



Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Московский государственный технический университет  
имени Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ ИУ «Информатика и системы управления»

КАФЕДРА ИУ7 «Программное обеспечение ЭВМ и информационные технологии»

**РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА**  
***К НАУЧНО-ИССЛЕДОВАТЕЛЬСКОЙ РАБОТЕ***  
***НА ТЕМУ:***  
***«Методы выделения составных частей научного***  
***текста»***

Студент      **ИУ7-74Б**

\_\_\_\_\_  
(Подпись, дата)      **К. А. Рунов**  
(И.О.Фамилия)

Руководитель

\_\_\_\_\_  
(Подпись, дата)      **Ю. В. Строганов**  
(И.О.Фамилия)

Консультант

\_\_\_\_\_  
(Подпись, дата)      **Ю. И. Бутенко**  
(И.О.Фамилия)

Рекомендованная руководителем НИР оценка: \_\_\_\_\_

## РЕФЕРАТ

Отчет 18 с., 4 рис., 1 табл., 11 источн., 1 прил.

DOCUMENT LAYOUT ANALYSIS, НАУЧНО-ТЕХНИЧЕСКИЙ ТЕКСТ, CONNECTED COMPONENT ANALYSIS, PROJECTION PROFILE ANALYSIS, RLSA, МАШИННОЕ ОБУЧЕНИЕ, КЛАССИФИКАЦИЯ

Цель работы — классификация методов выделения составных частей научного текста.

В данной работе был проведен анализ предметных областей научно-технических текстов и анализа структуры документов. Был проведен обзор существующих методов выделения составных частей научного текста. Были сформулированы критерии сравнения описанных методов и была проведена классификацию описанных методов по сформулированным критериям.

# СОДЕРЖАНИЕ

<b>РЕФЕРАТ</b>	<b>3</b>
<b>ВВЕДЕНИЕ</b>	<b>5</b>
<b>1 Анализ предметной области</b>	<b>6</b>
1.1 Анализ структуры документов (DLA) . . . . .	6
1.1.1 Этап предварительной обработки . . . . .	6
1.1.2 Этап анализа макета документа . . . . .	7
1.2 Типы макетов документов . . . . .	8
1.3 Структура научно-технического текста . . . . .	9
<b>2 Формализация задачи</b>	<b>10</b>
<b>3 Описание существующих методов</b>	<b>11</b>
3.1 Анализ связных компонент (ССА) . . . . .	11
3.2 Анализ проекционного профиля (РРА) . . . . .	12
3.3 Алгоритм размазывания по длине серии (RLSA) . . . . .	12
3.4 Методы на основе машинного обучения . . . . .	13
3.5 Гибридные методы на основе РРА и ССА . . . . .	14
<b>4 Классификация существующих методов</b>	<b>14</b>
<b>ЗАКЛЮЧЕНИЕ</b>	<b>15</b>
<b>СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ</b>	<b>17</b>
<b>ПРИЛОЖЕНИЕ А</b>	<b>18</b>

# ВВЕДЕНИЕ

Анализ структуры документов (Document Layout Analysis, DLA) играет ключевую роль в обработке научно-технических текстов. Такие документы обладают четкой структурой, включающей заголовки, авторов, аннотации, разделы, формулы, таблицы, графики и рисунки [1, 2, 3, 4]. Выявление этих элементов и их логических связей позволяет не только упрощать индексирование и поиск информации, но и улучшать автоматическую обработку текстов, включая аннотирование, реферирование и анализ содержимого.

Документ можно представить в виде иерархии физических модулей (страницы, колонки, абзацы, строки, слова, изображения) или логических модулей (заголовки, авторы, аффилиации, аннотации, разделы, библиография) [5].

Эффективный анализ структуры документов обеспечивает удобную навигацию по тексту, облегчает его разметку и позволяет быстро извлекать необходимые сведения [5].

Целью данной работы является классификация методов выделения составных частей научного текста.

Для достижения поставленной цели необходимо решить следующие задачи:

- провести анализ предметных областей анализа структуры документов и научно-технических текстов;
- провести обзор существующих методов выделения составных частей научного текста;
- сформулировать критерии сравнения описанных методов;
- провести классификацию описанных методов по сформулированным критериям.

