

УДК 004.021, 004.457

## **АНАЛИЗ МЕТОДОВ АВТОМАТИЗИРОВАННОЙ ГЕНЕРАЦИИ ДОКУМЕНТОВ ДЛЯ РЕШЕНИЯ ЗАДАЧИ СОЗДАНИЯ ГЕНЕРАТОРА КНИГИ АБСТРАКТОВ ДЛЯ КОНФЕРЕНЦИЙ ЛИТ ОИЯИ**

**Ильина Анна Владимировна<sup>1</sup>, Кореньков Владимир Васильевич<sup>2</sup>**

<sup>1</sup>Студент;

ГБОУ ВО МО «Университет «Дубна»,  
Институт системного анализа и управления;  
141980, Московская обл., г. Дубна, ул. Университетская, 19;  
e-mail: ttrnk75@yandex.ru.

<sup>2</sup>Директор, доктор технических наук, профессор;

Объединенный институт ядерных исследований,  
Лаборатория информационных технологий;  
141980, Московская обл., г. Дубна, ул. Жолио-Кюри, 6;  
ГБОУ ВО МО «Университет «Дубна»,  
Институт системного анализа и управления;  
141980, Московская обл., г. Дубна, ул. Университетская, 19;  
e-mail: korenkov@cv.jinr.ru.

*В работе рассматриваются методы автоматизированной генерации документов, выделены достоинства и недостатки каждого из них, приведено обоснование выбора метода генерации документов на основе шаблонов формата DOCX для решения задачи автоматизированной генерации книги абстрактов для конференций ЛИТ ОИЯИ.*

Ключевые слова: генерация документов, конференция, абстракт, программирование, DOCX, LaTeX.

## **ANALYSIS OF METHODS OF AUTOMATED GENERATION OF DOCUMENTS FOR THE CREATION OF A GENERATOR OF THE BOOK OF ABSTRACTS FOR LIT JINR CONFERENCES**

**Irina Anna<sup>1</sup>, Korenkov Vladimir<sup>2</sup>**

<sup>1</sup>Student;

Dubna State University,  
Institute of the system analysis and management;  
141980, Dubna, Moscow reg., Universitetskaya str., 19;  
e-mail: ttrnk75@yandex.ru.

<sup>2</sup>Director, Doctor of Technical Science, professor;

Joint institute for nuclear researches,  
Laboratory of Information Technologies;  
141980, Moscow reg., Dubna, Joliot-Curie, 6;  
Dubna State University,  
Institute of the system analysis and management;  
141980, Dubna, Moscow reg., Universitetskaya str., 19;  
e-mail: korenkov@cv.jinr.ru.

*The paper deals with the methods of automated document generation, highlighted the advantages and disadvantages of each of them, the rationale for the choice of the method of generating documents based on docx format templates to solve the problem of automated generation of abstracts for conferences LIT JINR.*

Keywords: document generation, conference, abstract, programming, DOCX, LaTeX.

## Введение

Проблема автоматизированной генерации документов особенно актуальна для систем, обрабатывающих огромное количество данных, которое в результате должно представлять из себя некоторые конкретизированные отчеты, оформленные специальным образом. Зачастую количество данных настолько большое, а формат выходного файла столь нетривиален, что вручную решить такую задачу не представляется возможным либо решение займет колоссальное количество времени.

Задачу автоматизированной генерации документов нельзя считать тривиальной и имеющей некое общее решение. Это связано в первую очередь с особенностями стилей конкретных генерируемых документов, а также с форматом входных данных, содержащих ключевую для генерации информацию. Однако существуют некоторые общие методы и подходы, составляющие основу алгоритмизации данного процесса.

Задача автоматизированной генерации документов актуальна и для ОИЯИ – международной межправительственной научно-исследовательской организации, перед которой задача автоматизации генерации различного рода документов стоит особенно остро.

### 1. Структура книги абстрактов

Одной из задач является задача генерации некоторого документа для конференций, которые постоянно организуются в институте. Рассмотрим этот документ более подробно.

Организация каждой конференции включает в себя создание книги абстрактов, определенной и стилизованной специальным образом (см. рис. 1). Книга содержит большое количество страниц, содержащих в себе данные о конференции и выступающих (см. рис. 2), представленные в исходном файле формата *XML*, который генерируется в информационной системе организации конференций *Indico* (см. рис. 3). Существует проблема автоматизированной генерации такой книги, подходы к решению которой будут рассмотрены в данной работе.

