



# Вопросы стандартизации

## Электронный документ в библиотеке: настоящее и будущее

**Настоящее**  
СЕГОДНЯ наша жизнь полна электронных документов: письма, заметки, инструкции, бланки, изображения... Библиотеки, как живой действующий организм, генерируют огромное число различных электронных документов, начиная со страниц официального сайта самой библиотеки и заканчивая электронными изданиями, базами данных, цифровыми копиями книг и т. д.

Поддавляющее большинство электронных документов нигде в библиотеке не регистрируется и в документации никак не отражается, хотя они активно используются, например, при обслуживании читателей/пользователей. При этом библиотечные электронные документы являются несомненным интеллектуальным активом, находящимся в распоряжении библиотеки. Библиотечные процессы сегодняшнего дня немыслимы без локального или удалённого использования электронных документов, хранящихся как в библиотеке, так и вне её.

Количество электронных документов огромно и растёт в геометрической прогрессии. Существует множество различных способов их создания, представления и хранения. Ценность электронных документов как совокупности информационных единиц определяется не только суммой их содержаний, но в значительной степени тем, каким образом они хранятся и сохраняются, обрабатываются и предоставляются пользователям, загружаются и перемещаются внутри библиотеки, между библиотекой и пользователем. Именно отлаженность и бесперебойное осуществление этих

процессов является сегодня одним из ключевых факторов экономической эффективности и степени социальной значимости библиотеки как культурно-информационного центра. По этой причине электронные документы, подобно любым другим типам библиотечных материалов, необходимо классифицировать, структурировать, проверять, оценивать, защищать, контролировать, измерять, то есть стандартизировать.

### Необходимость стандартизации

Создание национального библиотечного стандарта, отражающего основные характеристики и виды электронного документа крайне актуально на сегодняшний день. Особенно это касается вопросов терминологии: с одной стороны она должна быть единой, не противоречащей смежным областям деятельности, с другой стороны должна не противоречить международным терминам и стандартам. В первую очередь необходимо дать определение самому электронному документу, процессам его создания, а затем, по возможности, унифицировать типовые процедуры.

Говоря об актуальности вопросов стандартизации, нельзя не учитывать и «Концепцию развития национальной системы стандартизации Российской Федерации на период до 2020 года», одобренную распоряжением Правительства РФ от 24 сентября 2012 г. №1762-р.

Таким образом, основная задача для библиотек и органов научно-технической информации сводится к разработке национального стандарта, который однозначно определял бы понятие электронного документа и виды его существования. В дальнейшем это могло бы

**Электронные документы, как и любые другие, нуждаются в классификации, проверке, структурировании, оценке, защите, измерении — иными словами, в стандартизации.**



*Ольга Владимировна Барышева, ведущий программист Отдела перспективных электронных проектов Российской национальной библиотеки*



*Олег Николаевич Шорин, заместитель генерального директора Российской национальной библиотеки по информатизации*

привести к объединению усилий всех крупных поставщиков и владельцев электронных документов для разработки оптимальных по форме и формату представлений об электронных документах, используемых в библиотеках. Разработка стандарта на электронный документ должна стать базой для дальнейшей разработки как общих рекомендаций, так и стандартов на другие отдельные процессы и процедуры.

В то же время важно не утонуть в потоке электронных документов и чётко обозначить границы деятельности библиотек, связанные с электронными документами. Также необходимо обозначить типы и виды самих электронных документов, которые необходимо собирать, обрабатывать, хранить, предоставлять и т. д.

В данной статье мы представляем взгляд Российской национальной библиотеки на проблему стандартизации данных об электронном документе. В статье обобщаются предварительные предложения по формулировке определений и отбору понятий для дальнейшей совместной работы организаций, осуществляющих библиотечно-информационную деятельность, органов научно-технической информации, организаций, официально выпускающих в публичное обращение электронные документы с целью их массового использования и в научно-исследовательских, образовательных, культурно-просветительских целях.

### С чего начать?

Итак. Первое, что необходимо сделать, — установить термины и определения основных понятий в области электронных документов; набор характеристик, позволяющих проводить их идентификацию; базовые требования к формату представления электронных документов.

