.5 Desktop Edition, составляет от 1000 до 3000 страниц в день в зависимости от сложности форм. При объёме ввода от 100 000 форм в месяц и больше целесообразно использовать комплекс на базе клиент-серверной версии ABBYY FormReader Enterprise Edition.

Формой называется документ, имеющий фиксированную структуру и предназначенный для сбора определённой информации. Основные признаки формы - наличие чётко обособленных полей, пояснительных надписей, а также ряда служебных элементов. В обиходе формы часто называют бланками. Примерами таких форм, широко используемых в медицине и здравоохранении, являются рецепты .





Примеры оформления машиночитаемых форм

При заполнении формы вся значимая информация заносится в поля - особым образом разграфлённые ячейки. Именно эта информация подлежит дальнейшей обработке. Формы, в которых определение положения полей и отделение данных от разметки может быть выполнено автоматически, программными средствами, называют машиночитаемыми. Необходимо отметить, что практически любая форма может быть приведена к виду машиночитаемой.

Форма может быть заполнена одним из следующих способов:

- от руки (такой способ заполнения называется рукопечатным: все символы пишутся раздельно, каждый символ занимает одно знакоместо);
- на пишущей машинке;
- распечатана на принтере;
- типографским способом;
- комбинированно, сочетанием вышеперечисленных способов



Во избежание ошибок заполнения формы составляют таким образом, чтобы сделать этот процесс интуитивно понятным. Для этого используют следующие специальные элементы, позволяющие ясно указывать заполняющему, какую информацию, куда и как следует вносить. Это: Информационные поля. Предназначены для внесения собственно данных. Существует три вида информационных полей:

**Текстовые поля.** Каждое из них представляет собой группу знакомест, обычно с пояснительной надписью. Основное назначение знакомест - побудить заполняющего написать символы отдельно.

**Метки, или пункты.** Метка выглядит как одиночный замкнутый контур (квадрат, круг, многоугольник), снабжённый пояснительной надписью. Информация в такое поле вносится путём простановки условного знака («галочки», креста) внутри контура, либо путём его полного закрашивания.

**Группы меток.** Так называются несколько меток, расположенных рядом и объединённых по смыслу. Снабжаются пояснительной надписью возле каждой метки, а также общей пояснительной надписью, раскрывающей смысл вопроса. Как правило, метки внутри одной группы соответствуют взаимоисключающим вариантам ответа.



Сервисные поля. В них располагаются реперные блоки, используемые при распознавании. С их помощью программа определяет правильную ориентацию формы и выравнивает искажения при оцифровке изображения. Иногда сервисные поля служат для идентификации бланка при одновременной обработке нескольких различных форм. В качестве реперных блоков на формах могут выступать следующие элементы:

сплошные квадраты чёрного цвета, углы и кресты;

сплошные линии: горизонтальные и вертикальные;

статический текст, то есть любая пояснительная надпись.

Идентификационные поля. Эти элементы предназначены для автоматической идентификации самого бланка формы. Реперный блок типа «квадрат», «угол» или «крест» обычно также может использоваться как идентифицирующий элемент. Но, как правило, для целей идентификации используют номер, нанесенный на форму в процессе изготовления бланка, само название формы или штрих-код.

Области для размещения графических изображений. Используются для размещения графических (нераспознаваемых) объектов. В качестве примеров подобных объектов можно назвать блок введения подтверждающей записи, печать или штамп. Изображения из этих областей можно помещать в базу данных в каком-либо графическом формате (TIF, BMP, JPG, PCX или WMF).

Декоративные, необязательные элементы: логотипы, колонтитулы и прочие элементы стилизации. При автоматизированном вводе информации зачастую используются для идентификации форм - анализируя текст в логотипе, программа может определить, от какой организации поступил данный документ (например, счет).



