



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

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА
К ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЕ
НА ТЕМУ:

Метод выделения составных частей научного текста на основе
анализа распределения пикселей в сканирующей строке

Студент	ИУ7-84Б		К. А. Рунов
	(группа)	(подпись)	(инициалы, фамилия)
Руководитель ВКР		(подпись)	Ю. В. Строганов
			(инициалы, фамилия)
Консультант		(подпись)	(инициалы, фамилия)
Консультант		(подпись)	(инициалы, фамилия)
Нормоконтролер		(подпись)	А. С. Кострицкий
			(инициалы, фамилия)

2025 год

РЕФЕРАТ

СОДЕРЖАНИЕ

РЕФЕРАТ	5
ВВЕДЕНИЕ	8
1 Аналитический раздел	9
1.1 Анализ предметной области	9
1.1.1 Анализ структуры документов	9
1.1.2 Типы макетов документов	12
1.1.3 Структура научно-технического текста	13
1.2 Формализация предметной области	14
1.3 Описание существующих методов	15
1.3.1 Анализ связанных компонент	15
1.3.2 Анализ проекционного профиля	16
1.3.3 Алгоритм размазывания по длине серии	17
1.3.4 Методы на основе машинного обучения	17
1.3.5 Гибридные методы на основе РРА и ССА	18
1.4 Классификация существующих методов	18
1.5 Формализованная постановка задачи	20
2 Конструкторский раздел	22
2.1 Требования и ограничения метода	22
2.2 Описание разрабатываемого метода	22
2.2.1 Первичная разметка	23
2.2.2 Уточненная разметка	35
2.2.3 Объединенная разметка	43
2.3 Тестирование и классы эквивалентности	44
2.3.1 Тестирование первичной разметки	46
2.3.2 Тестирование уточненной разметки	48
2.3.3 Тестирование объединенной разметки	49
2.4 Структура разрабатываемого ПО	50
3 Технологический раздел	51
3.1 Выбор средств реализации	51
3.2 Реализация программного обеспечения	52

3.2.1	Точка входа	52
3.2.2	Функция обработки страниц	54
3.2.3	Интерфейсная функция для разметки	55
3.2.4	Выделение характеристик строки пикселей	55
3.2.5	Функция создания построчной разметки	56
3.2.6	Функция классификации строки	57
3.2.7	Первичная и уточненная разметки	57
3.2.8	Функция обновления состояния КА	58
3.2.9	Функция классификации сегмента	58
3.2.10	Обновление информации о сегменте	59
3.2.11	Слияние смежных сегментов одного класса	60
3.2.12	Объединенная разметка	60
3.3	Результаты тестирования	61
3.4	Пользовательский интерфейс	61
3.5	Демонстрация работы программы	64
3.6	Руководство пользователя	68
4	Исследовательский раздел	71
4.1	Описание исследования	71
4.2	Технические характеристики	72
4.3	Результаты исследования	72
4.3.1	Исследование времени выполнения	73
4.3.2	Исследование затрачиваемой памяти	74
	ЗАКЛЮЧЕНИЕ	76
	СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	79
	ПРИЛОЖЕНИЕ А	80

ВВЕДЕНИЕ

Современная научная деятельность характеризуется стремительным ростом объёма публикуемых материалов. По данным Clarivate [1], в 2023 году в базе Journal Citation Reports (JCR) насчитывалось более 21 500 научных журналов, охватывающих 254 научные дисциплины и представляющих 112 стран. Такой масштаб публикационной активности требует эффективных методов обработки и анализа научных текстов.

Научно тексты обладают определенной структурой, включающей текст, таблицы, графики, схемы, рисунки и листинги. Корректное выделение и идентификация этих элементов необходимы для таких задач, как автоматическое создание аннотаций, индексация информации, извлечение ключевых данных и повышение качества работы поисковых систем [2, 3].

Несмотря на существование различных методов анализа структуры документов, большинство из них либо основаны на ресурсоёмких подходах с использованием глубокого машинного обучения, либо не позволяют эффективно выделять все специфические компоненты научно-технических текстов. Это подтверждается сравнительным анализом методов, проведенным в аналитической части настоящей работы. Например, традиционные методы, такие как анализ проекционного профиля и связанных компонент, ограничены определёнными типами макетов и чувствительны к шумам и искажениям, в то время как методы на основе нейросетей требуют значительных вычислительных ресурсов и объёмных размеченных данных для обучения [4, 5].

Таким образом, актуальной задачей является разработка метода, который, опираясь на простые эвристические правила, позволял бы эффективно сегментировать и классифицировать компоненты научно-технических текстов без использования нейронных сетей. Это позволит существенно сократить затраты вычислительных ресурсов и повысить скорость обработки документов. На решение данной проблемы направлена настоящая дипломная работа.

1 Аналитический раздел

В данном разделе будет проведен анализ предметной области анализа структуры документов, будет приведена формализация задачи предметной области, будут описаны существующие методы и алгоритмы, будет проведена классификация существующих методов, будет обоснована потребность в разработке нового метода, будет сформулирована цель данной работы и формализована постановка задачи.

1.1 Анализ предметной области

В данном подразделе будет проведен анализ предметной области анализа структуры документов, будут рассмотрены существующие типы макетов документов, будут описаны структура научно-технического текста и его составные части.

1.1.1 Анализ структуры документов

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. [5]

Процесс анализа структуры документов состоит из двух основных этапов — этапа предварительной обработки и этапа анализа макета документа [5, 4].

На рисунке ниже приведена схема процесса анализа структуры документов.

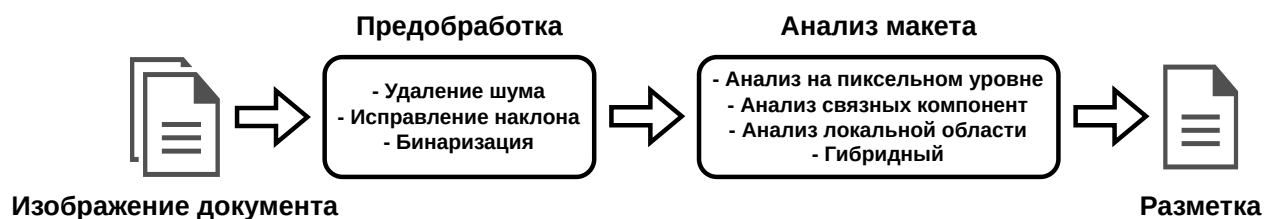


Рисунок 1 – Схема процесса анализа структуры документов [4]

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 [5, 4].

Цель этапа предварительной обработки — преобразовать входное изображение в соответствии с требованиями этапа анализа макета документа конкретного метода [5, 4].