# 1 Анализ предметной области

## 1.1 Анализ структуры документов (DLA)

Анализ структуры документов (Document layout analysis, DLA) — процесс сегментирования входного изображения документа на однородные компоненты, такие как блоки текста, рисунки, таблицы, графики и т.д., и их классификации [6].

В общем случае анализ структуры документа делится на два взаимосвязанных процесса: физический и логический анализ. Целью физического анализа является выявление структуры документа и определение границ его однородных областей. Целью логического анализа является разметка обнаруженных областей. Выявленные области классифицируются как элементы документа — рисунки, заголовки, абзацы, логотипы, подписи и другие. [2]

Процесс анализа структуры документов состоит из двух основных этапов — этапа предварительной обработки и этапа анализа макета документа [2, 5]. На рисунке ниже приведена схема процесса анализа структуры документов.

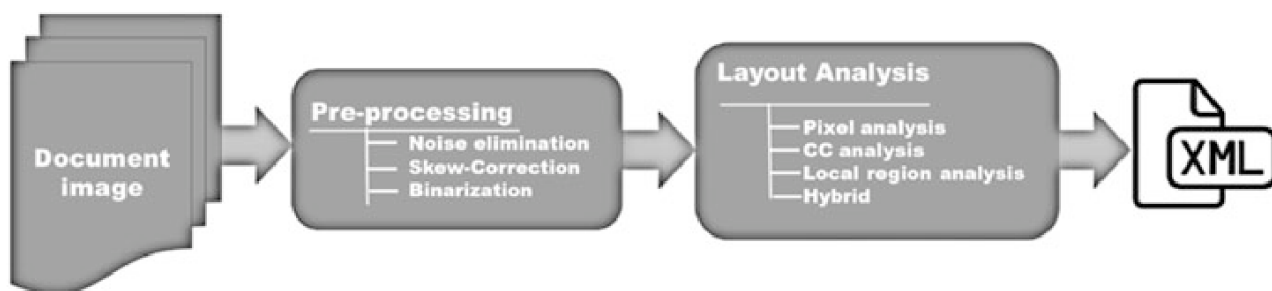


Рисунок 1 – Схема процесса анализа структуры документов [5]

### 1.1.1 Этап предварительной обработки

Этап анализа макета документа в любом методе анализа структуры документов (далее DLA) часто основывается на определённых предположениях о входных изображениях, таких как отсутствие шума, бинаризация, отсутствие наклона текста или все перечисленные факторы [2, 5].

Цель этапа предварительной обработки — преобразовать входное изображение в соответствии с требованиями этапа анализа макета документа кон-

кретного метода [2, 5].

В общем случае на этом этапе используются одна или несколько процедур предварительной обработки, таких как бинаризация, выравнивание и улучшение изображения [2, 7].

### **1.1.2 Этап анализа макета документа**

Анализ макета документа включает в себя определение границ и типов составляющих областей входного изображения документа. Процесс определения границ областей документа называется сегментацией областей документа, а классификация найденных областей по их типу — классификацией областей документа. [5]

Существуют три типа стратегий анализа макета документа: снизу вверх (bottom-up), сверху вниз (top-down) и гибридная (hybrid).

По стратегии снизу вверх (bottom-up) параметры анализа часто вычисляются на основе исходных данных. Анализ макета документа начинается с небольших элементов, таких как пиксели или связанные компоненты. Затем однородные элементы объединяются, создавая более крупные области. Процесс продолжается, пока не будут достигнуты заранее определённые условия остановки.

По стратегии сверху вниз (top-down) анализ макета документа начинается с крупных областей, например, на уровне всего документа. Затем эта большая область разбивается на более мелкие, такие как колонки текста, на основе определённых правил однородности. Анализ сверху вниз прекращается, когда дальнейшее разбиение областей становится невозможным или достигаются условия остановки.

Гибридная стратегия (hybrid) представляет собой комбинацию обеих стратегий (снизу вверх и сверху вниз). [2]

После сегментации областей происходит их классификация с помощью различных алгоритмов, в результате чего формируется логическая структура документа.