File View Option Tools About

Element properties	Element value
Group Number	0020
Element Number	4000
Group Description	Relationship
Element Description	Image Comments
VR	LT
VM	1
Length	0

(0002,0000) UL Group Length = <190>  
 (0002,0001) OB File Meta Information Version = <(#00)(#01)>  
 (0002,0002) UI Media Storage SOP Class UID = <1.2.840.10008.5.1.4.1.1.12.1>  
 (0002,0003) UI Media Storage SOP Instance UID = <1.3.12.2.1107.5.4.1.1921556504.2>  
 (0002,0010) UI Transfer Syntax UID = <1.2.840.10008.1.2.1>  
 (0002,0012) UI Implementation Class UID = <1.3.12.2.1107.5.4.1.2>  
 (0002,0016) AE Source Application Entity Title = <SIEMENS:DCR 1.60>  
 (0008,0005) CS Specific Character Set = <ISO\_IR 100>  
 (0008,0008) CS Image Type = <DERIVED\PRIMARY\SINGLE PLANE\SINGLE AS>  
 (0008,0016) UI SOP Class UID = <1.2.840.10008.5.1.4.1.1.1>  
 (0008,0018) UI SOP Instance UID = <1.3.12.2.1107.5.4.1.1>  
 (0008,0020) DA Study Date = <20060507>  
 (0008,0023) DA Content Date = <20060507>  
 (0008,0030) TM Study Time = <174906>  
 (0008,0033) TM Content Time = <181321>  
 (0008,0050) SH Accession Number = <10974>  
 (0008,0060) CS Modality = <XA>  
 (0008,0070) LO Manufacturer = <SIEMENS>  
 (0008,0080) LO Institution Name = <KKB>  
 (0008,0081) ST Institution Address = <>  
 (0008,0090) PN Referring Physician's Name = <>  
 (0008,1030) LO Study Description = <>  
 (0008,1050) PN Performing Physician's Name = <KEP>  
 (0008,1090) LO Manufacturer's Model Name = <POLYTRON>  
 (0009,0010) LO Private Creator = <POLYTRON-SMS 2.5>  
 (0009,1002) OB <Unknown Element> = <(#00)(#03)(#00)(#>  
 (0009,1004) OB <Unknown Element> = <(#00)(#03)(#00)(#>  
 (0009,1006) OB <Unknown Element> = <(#00)(#03)(#00)(#>  
 (0010,0010) PN Patient's Name = <AZEEV V.D.>  
 (0010,0020) LO Patient ID = <1057>  
 (0010,0030) DA Patient's Birth Date = <19480313>

Image view

Transfer Syntax: Explicit VR Little Endian      Format: Part10

Машиночитаемая форма медицинской карты  
 Всероссийской диспансеризации детей

# Системы оптического распознавания символов

Среди систем оптического распознавания символов выделяют два основных класса: OCR (англ. Optical Character Recognition - Оптическое распознавание символов) и ICR (англ. Intelligent Character Recognition – интеллектуальное распознавание символов). OCR-системы распознают печатные символы, нанесенные на бумагу типографским способом, при помощи принтера, плоттера или пишущей машинки. ICR-системы обрабатывают документы, заполненные печатными буквами и цифрами от руки, или, иначе говоря, распознают рукопечатные символы.



Принципы действия этих систем различны. OCR-система в процессе анализа выделяет на изображении блоки (текст, таблицы, иллюстрации), затем последовательно разделяет блоки на всё менее крупные объекты: абзацы, строки, слова, символы. Последние обрабатываются программными механизмами, осуществляющими собственно распознавание; эти механизмы называют классификаторами. Затем распознанные символы «собираются» в слова, слова – в строки, и так далее, вплоть до синтеза полного электронного аналога исходного документа.

ICR-система, нацеленная в первую очередь на обработку форм, функционирует иначе. На исходном изображении выделяются области, в которых должна содержаться смысловая информация, и затем именно эти фрагменты подвергаются дальнейшей обработке, в том числе и при помощи классификаторов. Иначе говоря, ICR-система не пытается построить точную электронную модель документа, а лишь извлекает информацию из чётко ограниченных областей. Эта информация передаётся в систему хранения.