На наш взгляд, библиотекам нет смысла распространять свои разработки на электронные документы типа бланков и шаблонов оформления; на документы, предназначенные исключительно для автоматизированной обработки, а также на электронные документы в составе баз данных и комплексных информационных ресурсов; оформленные в виде электронных изданий (по ГОСТ Р 70.83-2012); компьютерные

программы, электронные подписи и их аналоги, финансовые документы и документы ограниченного распространения (в том числе отчётно-учётную документацию).

Электронные документы, запакованные в контейнеры (архивы) с помощью программ архивации, необходимо рассматривать только после распаковки.

При всём многообразии существующих стандартов СИБИД и смежных областей, нормативная база, на которую можно реально опираться при создании новых стандартов, весьма невелика:

- ГОСТ Р 70.83-2012 Система стандартов по информации, библиотечному и издательскому делу. Электронные издания. Основные виды и выходные сведения
- ГОСТ Р 52292-2004 Информационная технология. Электронный обмен информацией. Термины и определения
- ГОСТ Р ИСО/МЭК 26300-2010 Информационная технология. Формат Open Document для офисных приложений (OpenDocument) v1.0
- ГОСТ 8.417-2002 Государственная система обеспечения единства измерений. Единицы величин.

### Термины и определения

Нами предлагаются следующие термины с соответствующими определениями:

**электронный документ:** созданный программными средствами наделённый самостоятельным контентом и оформлением нетиражный электронный объект, анализ содержания которого может быть представлен в формализованном виде, предназначенный для передачи во времени и пространстве в целях хранения и использования, которое может быть регламентировано административными, правовыми и другими нормами.

**электронный объект:** файл (совокупность файлов), формируемый в компьютерной программе пользователя или автоматизированной системе и содержащий в зафиксированном виде данные, предназначенные для восприятия компьютером и/или человеком с помощью соответствующего аппаратного и программного обеспечения.

*Примечание — понятие электронного объекта является родовым по отношению к электронному документу.*

**электронный (информационный) ресурс:** комплекс электронных источников информации, программного обеспечения и аппаратных средств, служащих для удовлетворения информационных потребностей.

*Примечание — понятие электронного документа является видовым по отношению к электронному ресурсу.*

**контент электронного документа:** содержимое, наполнение электронного документа в плане содержания (в отличие от формы).

**форма/оформление электронного документа:** выполнение формальной обработки содержимого документа в соответствии с целевым назначением и/или правилами/нормами его использования при неизменности контента.

*Примечание — форма, как правило, соотносится с определённым шаблоном электронного документа (например, электронный документ в форме письма), оформление — с изменением его внешнего вида (например, цветовой оформление). Соотнесение формы электронного документа с родом, и/или видом/жанром заключённого в нём произведения, контента (например, электронный документ в форме стихотворения, марша, сборника) не рекомендуется.*

**версия электронного документа:** формально идентифицированное уникальное качественное состояние контента электронного документа во временном ряду по отношению к электронным документам с тем же основным контентом.

*Примечание — качественно новый, оригинальный электронный документ не имеет версии.*

**идентификатор версии:** обозначение единицы в нумерованной/именованной последовательности качественных изменений контента электронного документа.

**эталонная версия электронного документа:** образец, шаблон контента электронного документа для создания последующих версий/редакций.

**редакция электронного документа:** результат процесса редактирования — создание обработанного и исправленного варианта существующего электронного документа или одной из его версий (в том числе локализация; изменение формы/оформления; не принципиаль-

ные изменения контента, не ведущие к качественному преобразованию, например исправление ошибок, перестановка абзацев, и т. д.).

**Примечание** — редакция не имеет обязательного формального идентификатора.

**компиляция:** способ составления контента электронного документа на основе использования/заимствования данных из уже существующих сторонних электронных документов

**копия электронного документа:** результат процесса копирования — дублирование, повторение электронного документа способом, отличным от способа его создания.

**Примечание** — при создании копии возможно изменение формы/оформления, формата, знаковой природы первичного документа, но не её контента. Копирование может быть произведено как с аналогового, так и с электронного документа; с оригинала, версии, редакции, другой копии.