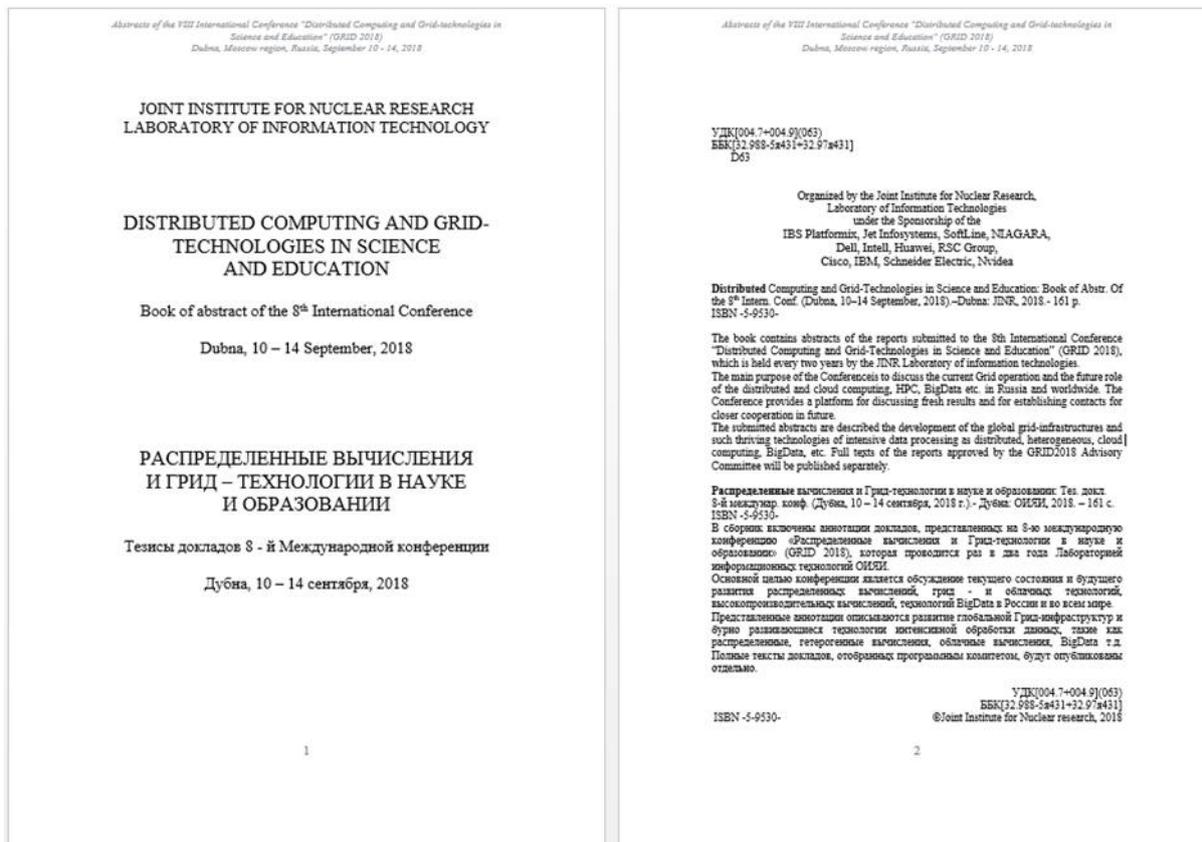


Рис. 1. Титульные страницы книги абстрактов

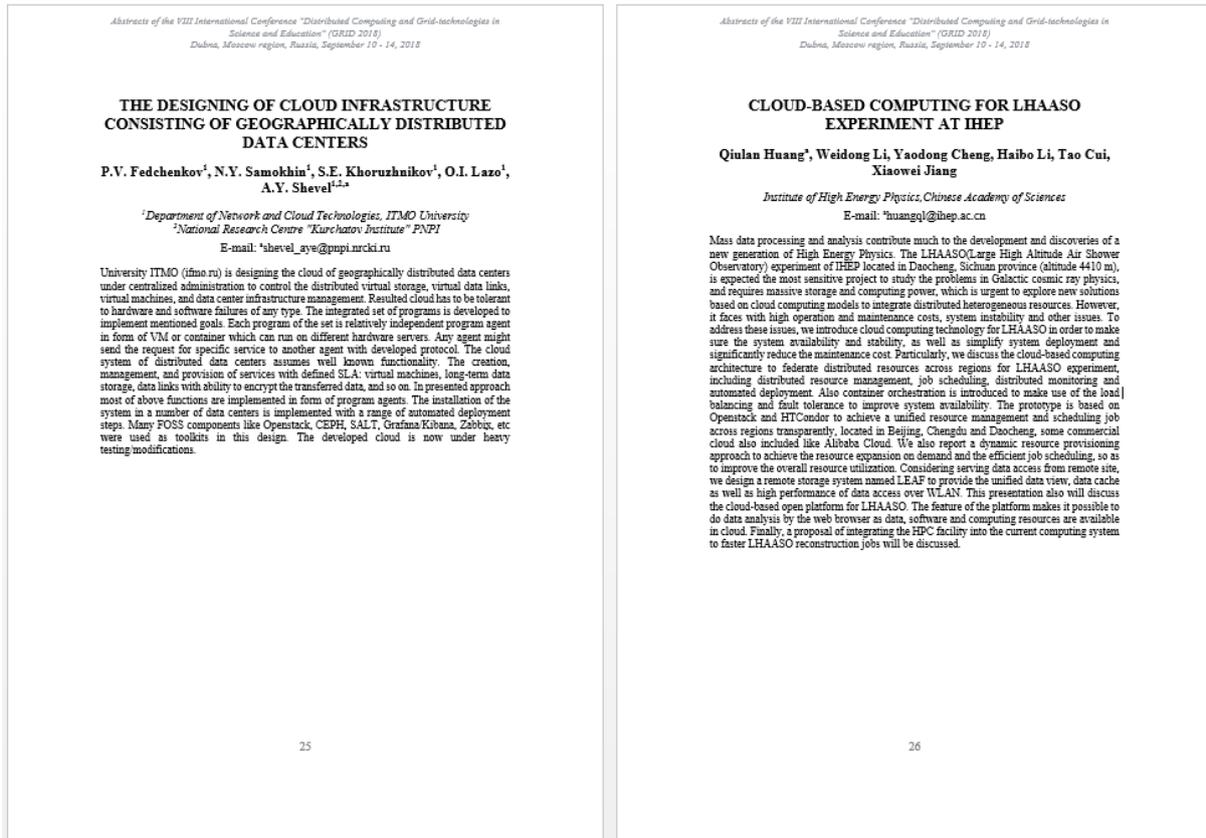


Рис. 2. Содержимое книги абстрактов

```

▼<AbstractBook>
  ▼<Conference>
    The 8th International Conference "Distributed Computing and Grid-technologies in Science and Education" (GRID 2018)
  </Conference>
  ▼<abstract>
    <Id>180</Id>
    ▼<Title>
      Грид и облачная инфраструктура дата-центра Института Физики НАН Азербайджана
    </Title>
    ▼<Content>
      Основными направлениями развития дата-центра Института физики НАН Азербайджана являются грид и облачные технологии. Пользуясь тем что грид-сегмент дата-центра интегрирован в инфраструктуру EGI/WLCG в качестве грид-сайта уровня Tier3, пользователи дата-центра получают возможность участвовать в таких международных проектах как ATLAS (CERN). Сотрудничество с международными научными центрами ОИЯИ, ЦЕРН в области информационных технологий способствует эффективному развитию дата центра, а также помогает решать задачи пользователей в таких научных областях, как физика высоких энергий, физика твердого тела и т. д.
    </Content>
    ▼<field id="content">