В общем случае на этом этапе используются одна или несколько процедур предварительной обработки, таких как бинаризация, выравнивание и улучшение изображения [5, 7].

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. [4]

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. Про-

цесс 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). [5]

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. [4]

1.1.2 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.

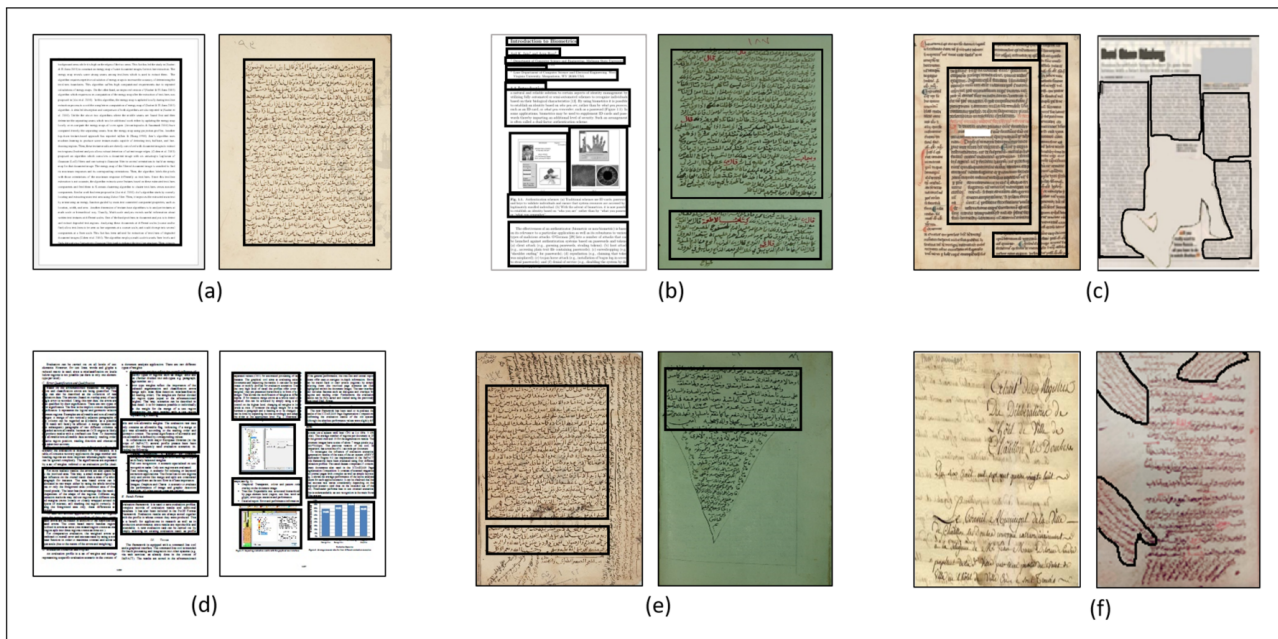


Рисунок 2 – Макеты документов: (а) Стандартный (прямоугольный), (b) Манхэттенский, (c) Не-Манхэттенский, (d) Многоколоночный Манхэттенский, (e) Произвольный (сложный), (f) С горизонтальным и диагональным наложением. [5]

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.

1.1.3 Структура научно-технического текста

Научно-технический текст обычно [9, 10, 11] следует четко определенному шаблону и имеет следующую структуру:

- 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, листинги);
- 6) Заключение;
- 7) Ссылки на литературу.

Содержимое научного текста часто не ограничивается текстом, а содержит также следующие составные части:

- 1) таблицы,
- 2) листинги,
- 3) схемы алгоритмов,
- 4) рисунки,

5) графики.

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

1.2 Формализация предметной области

Пусть D — документ, представленный в виде набора изображений, содержащих текст, листинги, таблицы, рисунки и прочие структурные элементы.

Документ

$$D = \{P_1, P_2, \dots, P_n\} \quad (1)$$

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 $(x_{i,j}, y_{i,j}, w_{i,j}, h_{i,j})$, где $(x_{i,j}, y_{i,j})$ — координаты верхнего левого угла, $w_{i,j}$ — ширина, $h_{i,j}$ — высота объекта.

Требуется построить отображение

$$F : D \rightarrow \{(O_{i,j}, C_{i,j})\}, \quad (2)$$

где каждому объекту $O_{i,j}$ ставится в соответствие класс

$$C_{i,j} = C_{i,j}(O_{i,j}), \quad (3)$$

область допустимых значений которого определяется исходя из требований к разметке.

Например, в случае задачи выделения составных частей научного текста, $C_{i,j} \subseteq \{\text{Фон, Текст, Таблица, Листинг, Схема алгоритма, Рисунок, График, Неопределенность}\}$; Объект классифицируется как «Неопределенность» в случае, когда не удалось распределить его ни в один из предыдущих классов.

Поставленную задачу можно решить, разбив на две подзадачи и решив каждую подзадачу соответственно: первая подзадача — нахождение объектов на страницах и выявление их геометрических свойств, вторая подзадача — классификация найденных объектов (определение $C_{i,j}$ для каждого объекта $O_{i,j}$).

Решением первой подзадачи является построение отображений

$$P_i \rightarrow \{O_{i,j}\}, \quad (4)$$

решением второй подзадачи является построение отображений

$$O_{i,j} \rightarrow C_{i,j}. \quad (5)$$

Далее будут рассмотрены существующие методы, позволяющие решить поставленную задачу.

1.3 Описание существующих методов

В данном подразделе будет описана суть методов на основе анализа связных компонент, методов на основе анализа проекционного профиля, алгоритма размазывания по длине серии, методов на основе машинного обучения и гибридных методов на основе анализа проекционного профиля и анализа связных компонент.

1.3.1 Анализ связных компонент

Методы на основе связных компонент (Connected Component Analysis, CCA) 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) [4].

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 связанных компонент могут работать в условиях скошенного текста при условии, что межстрочный интервал меньше пробела между абзацами [4].

1.3.2 Анализ проекционного профиля

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 [4].

1.3.3 Алгоритм размазывания по длине серии

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.