**сжатие/кодирование электронного документа:** алгоритмическое преобразование данных, производимое с целью оптимизации использования электронного документа.

**Примечание** — Сжатие/кодирование производится, как правило, для уменьшения объёма или ускорения загрузки электронного документа. Различие между ними состоит в том, что сжатый документ пригоден для непосредственного использования, а для кодированных электронных документов необходимо использование декодера. Степень сжатия и качественного изменения данных электронного документа зависят от используемого коэффициента/алгоритма. Все методы сжатия/кодирования данных делятся на два основных класса: с потерями и без потерь. Электронный документ может быть сжат/закодирован как целиком, так и частично, причём для разных составных частей могут быть использованы различные алгоритмы и коэффициенты сжатия.

**метаданные:** зафиксированный в определённой форме структурированный набор характеристик электронного документа, организованных в соответствии с определённой схемой, и предназначенный для идентификации, поиска, оценки и управления электронными документами.

**Примечание** — Метаданные формируются на основе схемы методом выделения общего для всех электронных документов и обязательного для использования при их обработке набора полей, правил структурирования областей и элементов, извлечения данных из электронного документа и приведения их в соответствии с предписанным синтаксисом.

**схема метаданных:** стандартизованный набор и структура представления метаданных, предназначенный для формального описания электронных документов.

**Примечание** — Схема метаданных включает в себя набор полей (атрибутов, свойств, элементов), отражающих характеристики электронного документа.



**формат данных:** конкретная форма представления данных, в которой установлены ограничения типа данных (по ГОСТ Р 52292-2004).

**Примечание** — формат файла является частной формой формата данных.

**формат (файла) электронного документа:** определённая спецификация, описывающая структуру файла, в соответствии с которой пакеты данных могут быть сохранены как файлы, переданы по сети в виде потока данных, и интерпретированы.

**Примечание** — в ряде операционных систем расширение имени файла яв-

ляется видимым для пользователя символьным идентификатором типа файла, недостаточным для полной идентификации формата электронного документа.

**размер (файла) электронного документа:** автоматически определяемое компьютером количество информации в стандартных единицах (по ГОСТ 8.417-2002).

**Примечание** — Фактический объём дискового пространства, занимаемого файлом, зависит от конкретной файловой системы.

**открытый формат:** свободная от лицензионных ограничений при использовании общедоступная спецификация (стандарт) хранения цифровых данных, позволяющая переносить их с одной программной платформы на другую без искажения формы, структуры, содержания.

**Примечание** — Не следует смешивать понятия открытого формата и свободной лицензии на использование. Открытость заключается в доступности спецификаций и соответствии открытого формата электронного документа стандарту, понятие свободы относится к передаче прав и является одной из моделей лицензирования.

**идентификация электронного документа:** анализ электронного документа по одному или нескольким характерным признакам с целью опознания, определения сходства/различия, отнесения к конкретному классу/виду, типу.

**идентификатор:** выбранный по какому-либо основанию деления признак, фиксирующий конкретную характеристику электронного документа, а также его обозначение.

**архивация/архивное хранение:** помещение электронного документа в условия, оптимальные для надёжного долговременного хранения с целью последующего обращения к нему в будущем.

## Виды электронных документов

### 1. По знаковой природе контента:

- текстовый электронный документ — электронный документ, контент-основу которого составляет читаемая информация преимущественно в виде слов;
- графический электронный документ — электронный документ, контент-основу которого составляет ▶

визуальное представление объектов/сущностей;

- звуковой (аудио) электронный документ — электронный документ, контент-основу которого составляет информация в форме, предназначенной для прослушивания.

## 2. По степени однородности:

- электронный документ, в котором объединены контент-элементы разной знаковой природы (например, текстово-визуальные, аудио-визуальные);
- электронный документ, в котором объединены контент-элементы разных динамических характеристик (например, текстово-звуковые, графико-звуковые).