      Основными направлениями развития дата-центра Института физики НАН Азербайджана являются грид и облачные технологии. Пользуясь тем что грид-сегмент дата-центра интегрирован в инфраструктуру EGI/WLCG в качестве грид-сайта уровня Tier3, пользователи дата-центра получают возможность участвовать в таких международных проектах как ATLAS (CERN). Сотрудничество с международными научными центрами ОИЯИ, ЦЕРН в области информационных технологий способствует эффективному развитию дата центра, а также помогает решать задачи пользователей в таких научных областях, как физика высоких энергий, физика твердого тела и т. д.
    </field>
    <field id="summary"/>
    ▼<PrimaryAuthor>
      <FirstName>Aleksey</FirstName>
      <FamilyName>Bondyakov</FamilyName>
      <Email>aleksey@jinr.ru</Email>
      <Affiliation>JINR (Joint Institute For Nuclear Research)</Affiliation>
    </PrimaryAuthor>
    ▼<Speaker>
      <FirstName>Aleksey</FirstName>
      <FamilyName>Bondyakov</FamilyName>
      <Email>aleksey@jinr.ru</Email>
      <Affiliation>JINR (Joint Institute For Nuclear Research)</Affiliation>
    </Speaker>
    <ContributionType>Sectional reports</ContributionType>
    <Track/>
  </abstract>
▼<abstract>
  ...