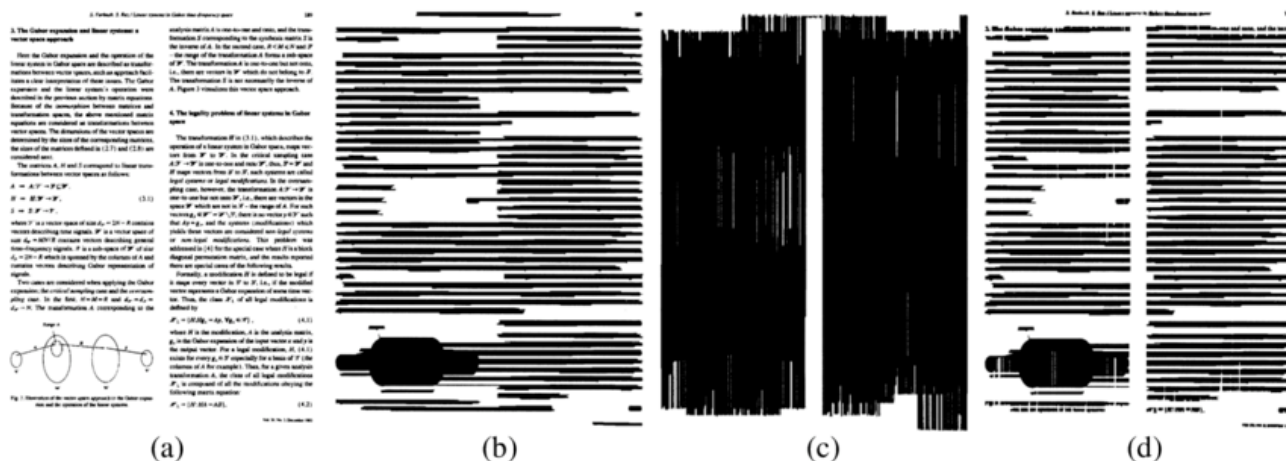


Рисунок 3 – Пример работы алгоритма RLSA. К начальному изображению документа (а) применяется горизонтальный (б) и вертикальный (с) RLSA, после чего в результате применения операции II к изображениям (б) и (с) формируется (д)

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 [4].

1.3.4 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. [5]

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. [5]

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

1.3.5 Гибридные методы на основе РРА и ССА

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

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

1.4 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.
- Специальное требование — позволяет сегментировать не только текст, но и такие составные части научного текста, как таблицы, листинги, схемы, рисунки, графики и прочее.

Ниже приведена таблица со сравнительным анализом рассмотренных методов.

Таблица 1 – Классификация методов DLA

Метод	Стратегия	Скорость	Гибкость	Ус-ть	СпецТреб
ССА	Снизу вверх	2	2	3	Нет
РРА	Сверху вниз	2	3	3	Нет
RLSA	Сверху вниз	1	3	3	Нет
ML	Снизу вверх	3	1	1	Да
РРА + ССА	Гибридный	2	3	2	Нет

Как можно видеть по сравнительной таблице, ни один из рассмотренных методов не позволяет «быстро», на основе лишь различных эвристик, произвести разметку документа, сегментирующую не только текст, но и другие составные части научного текста.

1.5 Формализованная постановка задачи

Целью данной работы является разработка метода, позволяющего выделять составные части научного текста на основе простых правил, без использования нейросетей, а также разработка алгоритма, реализующего данный метод.

Формализованная в виде IDEF0 диаграммы постановка задачи представлена на рисунке 4 ниже.

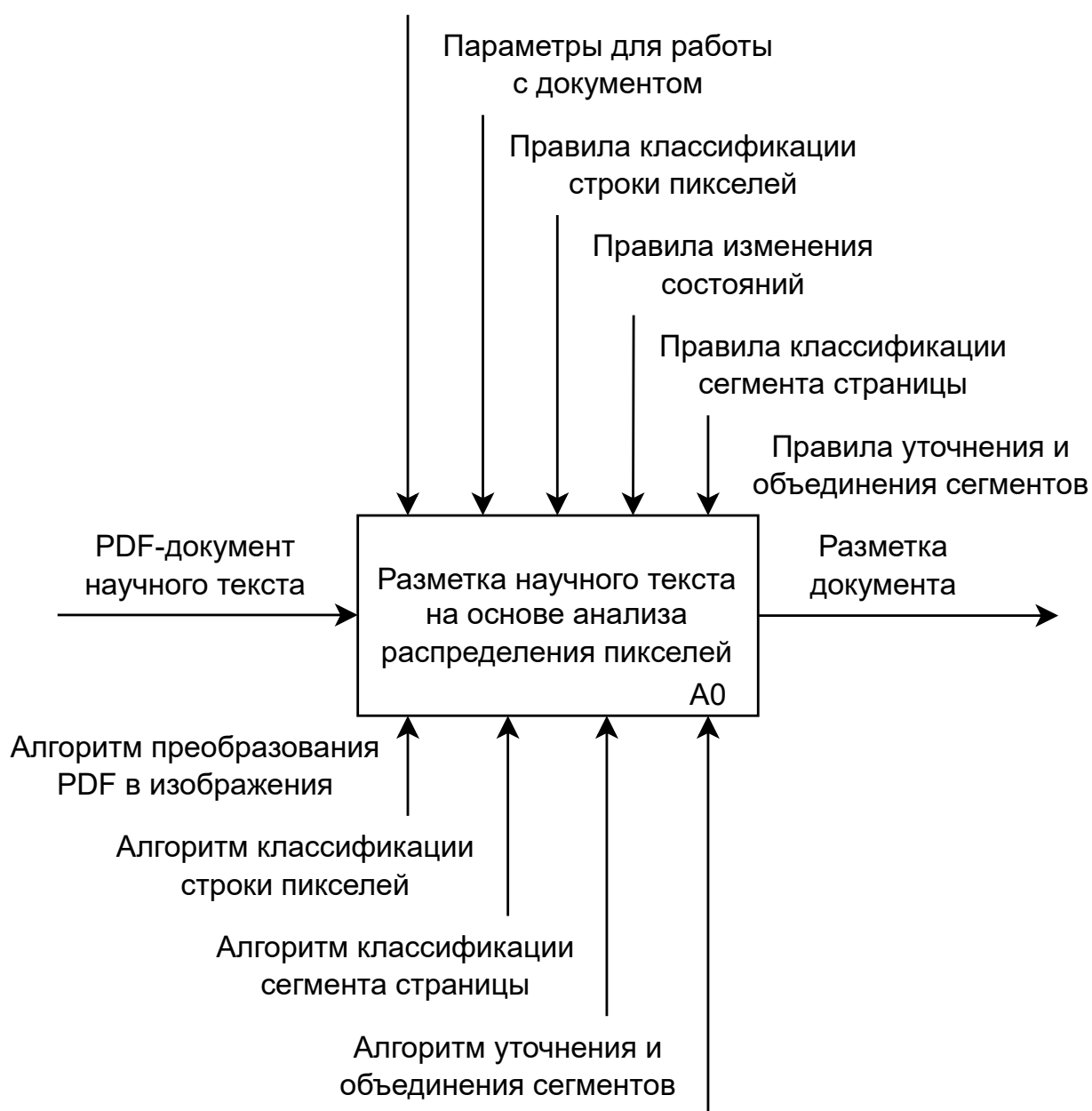


Рисунок 4 – Постановка задачи

Вывод

В данном разделе был проведен анализ предметной области анализа структуры документов, была приведена формализация задачи предметной области, были описаны существующие методы и алгоритмы, была проведена классификация существующих методов, была обоснована потребность в разработке нового метода, была сформулирована цель данной работы и формализована постановка задачи.