По завершении данного этапа извлеченные геометрическая и логическая структуры сохраняются для последующей реконструкции. Для этого, как правило, используется иерархическая древовидная структура данных. [5]

## 1.2 Типы макетов документов

Макеты документов могут иметь различные структуры. Печатные документы можно разделить на шесть типов [8]: прямоугольные, Манхэттенские, не-Манхэттенские, многоколоночные Манхэттенские, с горизонтальным наложением и с диагональным наложением.

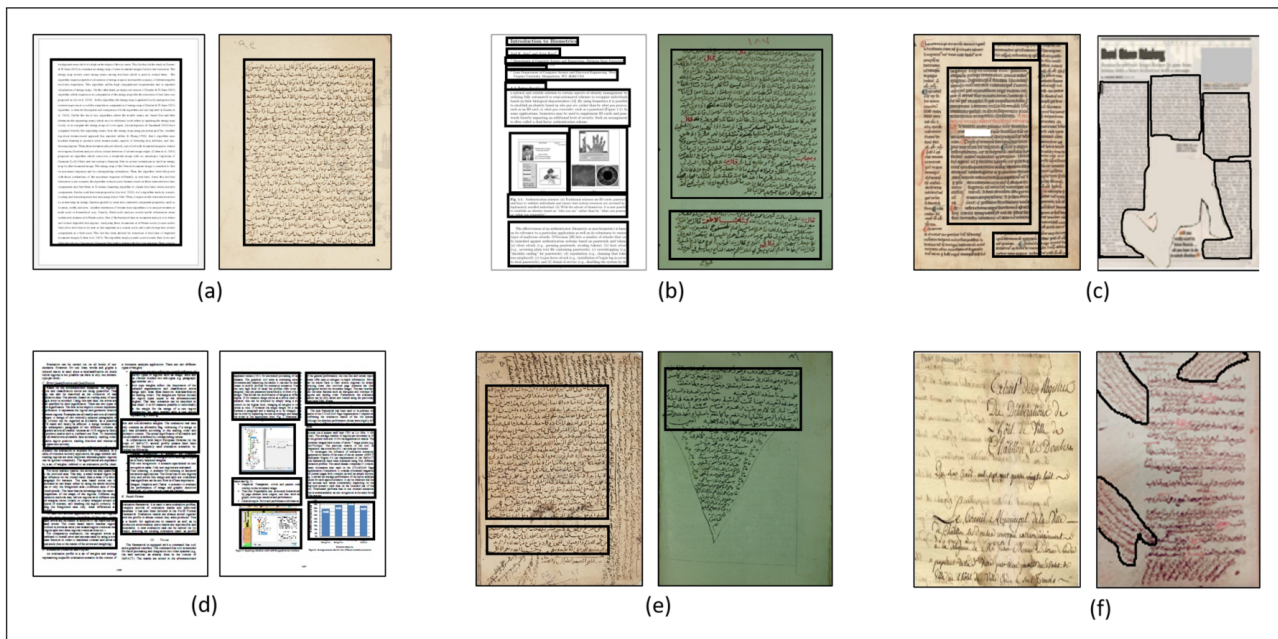


Рисунок 2 – Макеты документов: (а) Стандартный (прямоугольный), (b) Манхэттенский, (c) Не-Манхэттенский, (d) Многоколоночный Манхэттенский, (e) Произвольный (сложный), (f) С горизонтальным и диагональным наложением. [2]

На рисунке выше показаны примеры описанных типов макетов документов:

- Стандартный макет характеризуется большими прямоугольными текстовыми блоками, расположенными в одной или нескольких колонках, при этом каждая колонка содержит по одному абзацу.
- Если документ содержит несколько абзацев в колонках, его можно отнести к Манхэттенскому макету. Примеры таких документов — научно-технические статьи, журналы и другие.
- Не-Манхэттенские макеты включают зоны непрямоугольной формы.

- Макеты с наложением содержат элементы, такие как текст, который перекрывает другие элементы документа. Наложение может возникать, например, из-за просвечивания (см. Рисунок 2(f)).
- Документы с произвольными (или сложными) макетами могут включать рукописный и/или печатный текст, содержащий различные стили, типы и размеры шрифтов.