```

Рис. 3. Входной файл формата XML

## 2. Постановка задачи

Цель: автоматизировать генерацию документа абстрактов для конференций, проходящих в ЛИТ ОИЯИ.

Исходные данные:

1. Входной документ формата *XML*, содержащий данные абстрактов и конференции.
2. Пример выходного файла в формате *DOCX*.
3. Справочная литература по языкам объектно-ориентированного программирования высокого уровня.

Модельные представления:

1. Автоматический генератор документов должен представлять из себя программное решение, реализованное под ОС *Linux* и *Windows* (далее – программный продукт (ПП)).
2. ПП должен принимать на вход документ в формате *XML*, автоматический сгенерированный веб-сервисом для организации конференций *Indico*.
3. ПП должен генерировать конечный документ в формате *DOCX* (преимущественно) или *PDF*.

Ожидаемый результат: ПП, позволяющий автоматически сгенерировать книгу абстрактов для конференций, проходящих в ЛИТ ОИЯИ, в формате *DOCX* или *PDF* из документа формата *XML*.

Критерий оценки результата: стиль и форматирование автоматически сгенерированного документа должны быть идентичны стилю и форматированию документа, содержащего пример выходного файла.

Рассмотрим некоторые общие подходы к автоматизированной генерации документов.

## 3. Метод генерации документов с помощью системы *LaTeX*

*LaTeX* – это издательская система на базе *TeX* для верстки текстов с формулами. Сам по себе *TeX* представляет собой специализированный язык программирования, на котором пишутся издательские системы, используемые на практике. Точнее говоря, каждая издательская система на базе *TeX* представляет собой пакет макроопределений (макропакет) этого языка (см. рис. 4) [1].

Достоинства метода:

1. Напечатанный документ выглядит как в книге.
2. Указав с помощью простых средств логическую структуру текста, автор может не вникать в детали оформления, причем эти детали при необходимости нетрудно изменить.
3. Высокое качество и гибкость верстки абзацев и математических формул.
4. Система позволяет генерировать один большой документ из многих маленьких.

Недостатки метода:

1. Наличие определенных специальных навыков верстки и работы в системе *LaTeX*.
2. Работа с исходным текстом и просмотр того, как текст будет выглядеть на печати, – разные операции (не поддерживается свойство *WYSIWYG* (*What you see is what you get*)) [2].
3. *TeX* плохо приспособлен для верстки страниц со сложным взаимодействием текста и графики.
4. Чтобы изменить какую-либо информацию в документе, необходимо править исходный *.tex*-файл и осуществлять генерацию конечного документа заново.