2 Конструкторский раздел

В данном разделе будут приведены требования и ограничения разрабатываемого метода, будут описаны основные этапы разрабатываемого метода, сценарии тестирования, классы эквивалентности тестов, а также структура разрабатываемого ПО.

2.1 Требования и ограничения метода

Метод выделения составных частей научного текста на основе анализа распределения пикселей в сканирующей строке должен:

- 1) Работать с одноколонными Манхэттенскими макетами документов;
- 2) Выделять текстовые блоки;
- 3) Выделять таблицы;
- 4) Выделять листинги;
- 5) Выделять схемы алгоритмов;
- 6) Выделять рисунки;
- 7) Выделять графики;
- 8) Работать на основе простых правил и эвристик, без использования нейросетей.

2.2 Описание разрабатываемого метода

Поставленная задача решается в четыре этапа:

- 1) Преобразование PDF документа в изображения;
- 2) Первичная разметка страниц;
- 3) Создание уточненной разметки на основе первичной;
- 4) Объединение уточненной разметку в более крупные блоки.

Разметка, ее уточнение и объединение происходят на основе определенных правил, которые будут описаны в данном разделе далее.

Основные этапы разрабатываемого метода представлены на IDEF0 диаграмме первого уровня (см. Рисунок 5).

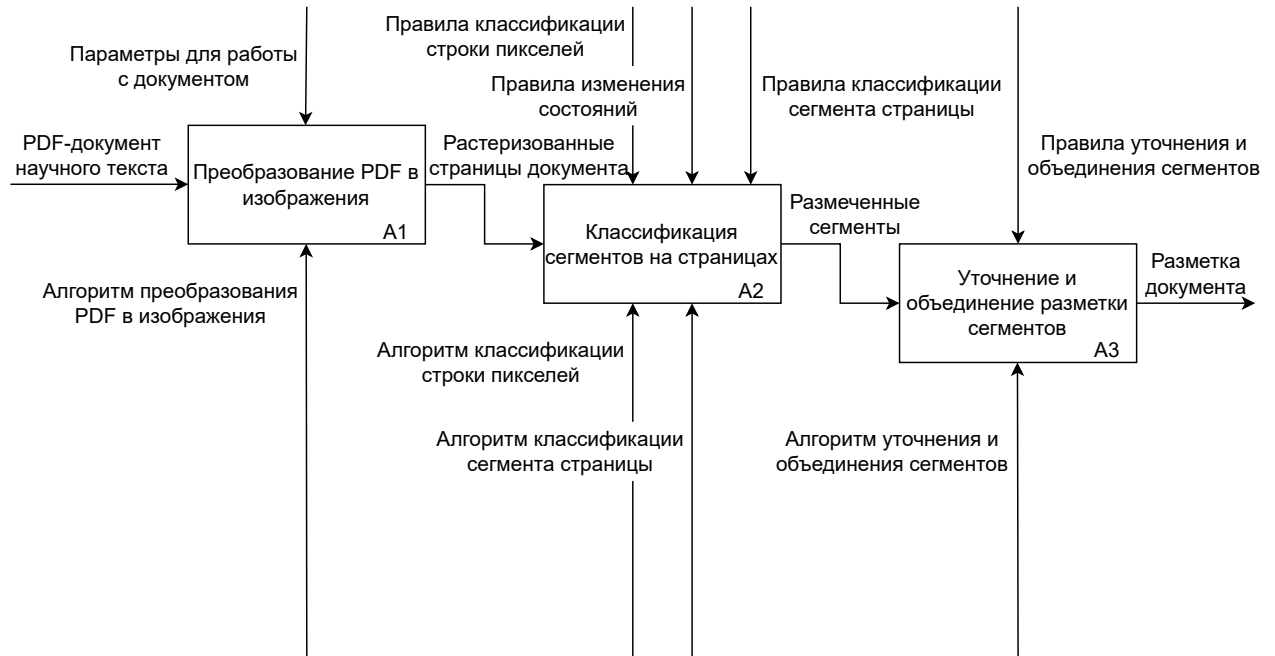


Рисунок 5 – IDEF0-диаграмма метода выделения составных частей научного текста на основе анализа распределения пикселей в сканирующей строке

2.2.1 Первичная разметка

Исходные данные — изображение страницы документа. Получаемый результат — разметка страницы и информация о каждом сегменте на ней.

Разметка страницы — массив кортежей типа

$$(y_start, y_end, C), \quad (6)$$

где y_start — y -координата начала сегмента в пространстве изображения, y_end — y -координата конца сегмента в пространстве изображения, C — класс сегмента, где $C \subseteq \{\text{Фон, Немного текста, Много текста, Цвет, Черная линия средней длины, Длинная черная линия, Не определено}\}$.

Такое множество классов было выделено по двум причинам. Во-первых, на основе распределения пикселей в строке ее уже можно причислить к одному из перечисленных классов. Во-вторых, на основе принадлежности строки

к одному из классов можно делать предположения насчет класса сегмента, которому строка принадлежит.

Так, если строка классифицируется как «Много текста» можно предположить, что сегмент вероятно принадлежит классу «Текст».

Если строка классифицируется как «Цвет», можно предположить, что сегмент вероятно принадлежит классу «Рисунок» или «График».

Если строка классифицируется как «Черная линия средней длины», можно предположить, что сегмент вероятно принадлежит классу «Схема алгоритма».

Если строка классифицируется как «Длинная черная линия», можно предположить, что сегмент вероятно принадлежит классу «Таблица» или «Листинг».