## 3. По составу элементов:

- однородный (гомогенный) электронный документ — состоящий из контент-объектов одной знаковой природы;
- разнородный (гетерогенный) электронный документ — имеющий в своём составе контент-объекты различной знаковой природы.

## 4. По динамическим характеристикам:

- статический (неподвижный) электронный документ: статическое визуальное представление контент-элементов (например, фотография);
- динамический (движущийся) электронный документ: серия последовательного представления контент-элементов, которая приводит к эффекту движения / воспроизведению сигналов во времени (например, музыка, видео).

## 5. По количеству элементов:

- простой (односоставный) электронный документ — электронный документ, состоящий из единственного контент-элемента (например, фотография);
- составной (многосоставный) электронный документ — электронный документ, состоящий из более чем одного контент-элемента (например, слайд-шоу).

## 6. По структуре контента:

- плоский электронный документ — электронный документ с последовательной линейной связью контент-элементов;
- объёмный электронный документ — электронный документ с простран-

ственной нелинейной связью контент-элементов.

## 7. По процессам деривации (порождения электронного документа):

- впервые созданный электронный документ;
- электронный документ с изменённым контентом;
- электронный документ с изменённой формой/оформлением без качественного изменения контента электронного документа с изменённой знаковой природой;
- электронный документ с изменённым форматом (файла);
- электронный документ с изменённым размером (файла);
- дубликат электронного документа (полная идентичная копия, отличная по времени создания).

## 8. По происхождению контента:

- новый электронный документ (в том числе копия аналогового документа, ранее не представленная в электронной форме);
- редакция;
- версия;
- компиляция.

## 9. По производности:

- оригинальный, созданный впервые в электронной форме;
- копия электронного документа;
- конвертированный электронный документ (переведённый из одного формата в другой);
- трансформированный электронный документ (переведённый из одной знаковой системы в другую методом синтеза или анализа).

*Примечание — К наиболее распространённым способам трансформации относятся: автоматическое распознавание текста, речи, знаков; автоматический перевод; автоматический синтез речи.*

## Идентификация электронного документа

Идентификация производится на основе анализа блока постоянных характеристик электронного документа.

Определение постоянных характеристик электронного документа осуществляется в процессе его обработки, — комплекса документальных и информационных процессов, в основе которых лежит формально-содержательный анализ.

Результатом обработки является создание метаданных по определённой схеме. Метаданные могут формироваться полностью или частично автоматически при создании и/или автоматизированной обработке электронного документа.

Схема метаданных представляет собой набор элементов метаданных, предназначенных для конкретного практического применения, например, описания электронного документа. Определение значений самих элементов называется семантикой схемы. Содержание, присваиваемое элементам метаданных, называется значением. Схема метаданных в целом определяет имена элементов и их семантику, а также правила приведения значений (например, правила оформления, перечень допустимых значений) и синтаксические правила, определяющие кодировку элементов и их значений. Схема метаданных, в которой не установлены правила синтаксиса, называется синтаксически независимой, то есть метаданные могут кодироваться в любой определяемой синтаксической системе.

Выбор схемы метаданных зависит от условий, в которых осуществляется функционирование электронного документа, то есть от пользовательской среды, информационных процессов, целевого и пользовательского назначения, объектов и субъектов информационно-го взаимодействия.

В зависимости от конкретного назначения схема метаданных может модифицироваться с помощью расширения и профиля.

Расширение — добавление элементов к уже разработанной схеме для поддержки метаданных конкретного вида электронных документов или создание метаданных для конкретной группы пользователей.

Профиль используется для ограничения числа используемых элементов метаданных, для уточнения определения элементов при описании конкретного вида электронных документов, для определения значений, который может принимать тот или иной элемент.

Модификация схемы выполняет функцию разделения метаданных на универсальные (для всех электронных документов) и специальные (отдельные виды / ориентация на группы пользователей).

Универсальный набор метаданных для электронных документов содержит следующие блоки:

- данные об электронном документе как интеллектуальном объекте (сведения о создателе, заглавии, ответственности, содержании и языковой принадлежности);
- сведения об электронном документе как о физическом объекте (формат, размер, компоненты, адресная информация);
- характеристики жизненного цикла электронного документа (даты и иные параметры времени);
- данные о связи электронного документа с другими (сведения о версии, взаимном цитировании, отношениях «род-вид» и «часть-целое»);
- сведения о доступе к электронному документу (условия, права и правила использования).