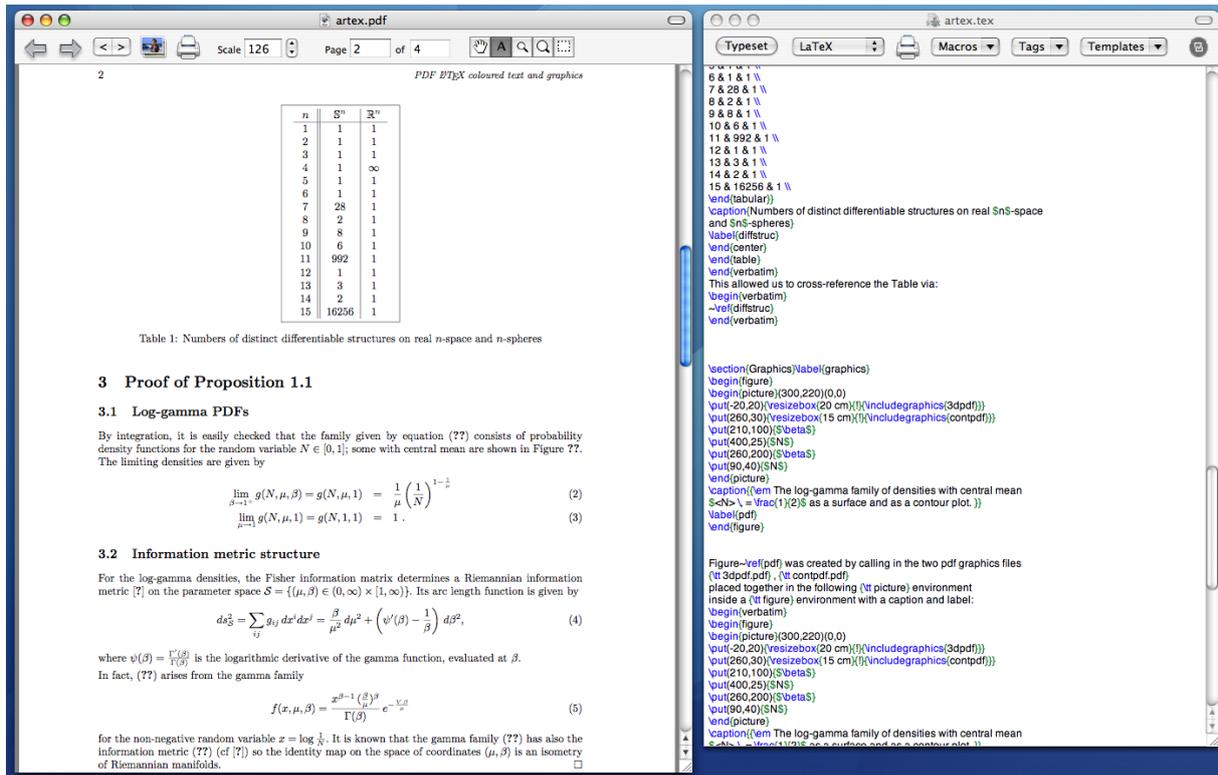


Рис. 4. Пример верстки документа в системе LaTeX

#### 4. Метод генерации документов на основе шаблонов формата DOCX

В основе метода лежит достаточно быстрый процесс создания текстового шаблона, содержащего некоторые переменные (или поля) (см. рис. 5), которые будут автоматически заменены необходимыми конкретными данными с помощью специальных библиотек по генерации документов из шаблонов. Недостатком этого метода можно считать необходимое написание программного продукта с подключением соответствующей библиотеки. Такой продукт должен уметь считать информацию из входного файла и с помощью библиотеки заменить ей соответствующие переменные. К недостаткам также следует отнести и тот факт, что при некорректной работе соответствующей библиотеки конечный документ будет так же некорректно сформирован (существует зависимость от разработчиков стороннего ПО) [3]. Данный подход является наиболее практичным для генерации документов с большими объемами информации, которые должны быть оформлены специальным образом, поскольку:

1. Поддерживает свойство *WYSIWYG* (*What you see is what you get*) [2].
2. Не требует написания *tex*- или *html*-представлений.
3. Конечный документ имеет редактируемый формат *DOCX*.

В работе были рассмотрены некоторые общие подходы к автоматизированной генерации документов, а также выделены достоинства и недостатки каждого из них.