Информация о сегменте содержит следующие данные:

- 1) `start` — ордината начала сегмента;
- 2) `end` — ордината конца сегмента;
- 3) `count_long_black_line` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Длинная черная линия»;
- 4) `count_single_long_black_line` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Длинная черная линия», считая несколько подряд идущих «Длинных черных линий» за одну;
- 5) `count_medium_black_line` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Черная линия средней длины»;
- 6) `count_single_medium_black_line` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Черная линия средней длины», считая несколько подряд идущих «Черных линий средней длины» за одну;
- 7) `count_total_medium_black_line` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Черная линия

средней длины», с учетом всех «Черных линий черной длины» если таких было зафиксировано несколько внутри одной сканирующей строки;

- 8) `count_many_text` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Много текста»;
- 9) `count_few_text` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Немного текста»;
- 10) `count_color` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Цвет»;
- 11) `count_undefined` — количество раз, когда при разметке сегмента встретилась строка, идентифицированная, как «Не определено»;
- 12) `count_white_px` — количество белых пикселей в сегменте;
- 13) `count_color_px` — количество цветных пикселей в сегменте;
- 14) `count_gray_px` — количество черных пикселей в сегменте (сумма трех данных счетчиков дает общее количество пикселей в сегменте);
- 15) `heatmap_black` — массив, i -й элемент которого отражает количество черных пикселей в i -й колонке пикселей сегмента;
- 16) `heatmap_color` — массив, i -й элемент которого отражает количество цветных пикселей в i -й колонке пикселей сегмента.

Данная информация будет использоваться для уточнения разметки в следующем этапе.

Первичная разметка создается в результате классификации строк на основе распределения пикселей в них и изменения состояний конечного автомата первичной разметки.

Диаграмма изменения состояний конечного автомата изображена на рисунке 6. Классификация любого сегмента начинается из состояния «Фон», и заканчивается в состоянии «Фон», причем в результате классификации сегменту присваивается класс в соответствии с предпоследним состоянием, в котором находился конечный автомат (до последнего «Фона»).

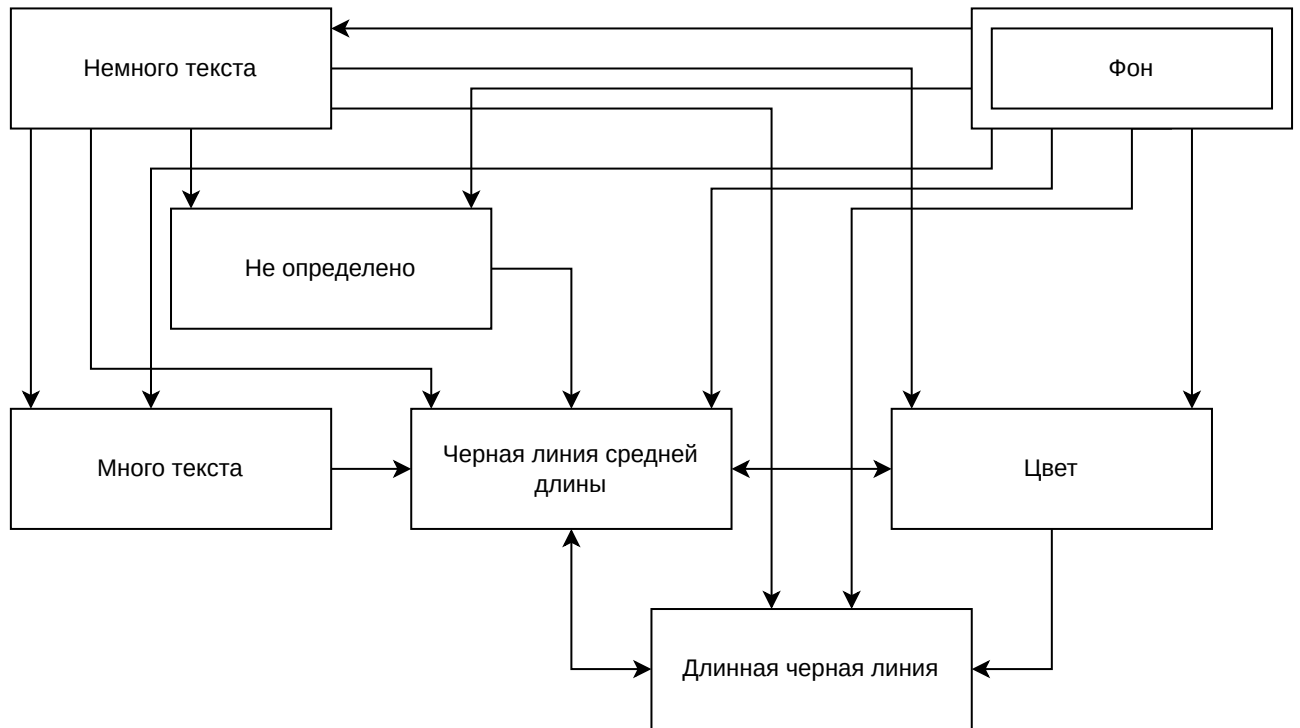


Рисунок 6 – Конечный автомат, все состояния

Из состояния «Фон» можно попасть в любое состояние. Из состояния «Немного текста» можно попасть в любое состояние, кроме «Фона».

На рисунке 7 изображена редуцированная диаграмма конечного автомата, опускающая состояния «Фон» и «Немного текста».

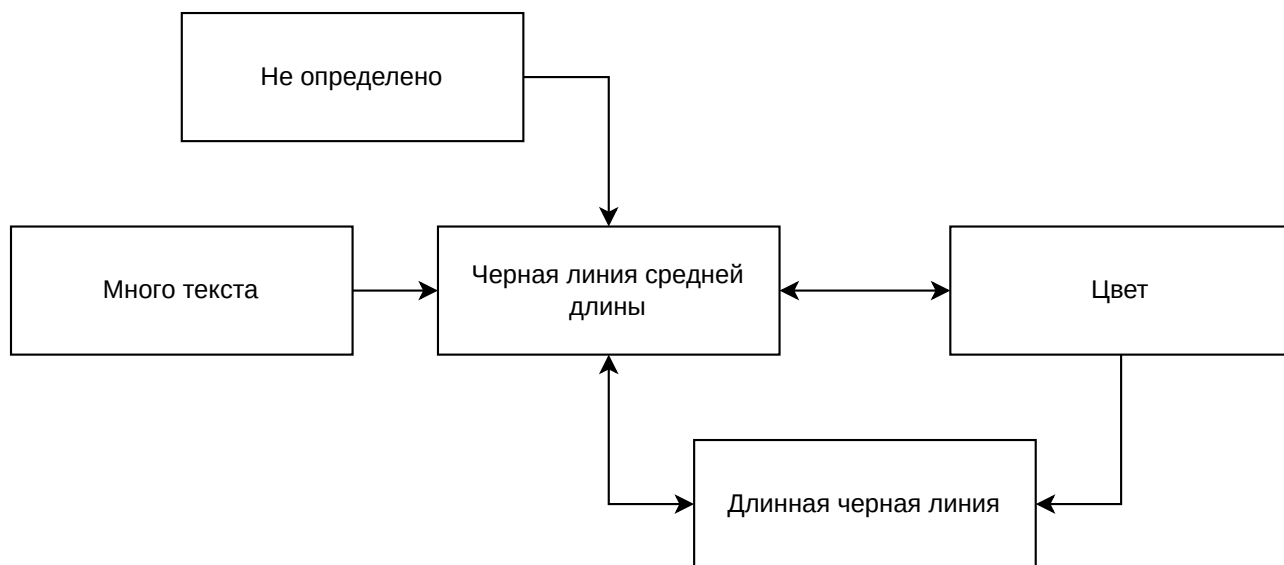


Рисунок 7 – Конечный автомат, состояния кроме «Немного текста» и «Фон»

На рисунке 8 ниже изображена схема алгоритма классификации конкретной сканирующей строки.