К ICR-системам предъявляется также требование по распознаванию специальных объектов – меток. Их использование в формах позволяет упростить заполнение форм и значительно, до 99,9%, повысить качество ввода. Для распознавания отмеченной метки система анализирует и сохраняет информацию о распределении чёрного цвета в указанной окрестности метки. Если уровень затемнения отчетливо превышен, выносится решение о том, что данное знакоместо отмечено.





Подготовка машиночитаемой формы является весьма непростой проблемой, требующей одномоментного и эффективного решения двух задач: простого и удобного заполнения формы и эффективного ее распознавания. В силу этого структура формы, ее компоновка, формат бумаги, название и длина полей, использование разделителей и меток должно быть тщательно продумано. Очень важным является правильный выбор и расстановка реперных и идентификационных элементов формы, от которых зависит эффективность автоматического ввода данных. По этой причине машиночитаемый документ по внешнему виду может сильно отличаться от учетной формы, лежащей в его основе (см. рис. 77).

Для создания несложных форм можно использовать общеизвестный текстовый редактор Microsoft Word или пакет Microsoft Visio. Несмотря на то, что последний предназначен для рисования графиков и схем, с его помощью можно делать неплохие формы. Специальным инструментом для создания форм является программа FormDesigner из комплекта поставки ABBYY FormReader. Этот простой и удобный инструмент позволяет быстро и безошибочно изготовить форму любой сложности.

На следующем этапе программа распознавания настраивается на работу с созданной или имеющейся формой. В случае, если бланк формы подготовлен с использованием ABBYY FormDesigner, для этого достаточно импортировать в ABBYY FormReader шаблон в формате XFD, созданный при помощи ABBYY FormDesigner. В таком шаблоне на изображении формы уже будут содержаться размеченные блоки.



И, наконец, если сканирование документов производится впервые, требуется приобрести сканер. Именно от параметров выбранного сканера будет зависеть и скорость, и качество обработки данных. Следует отметить, что при большом количестве форм (свыше 100 ежедневно) обычные планшетные сканеры неприменимы. Эти устройства широко распространены в офисах, они неплохо справляются с оцифровкой фотографий и обычных документов, но для потокового ввода непригодны: у них невысокое быстродействие и относительно небольшой ресурс.

При выборе сканера для автоматизированного ввода форм в первую очередь следует обратить внимание на следующие основные его параметры и возможности:

Формат бумаги. Чаще всего для ввода форм используются устройства, способные сканировать листы формата А3, А4 и А5.