Таким образом, документы, содержащие научно-технические тексты, обычно используют Манхэттенский макет.

### **1.3 Структура научно-технического текста**

Научно-технический текст обычно [9, 10, 11] следует четко определенному шаблону и состоит из следующих частей:

- 1) Название;
- 2) Информация об авторах;
- 3) Аннотация и ключевые слова;
- 4) Введение;
- 5) Основная часть (кроме текста содержащая в том числе таблицы, рисунки, графики, уравнения);
- 6) Заключение;
- 7) Ссылки на литературу.

На рисунке ниже показана структура научной статьи.





Рисунок 3 – Структура научной статьи

Зная структуру научного текста и его основные части можно перейти к формализации задачи выделения составных частей научного текста.

## 2 Формализация задачи

Путь  $D$  — документ, представленный в виде набора изображений, содержащих текст, формулы, таблицы, рисунки и прочие структурные элементы.

Документ  $D = \{P_1, P_2, \dots, P_n\}$  состоит из страниц  $P_1, P_2, \dots, P_n$ , а каждая страница  $P_i$  в свою очередь содержит множество объектов  $O_{i,1}, O_{i,2}, \dots, O_{i,m}$ . Объект  $O_{i,j}$  — кортеж  $(G_{i,j}, T_{i,j})$ , где  $G_{i,j}$  — геометрические свойства объекта ( $G_{i,j} = (x_{i,j}, y_{i,j}, w_{i,j}, h_{i,j})$ , где  $(x_{i,j}, y_{i,j})$  — координаты верхнего левого угла,  $w_{i,j}$  — ширина,  $h_{i,j}$  — высота объекта),  $T_{i,j}$  — текстовые свойства объекта (строка символов, соответствующая текстовому содержимому объекта, возможно, пустая).

Требуется построить отображение

$$F : D \rightarrow \{(O_{i,j}, C_{i,j})\},$$

где каждому объекту  $O_{i,j}$  ставится в соответствие класс  $C_{i,j} = C_{i,j}(O_{i,j})$ ,  $C_{i,j} \subseteq \{\text{Название, Информация об авторах, Аннотация и ключевые слова, Введение, Основная часть, Таблица, Рисунок, График, Уравнение, Заключение, Ссылки на литературу}\}$ .

Поставленную задачу можно решить, разбив на две подзадачи и решив каждую подзадачу соответственно: первая подзадача — нахождение объектов на страницах и выявление их геометрических и текстовых свойств, вторая подзадача — классификация найденных объектов (определение  $C_{i,j}$  для каждого объекта  $O_{i,j}$ ). Решением первой подзадачи является построение отображений  $P_i \rightarrow \{O_{i,j}\}$ , решением второй подзадачи является построение отображений  $O_{i,j} \rightarrow C_{i,j}$ .

Далее будут рассмотрены существующие методы, позволяющие решить поставленную задачу.

### 3 Описание существующих методов

#### 3.1 Анализ связных компонент (ССА)

Методы на основе связных компонент (Connected Component Analysis, ССА) анализируют и объединяют связные компоненты для формирования однородных областей.

Определение начальных компонент, которые впоследствии объединяются, происходит, как правило, следующим образом. Изображение проходит стадию бинаризации (преобразование к черно-белому формату и назначение каждому пикселю интенсивности 0 или 1), после чего смежные пиксели объединяются на основе 4- или 8-связности. При 4-связности два пикселя считаются связными, если они расположены друг за другом по горизонтали или вертикали. При 8-связности два пикселя считаются связными, если они являются 4-связными, либо расположены друг за другом по диагонали.

После определения начальных компонент, компоненты объединяются в

однородные области путем применения специальных алгоритмов. В качестве таких алгоритмов могут выступать, например, преобразование Хафа (Hough transform) или алгоритм К-ближайших соседей (K-nearest neighbor, KNN) [5].

Далее происходит классификация однородных областей. Для классификации области могут использоваться эвристические алгоритмы (классификация на основе ширины штриха, размера или формы компонента) и алгоритмы на основе машинного обучения.