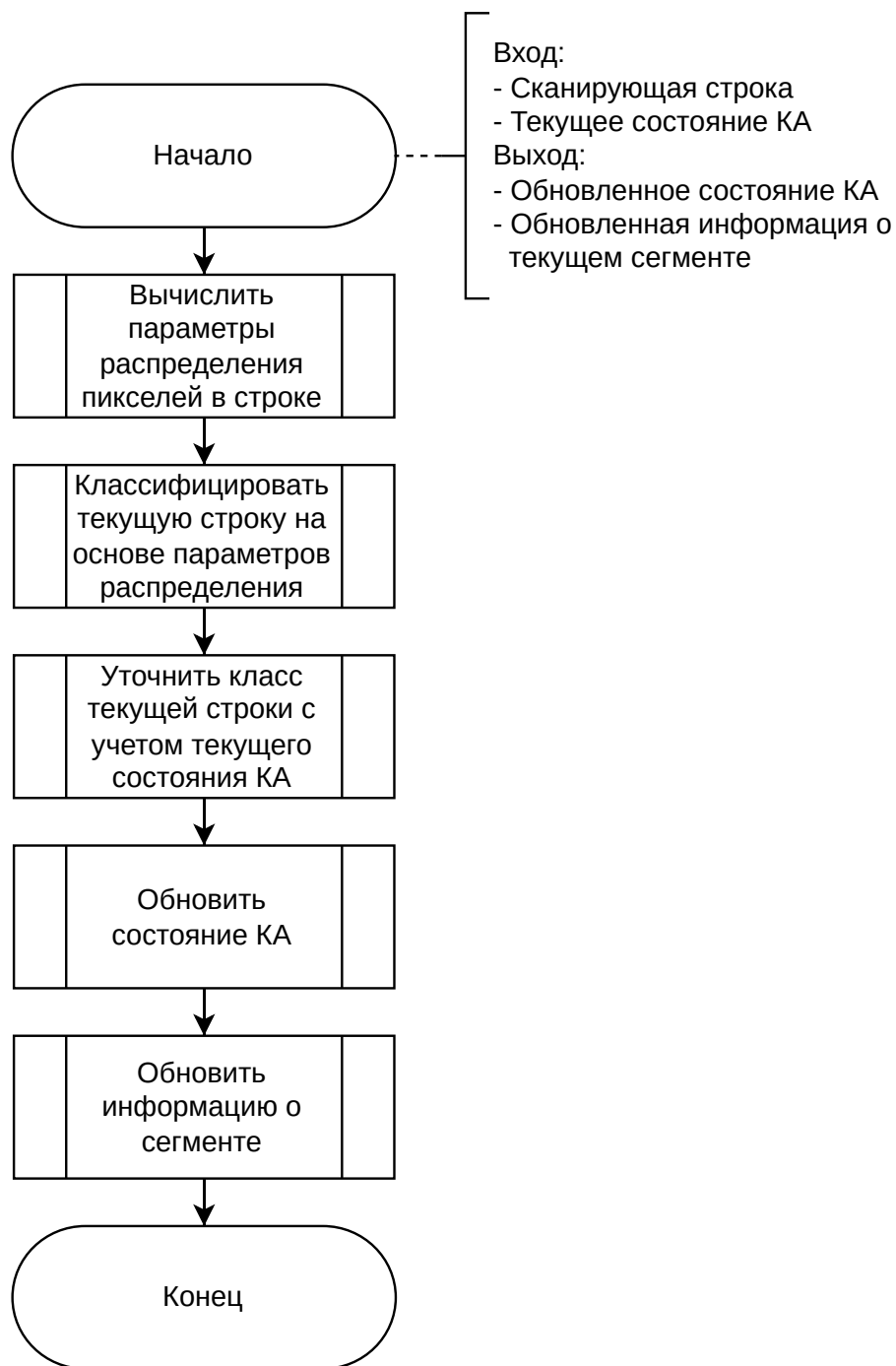


Рисунок 8 – Разметка сканирующей строки

Классификация строки

Классификация строки производится на основе следующих параметров распределения пикселей в ней:

- 1) `count_white` — количество белых пикселей в строке;
- 2) `count_color` — количество цветных пикселей в строке;
- 3) `count_gray` — количество серых (почти черных) пикселей в строке;
- 4) `comp_lengths` — массив длин участков подряд идущих не белых пикселей;
- 5) `gap_lengths` — массив длин промежутков (белых пикселей) между участками подряд идущих не белых пикселей;
- 6) `gray_comp_lengths` — массив длин участков подряд идущих серых пикселей;
- 7) `color_comp_lengths` — массив длин участков подряд идущих цветных пикселей;
- 8) `first_nonwhite_index` — индекс первого не белого пикселя в строке.

Классификация как «Фон»

Строка классифицируется как «Фон», если выполнено условие «`first_nonwhite_index` не установлен» — в строке не нашлось не белых пикселей.

Классификация как «Длинная черная линия»

Строка классифицируется как «Длинная черная линия», если: строка содержит единственную серую компоненту И эта компонента достаточно длинная.

Содержание единственной серой компоненты определяется на основе одновременного выполнения следующих условий:

- Длина `comp_lengths` равна 1;

- Длина `gap_lengths` равна 0;
- Длина `color_comp_lengths` равна 0.

Длина серой компоненты для классификации строки как «Длинная черная линия» считается достаточно большой, если `count_gray` больше некоторого параметра, который является произведением длины строки на некоторую наперед заданную константу, принимающую значения от 0 до 1, например $1/2$.

Классификация как «Черная линия средней длины»

Строка классифицируется как «Черная линия средней длины», если строка содержит компоненту, длина которой больше некоторого параметра, который является произведением длины строки на некоторую наперед заданную константу, принимающую значения от 0 до 1, например $1/16$.