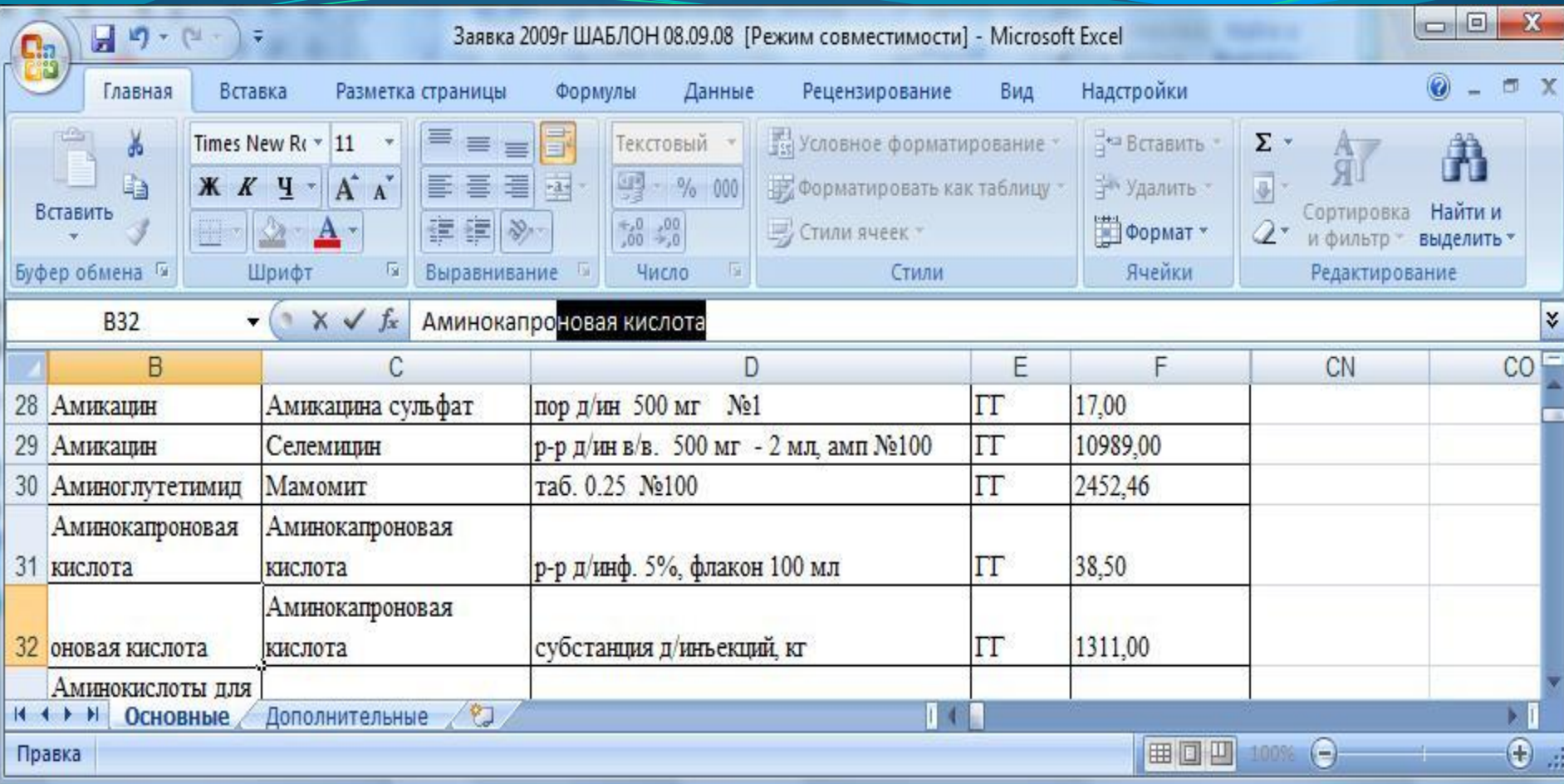
Разрешение изображений. Для ввода форм требуется разрешение 200-300 dpi. Сканирование с более высоким разрешением приводит к неизбежному замедлению, а скорость сканирования может быть одним из самых критичных параметров при потоковом вводе форм.

Двустороннее сканирование. Для многих проектов необходимо применять сканеры, которые могут осуществлять как одностороннее, так и двустороннее сканирование в черно-белом или цветном режимах. Последний режим необходимо использовать, например, при очистке изображения от цветных печатей и сохранении цветных фотографий с анкет.

Наличие устройства для автоматической подачи бумаги – автоподатчика. Это устройство, позволяющее загружать формы в сканер пачками обычно по 25, 50 или 100 документов, необходимо практически в любом случае, иначе работа оператора ввода будет на 90% состоять из манипуляций с бумагой и сканером.

Производительность. Часто скорость работы всей системы автоматизированного ввода зависит именно от быстродействия выбранного сканера. Выделяют три основные категории сканеров: офисный малой производительности (500 листов в день), офисный средней производительности и высокопроизводительные промышленные сканеры (более 20 тысяч листов).





Пример офисного сканера для потоковой обработки форм: Trüper 3200 фирмы Böwe Bell + Howell (США)

Для потокового ввода документов в системе ABBYY FormReader практически не требуется специально обученного персонала. Обычно привлекаются операторы, которые выполняют потоковый ввод форм, и администратор комплекса, который занимается настройками и выполняет контрольные функции. Однако, качество работы системы в целом определяется, как всегда, работоспособностью всех ее элементов, главные из которых были нами рассмотрены.