Связь между метаданными и электронным документом, который они описывают, может осуществляться двумя способами:

- метаданные могут содержаться в записи, хранящейся отдельно от описываемого электронного документа;
- метаданные могут храниться непосредственно в теле электронного документа и извлекаться по мере необходимости (например, для построения поискового индекса).

### Типы метаданных:

- описательные метаданные (данные для поиска и идентификации контента электронного документа);
- структурные метаданные (данные о том, каким образом расположены и соединены элементы контента);
- административные метаданные (данные для управления и обеспечения сохранности электронного документа, включая технические и правовые аспекты).

### Базовые требования к форматам представления электронных документов

Цели определения базовых требований к форматам электронных документов:

1. Поддержка стратегического планирования в отношении цифрового контента электронных документов;
2. Обеспечение долгосрочного хранения и инвентаризации электронных до-

кументов, включая выявление инструментов и документации, необходимых для управления их контентом;

3. Разработка стратегии для поддержки стабильных форматов электронных документов в устойчивых технологических и пользовательских средах;

4. Разработка стратегии для конвертирования электронных документов неустойчивых форматов в целях сохранения их контента в неустойчивых технологических и пользовательских средах;

5. Определение политики сбора электронных документов;

6. Разработка политики развития и использования медиа-независимых форматов электронных документов;

7. Обеспечение экономической эффективности создания, хранения и работы с электронными документами;

8. Выявление форматов электронных документов, оптимальных для использования с конкретными видами контента;

9. Разработка механизмов технической защиты и миграции файлов электронных документов, а также стратегии поддержки парка аппаратного обеспечения (для аппаратно-зависимых форматов электронных документов при невозможности их конвертации) или стратегии портирования программного обеспечения.

### Классификация форматов электронных документов

Общепринятой единой классификации форматов в настоящее время не существует. Наиболее распространёнными основаниями для деления форматов на классы являются:

- расширение имени файла электронного документа (например, \*.doc или \*.docx, txt);
- тип информации интернет медиа-типов типа MIME (например, текст/HTML);
- цель использования (например, форматы электронных книг);
- служебное назначение или область применения (например, коммуникационные форматы ГИС);
- конкретные устройства (например, \*.raw цифровых камер);
- операционные системы и носители (например, \*.iso образ диска);
- алгоритм сжатия (например, jpg —

формат сжатия графических файлов);

- степени защиты контента (например, pdf-файлы с AdobeDRM)

*Примечание: различные форматы файлов различаются степенью детализации, один формат может накладываться на другой или использовать элементы других форматов.*

### Требования к форматам:

Требования, перечисленные ниже, являются универсальными, то есть относятся к цифровым форматам для всех видов электронных документов.

Для электронных документов отдельных видов могут дополнительно применяться специфические требования.

### открытость

Разработка открытых форматов электронных документов вызвана необходимостью создания условий для эффективного обмена электронными документами между всеми участниками процессов коммуникации на основе гармонизации способов и средств взаимодействия между информационными системами различных производителей и унификации существующих форматов электронных документов.

Соблюдение требования открытости формата представления электронного документа позволяет обеспечивать взаимодействие различных информационных систем; поддерживать возможность коллективной работы с электронными документами; предоставлять электронный документ в государственные органы, физическим и юридическим лицам; обеспечивать унифицированную обработку по стандартной схеме метаданных и надёжность долговременного хранения электронных документов.

Использование открытых форматов файлов необходимо для организации публичных сервисов; создания электронных документов, которые вводятся в публичное обращение; при проведении государственных тендеров на разработку или закупку программного обеспечения и т.п.

### Основные разработчики стандартов открытых форматов:

- ISO — International Organization for Standardization;
- ECMA — European Computer Manufacturers Association; Ecma International — European association for standardiz-

ing information and communication systems;

- NISO — National Information Standards Organization, a non-profit association accredited by the American National Standards Institute (ANSI);
- OASIS — Organization for the Advancement of Structured Information Standards;
- W3C — The World Wide Web Consortium;
- ITU-T — Telecommunication Standardization Sector of the International Telecommunication Union.