Классификация как «Много текста»

Строка классифицируется как «Много текста», если она либо содержит очень много черных компонент, либо содержит много черных компонент и не содержит цвета.

Считается, что строка содержит очень много черных компонент, если длина `comp_lengths` превышает некоторую наперед заданную константу, например 100.

Считается, что строка содержит много черных компонент, если длина `comp_lengths` превышает некоторую наперед заданную константу, например 80.

Таким образом, если строка содержит очень много черных компонент, она будет классифицирована как «Много текста» вне зависимости от наличия в ней цвета.

Классификация как «Цвет»

Строка классифицируется как «Цвет», если `count_color` больше нуля.

Классификация как «Немного текста»

Строка классифицируется как «Немного текста», если содержит немного компонент И компоненты преимущественно небольшого размера И промежутки между компонентами преимущественно небольшого размера И отсутствуют большие промежутки.

Считается, что компоненты и/или промежутки между компонентами преимущественно небольшого размера, если среднее арифметическое `comp_lengths` и/или `gap_lengths` меньше некоторой наперед заданной константы, например 20 пикселей.

Считается, что большие промежутки отсутствуют, если в массиве `gap_lengths` отсутствуют элементы с индексом стандартного отклонения больше шести.

Индекс стандартного отклонения Z для элемента массива x вычисляется по формуле:

$$Z(x) = \frac{x - \mu}{\sigma}, \quad (7)$$

где μ — среднее значение элементов массива, σ — стандартное отклонение элементов массива.

Классификация как «Не определено»

В остальных случаях строка классифицируется как «Не определено».

Изменение состояний КА

Далее будут представлены таблицы переходов состояний конечного автомата, в которых некоторые правила были подобраны исходя из логических соображений, а некоторые — эмпирическим путем.

Таблица 2 – Таблица переходов состояний из начального состояния «Фон»

Состояние текущей строки	Конечное состояние
Не определено	Не определено
Много текста	Много текста
Немного текста	Немного текста
Длинная черная линия	Длинная черная линия
Черная линия средней длины	Черная линия средней длины
Цвет	Цвет
Фон	Фон

Таблица 3 – Таблица переходов состояний из начального состояния «Не определено»

Состояние текущей строки	Конечное состояние
Не определено	Не определено
Много текста	Много текста
Немного текста	Не определено
Длинная черная линия	Не определено
Черная линия средней длины	Черная линия средней длины
Цвет	Цвет
Фон	Фон

Таблица 4 – Таблица переходов состояний из начального состояния «Много текста»

Состояние текущей строки	Конечное состояние
Не определено	Много текста
Много текста	Много текста
Немного текста	Много текста
Длинная черная линия	Много текста
Черная линия средней длины	Черная линия средней длины
Цвет	Много текста
Фон	Фон

Таблица 5 – Таблица переходов состояний из начального состояния
«Немного текста»

Состояние текущей строки	Конечное состояние
Не определено	Не определено
Много текста	Много текста
Немного текста	Немного текста
Длинная черная линия	Длинная черная линия
Черная линия средней длины	Черная линия средней длины
Цвет	Цвет
Фон	Фон

Таблица 6 – Таблица переходов состояний из начального состояния
«Длинная черная линия»

Состояние текущей строки	Конечное состояние
Не определено	Длинная черная линия
Много текста	Длинная черная линия
Немного текста	Длинная черная линия
Длинная черная линия	Длинная черная линия
Черная линия средней длины	Черная линия средней длины
Цвет	Длинная черная линия
Фон	Фон

Таблица 7 – Таблица переходов состояний из начального состояния «Черная
линия средней длины»

Состояние текущей строки	Конечное состояние
Не определено	Черная линия средней длины
Много текста	Черная линия средней длины
Немного текста	Черная линия средней длины
Длинная черная линия	Длинная черная линия
Черная линия средней длины	Черная линия средней длины
Цвет	Цвет
Фон	Фон

Таблица 8 – Таблица переходов состояний из начального состояния «Цвет»

Состояние текущей строки	Конечное состояние
Не определено	Цвет
Много текста	Цвет
Немного текста	Цвет
Длинная черная линия	Длинная черная линия
Черная линия средней длины	Черная линия средней длины
Цвет	Цвет
Фон	Фон

Пример первичной разметки

Пусть при сканировании документа встретилось данное предложение.
Разметка будет происходить следующим образом:

- 1) Начальное состояние конечного автомата — «Фон».
- 2) Встретилась строка с не белым пикселем — вершина буквы «П». В ней содержится одна сплошная черная компонента, но она недостаточно длинная для классификации, как «Черная линия средней длины». Также не хватает информации для причисления ее к какому-либо другому классу, поэтому строка классифицируется как «Не определено».
- 3) Конечный автомат меняет состояние: из «Фона» по «Не определено» переходит в «Не определено».
- 4) Обновляется информация о сегменте.
- 5) Сканируется следующая строка, в которой встречается два очень маленьких компонента с небольшим расстоянием между ними — вертикальные линии буквы «П». Этой информации оказывается достаточно для классификации строки как «Немного текста».
- 6) Конечный автомат не меняет состояние: из «Не определено» по «Немного текста» остается в «Не определено».
- 7) Обновляется информация о сегменте.

- 8) Подобные действия продолжаются, пока в сканирующей строке не встретятся вершины других букв.
- 9) Сканируется строка, в которой встречаются вершины других букв — много маленьких компонентов на протяжении почти всей длины строки. Этой информации достаточно для классификации строки как «Много текста».
- 10) Конечный автомат меняет состояние: из «Не определено» по «Много текста» переходит в «Много текста».
- 11) Обновляется информация о сегменте.
- 12) На протяжении следующих итераций строки продолжают классифицироваться как «Много текста» до тех пор, пока не начнут встречаться «хвосты» букв «р» и «у».
- 13) При сканировании «хвостов» строки будут классифицироваться как «Немного текста», но это не повлияет на состояние конечного автомата «Много текста».
- 14) Обновляется информация о сегменте.
- 15) После «хвостов» наконец встретится «Фон», конечный автомат перейдет из состояния «Много текста» в состояние «Фон», и сегмент будет классифицирован в соответствии с предпоследним состоянием конечного автомата — «Много текста».

2.2.2 Уточненная разметка

Исходные данные — разметка страницы и информация о каждом сегменте на ней. Получаемый результат — уточненная разметка страницы.