Методы на основе связных компонент могут работать в условиях скошенного текста при условии, что межстрочный интервал меньше пробела между абзацами [5].

### **3.2 Анализ проекционного профиля (PPA)**

Суть методов на основе анализа проекционного профиля (Projection Profile Analysis, PPA) заключается в следующем. Пиксели документа проецируются на вертикальную и горизонтальную ось, после чего строятся гистограммы распределения пикселей. Далее анализируются пики и впадины на гистограммах. Впадины на вертикальном профиле указывают на пробелы между колонками текста. Впадины на горизонтальном профиле указывают на пробелы между строками или блоками текста.

На основе проведенного анализа можно разделить документ на логические компоненты — текстовые блоки, заголовки, таблицы, изображения.

Данные методы работают только с Манхэттенскими макетами документов и чувствительны ко скошенности текста в документе [5].

### **3.3 Алгоритм размазывания по длине серии (RLSA)**

Методы на основе RLSA преобразуют бинарное изображение документа путем «размазывания» черных пикселей горизонтально и/или вертикально для формирования однородных областей.

На рисунке ниже можно видеть пример работы RLSA.

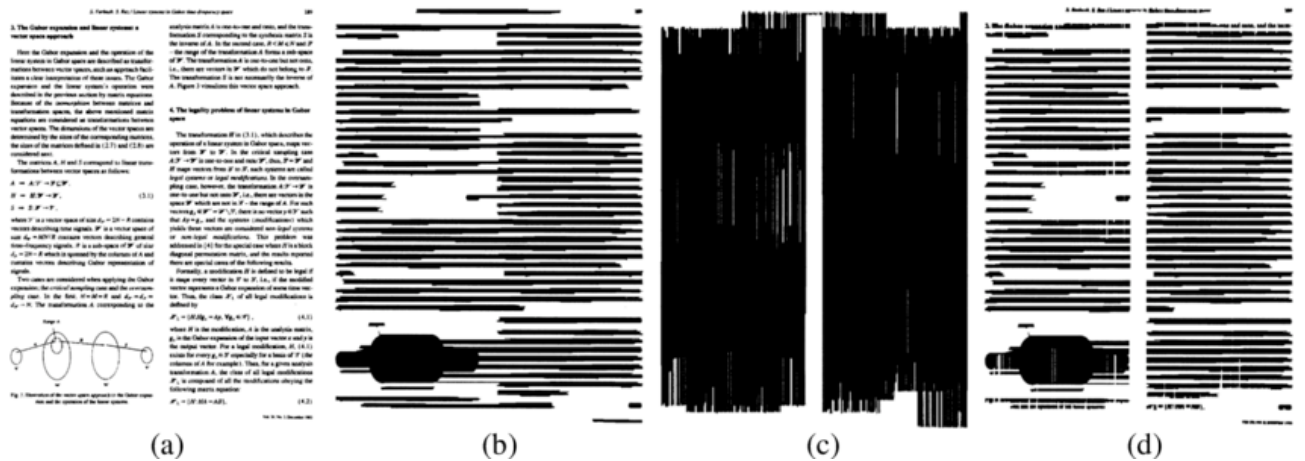


Рисунок 4 – Пример работы алгоритма RLSA. К начальному изображению документа (a) применяется горизонтальный (b) и вертикальный (c) RLSA, после чего в результате применения операции II к изображениям (b) и (c) формируется (d)

Для классификации полученных областей также используются эвристические алгоритмы и алгоритмы на основе машинного обучения.

Данные методы, как и методы на основе анализа проекционного профиля, работают преимущественно с Манхэттенскими макетами документов и чувствительны к скошенному тексту [5].

### 3.4 Методы на основе машинного обучения

Методы, не основанные на глубоком обучении, используют простые архитектуры нейросетей для обучения. Анализ с использованием нейросети происходит на трех уровнях: уровне пикселей, уровне блоков текста и уровне страниц.