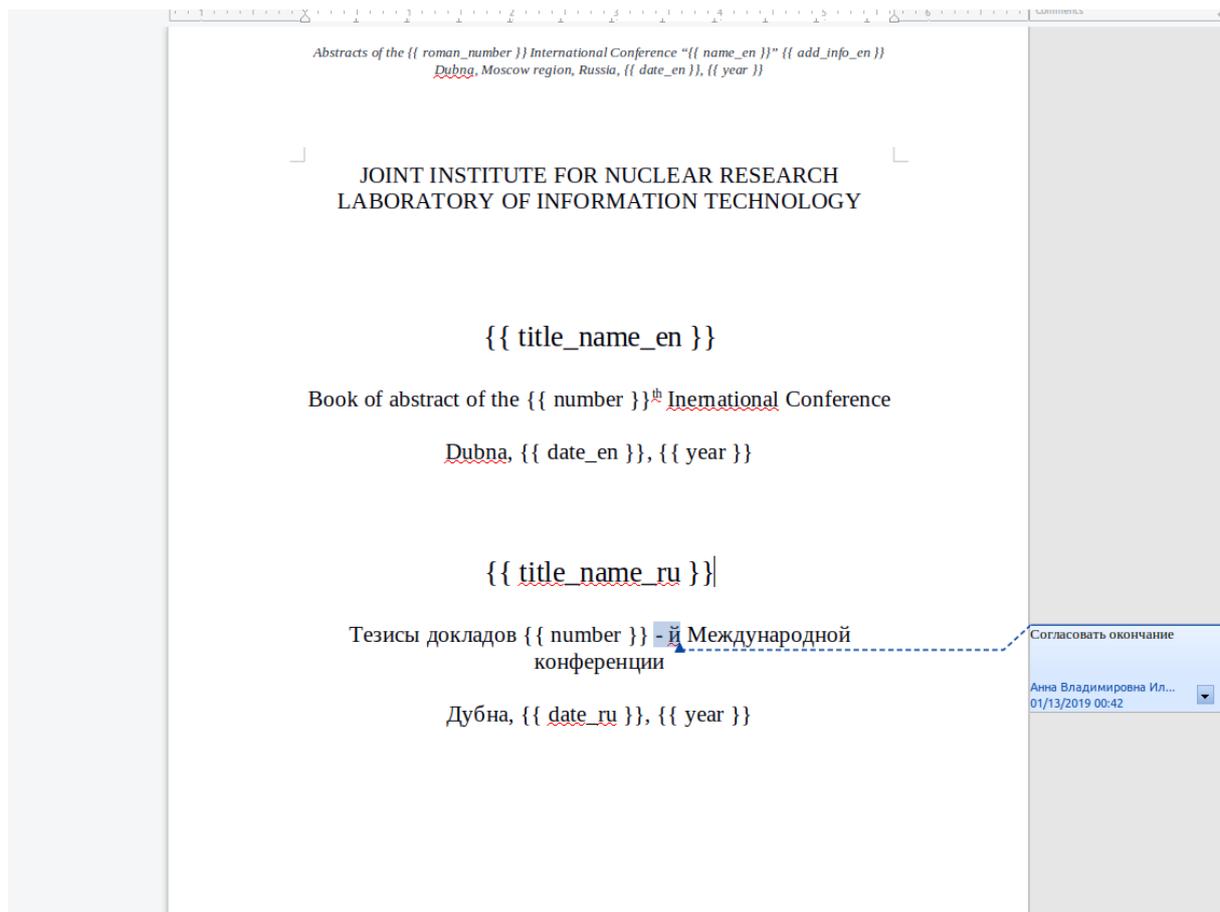


Рис. 5. Пример шаблона для генерации документа

## Заключение

Наиболее практичным подходом для генерации книги абстрактов выбран метод генерации документа на основе шаблона формата *DOCX* с написанием соответствующего ПП, поскольку:

1. Позволяет достичь поставленной цели.
2. Не требует специальных навыков верстки и программирования для конечных пользователей.
3. Конечный документ имеет редактируемый формат *DOCX*, что гарантирует возможность правки информации в случае её изменения.
4. Документ формата *DOCX* легко экспортировать в документ формата *PDF* и тут же отправить в печать.

## Список литературы

1. TeX – система для вёрстки текстов с формулами. – [Электронный ресурс]. URL: [http://www.wikiznanie.ru/ru-wz/index.php/TeX\\_-\\_система\\_для\\_вёрстки\\_текстов\\_с\\_формулами](http://www.wikiznanie.ru/ru-wz/index.php/TeX_-_система_для_вёрстки_текстов_с_формулами).
2. WYSIWYG. – [Электронный ресурс]. URL: <https://ru.wikipedia.org/wiki/WYSIWYG>.
3. Автоматизируем генерацию docx документов по шаблону с помощью python и docxtempl. – [Электронный ресурс]. URL: <http://tut-admin.ru/2017/05/31/автоматизируем-генерацию-docx-документов/>.