### распространённость

Распространённость предполагает максимально широкий круг пользователей электронных документов, представленных в данном формате, включая первичных создателей, распространителей, пользователей электронных документов.

Распространённость включает в себя использование формата:

- в качестве мастер-формата электронного документа;
- для доставки электронного документа конечным пользователям;
- как средство обмена между системами.

Распространённость формата замедляет его устаревание, и предполагает при разработке новых форматов электронных документов параллельное развитие инструментов их конвертации без дополнительных экономических затрат на миграцию, портирование и эмуляцию систем для работы с электронными документами устаревших распространённых форматов.

Свидетельством распространённости формата электронного документа является 1) поддержка его максимальным количеством конкурирующих программных средств для создания, просмотра, поиска, воспроизведения (вне зависимости от производителя и его лицензионной политики) и 2) выбор электронных документов в данном формате максимальным количеством конечных пользователей при наличии альтернативных форматов представления тех же электронных документов.

Распространённость формата электронных документов является основой для совместной работы организаций, осуществляющих библиотечно-информационную деятельность, органов научно-технической информации, органи-

заций, официально выпускающих в публичное обращение электронные документы с целью их массового использования.

### прозрачность

Прозрачность предполагает минимальное количество аппаратно-программных средств, которое требуется задействовать для того, чтобы контент электронного документа стал доступен конечному пользователю. Соответственно, форматы, которые предполагают различные виды пост-обработки электронных документов с целью оптимизации (особенно шифрование и сжатие), будут обладать меньшей прозрачностью, чем электронные документы в формате, где сжатие не использовалось.

Прозрачность усиливается, если текстовое содержание (в том числе текст метаданных, внедрённых в файлы электронных документов с графическим (нетекстовым) контентом) кодируется в стандартных кодировках (например, UNICODE в кодировке UTF-8) и хранится в естественном порядке чтения.



Существуют некоторые виды контента, которые даже при создании электронного документа не могут быть сохранены в несжатом виде. Для унификации требований прозрачности необходимо определение коэффициента прозрачности форматов, используемых для разных видов электронных документов, для разных целей использования и разных видов деятельности.

### внедрение метаданных

Форматы, позволяющие внедрение метаданных, то есть создание и хранение их непосредственно в теле документа, являются наиболее предпочтительными для создания/перевода в них электронных документов. При этом оптимальным является использование форматов с авто-метаданными, то есть тех, в которых часть значений полей автоматически формируется программными

средствами. Максимальная полнота авто-метаданных, документирующих жизненный цикл электронного документа с момента его создания, облегчает его поддержку в долгосрочной перспективе и обеспечивает наименьшую уязвимость во всех неинформационных процессах, связанных с обеспечением сохранности.

### Будущее электронного документа: заключение

В условиях современного динамического развития общества электронные документы становятся таким же стратегическим ресурсом, как традиционные материальные и энергетические ресурсы. Становление современного информационного общества немислимо без использования информационных ресурсов в электронном виде. В свою очередь, трудно использовать электронные информационные ресурсы без стандартизации процессов перевода аналоговых документов в электронную форму, создания электронных информационно-насыщенных систем.

Электронные документы приобретают новый статус, при котором реализуется качественно иной уровень производства, хранения, организации и распространения самой разнообразной информации, что обеспечивает её ещё более широкое распространение и эффективное использование.

При разработке стандартов на электронные ресурсы необходима координация международных и национальных стандартов с участием российских специалистов и технических комитетов по стандартизации.

С авторами можно связаться:

**shorin@nlr.ru**  
**barysh@nlr.ru**

Статья посвящена вопросам стандартизации при регистрации, обработке, хранении, копировании электронных документов в библиотеках.

**Электронные документы, библиотечные стандарты**

The article is devoted to standardization in the acquisition, processing, storage and copying of electronic documents in libraries.

**Electronic documents, library standards**