Методы на основе машинного обучения в области анализа макетов документов страдают от несбалансированности данных и отсутствия контекстной информации. Если модель обучалась на документах, в наборе данных для обучения количество текстовой и фоновой информации сильно превосходит количество информации о рисунках и графиках. В связи с этим обученная модель может склоняться в сторону текстовых или фоновых пикселей. [2]

Обучение моделей лишь на основе информации о пикселях чревато потерей контекстной информации. Поэтому при обучении на уровнях блоков текста и страниц прибегают к использованию методов извлечения признаков для создания более надежных моделей. [2]

Методы на основе машинного обучения работают с любыми макетами документов и не чувствительны к скошенному тексту.

### 3.5 Гибридные методы на основе PPA и CCA

Гибридные методы на основе анализа проекционного профиля и связанных компонент используют работают следующим образом. Начальные компоненты определяются применением метода из PPA, после чего происходит их уточнение применением методов из CCA.

Такой комбинированный подход позволяет лучше сегментировать текст, чем каждый метод по отдельности.

## 4 Классификация существующих методов

Для сравнения рассмотренных методов можно выделить следующие критерии:

- Стратегия анализа макета документа;
- Скорость работы — требования метода к вычислительным ресурсам;
- Гибкость — способность метода адаптироваться к различным типам макетов документов;
- Устойчивость — способность метода адаптироваться к шумам и искажениям текста.

Ниже приведена таблица со сравнительным анализом рассмотренных методов.

Таблица 1 – Классификация методов DLA

Метод	Стратегия	Скорость	Гибкость	Устойчивость
CCA	Снизу вверх	2	2	3
PPA	Сверху вниз	2	3	3
RLSA	Сверху вниз	1	3	3
ML	Снизу вверх	3	1	1
PPA + CCA	Гибридный	2	3	2

## ЗАКЛЮЧЕНИЕ

В ходе данной научно-исследовательской работы был проведен анализ предметных областей научно-технических текстов и анализа структуры документов, проведен обзор существующих методов выделения составных частей научного текста, были сформулированы критерии сравнения описанных методов и была проведена классификацию описанных методов по сформулированным критериям.

Таким образом, все задачи для достижения цели данной работы были решены, и цель работы — классификация методов выделения составных частей научного текста — была достигнута.

## СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Soto C., Yoo S. Visual Detection with Context for Document Layout Analysis // Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019. С. 3464–3470.
2. Binmakhashen G.M., Mahmoud S.A. Document Layout Analysis: A Comprehensive Survey // ACM Comput. Surv. 2019. Т. 52, № 6.
3. Song M., Rosenfeld A., Kanungo T. Document structure analysis algorithms: A literature survey // Proceedings of SPIE — The International Society for Optical Engineering. 2003. Т. 5010. С. 197–207.
4. Arlazarov et al. Document image analysis and recognition: a survey // Computer Optics. 2022. Т. 46. С. 567–589.
5. Bhowmik S. Document Layout Analysis. — Springer Singapore, 2023 — 86 с.
6. Bhowmik et al. Text and non-text separation in offline document images: a survey // International Journal on Document Analysis and Recognition (IJDAR). 2018. Т. 21.
7. Kasturi R., O’Gorman L., Govindaraju V. Document image analysis: A primer // Sadhana — Academy Proceedings in Engineering Sciences. 2002. Т. 27. С. 3–22.
8. Kise K. Page Segmentation Techniques in Document Analysis // Doermann D., Tombre K. Handbook of Document Image Processing and Recognition. — Springer London, 2014. С. 135–175.
9. Бутенко Ю.И. Модель текста научно-технической статьи для разметки в корпусе научно-технических текстов // Вестник Новосибирского государственного университета. Серия: Информационные технологии. 2022. Т. 20, № 1. С. 5–13.
10. Романов Д.А. Кратко о структуре экспериментальной научной статьи на английском языке // Вестник Казанского технологического университета. 2014. Т. 17, № 6. С. 325–327.

11. Раицкая Л.К. Структура научной статьи по политологии и международным отношениям в контексте качества научной информации // Полис. Политические исследования. 2019. № 1. С. 167–181.



## **ПРИЛОЖЕНИЕ А**

Презентация к научно-исследовательской работе содержит 7 слайдов.