Уточненная разметка страницы — массив кортежей типа

$$(y_start, y_end, C), \quad (8)$$

где y_start — y -координата начала сегмента в пространстве изображения, y_end — y -координата конца сегмента в пространстве изображения, C —

класс сегмента, где $C \subseteq \{\text{Фон, Текст, Таблица, Листинг, Схема алгоритма, Рисунок, График, Не определено}\}$, причем координаты начала и конца сегментов совпадают с соответствующими координатами сегментов первичной разметки, уточняется только класс на основе информации о сегментах.

Далее будут перечислены правила, на основе которых сегмент меняет класс из первичного, соответствующего одному из состояний конечного автомата, в уточненный из множества $\{\text{Фон, Текст, Таблица, Листинг, Схема алгоритма, Рисунок, График, Не определено}\}$.

Из «Не определено» в...

«Текст»

Правило: Высота сегмента небольшая ИЛИ Много строк в сегменте было классифицировано как «Немного текста».

Высота сегмента считается небольшой, если

$$height < C1, \quad (9)$$

где $height = end - start$, а $C1$ — наперед заданная константа, соответствующая максимальному количеству пикселей, при котором сегмент считается небольшим.

Считается, что много строк в сегменте было классифицировано как «Немного текста», если

$$\frac{count_few_text}{height} > C2, \quad (10)$$

где $C2$ — наперед заданная константа, соответствующая минимальному отношению количества строк в сегменте, классифицированных как «Немного текста», к высоте сегмента, при котором считается, что сегмент содержит много строк, классифицированных как «Немного текста».

Обоснование: Если сегмент находится в состоянии «Не определено», значит на протяжении его разметки не встретилось ни одной строки, позволяющей классифицировать его, как «Много текста», «Черная линия средней длины» или «Цвет». Следовательно, данный сегмент вероятно относится либо к классу «Текст», либо к классу «Не определено». Если высота сегмента

небольшая, вероятно он все же относится к тексту. Такое возможно, например, в случае использования буквы «й» — сегмент, содержащий ее дугу будет классифицирован, как «Не определено».

«Таблица»

Сегмент может получить класс «Таблица» только из состояний «Длинная черная линия» или «Много текста».

«Листинг»

Правило: Сегмент содержит ровно две вертикальные черные линии высотой с весь сегмент (вычисляется на основе `heatmap_black`).

«Схема алгоритма»

Сегмент может получить класс «Схема алгоритма» только из состояний «Черная линия средней длины» или «Длинная черная линия».

«Рисунок»

Правило: Сегмент большой.

Сегмент считается большим, если

$$height > C3, \quad (11)$$

где $C3$ — наперед заданная константа, соответствующая минимальной высоте «Не определенного сегмента», чтобы классифицировать его, как «Рисунок».

Обоснование: Если на протяжении длительного отрезка документа ни разу не встретилось линии, классифицируемой, как «Много текста», «Черная линия средней длины» или «Цвет», то, вероятно, перед нами ни «Текст», ни «Схема алгоритма», ни «График», а «Рисунок».

«График»

Правило: Сегмент содержит одну высокую вертикальную линию.

Считается, что сегмент содержит высокую вертикальную линию, если в массиве `heatmap_black` содержится единственный элемент, значение которого превышает

$$height * C4, \quad (12)$$

где $C4$ — наперед заданная константа такая, что $height * C4$ не выше оси ординат графика.

Такая корректирующая константа нужна, потому что часто высота сегмента графика превышает высоту его оси абсцисс, так как на оси ординат

есть «засечки», которые увеличивают размер сегмента, но не увеличивают высоту вертикальной оси.

«Не определено»

Сегмент получает класс «Не определено», если ранее не удалось причислить его ни к одному другому классу.

Порядок проверки принадлежности к уточненному классу:

- 1) Текст,
- 2) Листинг,
- 3) Рисунок,
- 4) График,
- 5) Не определено.

Из «Немного текста» в...

«Текст»

Правило: Всегда.

Обоснование: В соответствии с конечным автоматом, сегмент получает класс «Немного текста» только в том случае, если на протяжении его разметки встречались только строки, классифицируемые, как «Немного текста». Следовательно, вероятно, данный сегмент относится к «Тексту».

Из «Много текста» в...

«Таблица»

Правило: Сегмент достаточно большой И содержит больше двух вертикальных линий И минимальное расстояние между двумя вертикальными линиями не слишком маленькое.

«Листинг»

Правило: Сегмент достаточно большой И содержит ровно две вертикальные линии И не содержит черных линий среднего размера И (содержит хотя бы одну строку, классифицированную как «Много текста» ИЛИ не содержит строк, классифицированных как «Цвет»)

Обоснование: В листингах кода как правило не встречается черных линий среднего размера, но при этом бывают схемы алгоритмов, например IDEF0 диаграммы, которые, как и листинг, ограничены двумя вертикальными черными линиями.

Если сегмент содержал хотя бы одну строку, классифицированную как «Много текста» (что часто бывает в случае с листингами), то, вероятно, данный сегмент относится к классу «Листинг» несмотря на наличие цвета, которым, бывает, выделяют ключевые слова, относящиеся к конкретному языку программирования.

Если же таких строк нет, при этом есть строки, классифицированные, как «Цвет», с большой вероятностью утверждать, к какому классу относится сегмент, становится затруднительно.

«Текст»

Сегмент получает класс «Текст», если ранее не удалось причислить его ни к одному другому классу.

Порядок проверки принадлежности к уточненному классу:

- 1) Таблица,
- 2) Листинг,
- 3) Текст.

Из «Цвет» в...

«График»

Правило: Сегмент имеет единственную высокую вертикальную линию И отношение цветных пикселей в сегменте к белым мало.

Считается, что отношение цветных пикселей в сегменте к белым мало, если

$$\frac{count_color_px}{count_white_px} < C5, \quad (13)$$

где $C5$ — наперед заданная константа, соответствующая минимальному отношению цветных пикселей в сегменте к белым, чтобы данное отношение считалось малым.

Обоснование: Как правило, графики в документах имеют ось ординат, а сами изображены цветом на белом фоне, следовательно, белые пиксели в

сегменте с графиком преобладают над цветными.

«Не определено»

Правило: Сегмент небольшой.

Считается, что сегмент небольшой, если его высота меньше некоторой наперед заданной константы C_6 , соответствующей максимальной высоте сегмента, при которой сегмент можно считать небольшим.

«Рисунок»

Сегмент получает класс «Рисунок», если ранее не удалось причислить его ни к одному другому классу.

Порядок проверки принадлежности к уточненному классу:

- 1) График,
- 2) Не определено,
- 3) Рисунок.

Из «Черная линия средней длины» в...

«Текст»

Правило: В сегменте встречалось много строк, идентифицируемых, как «Много текста» ИЛИ (сегмент небольшой И (в сегменте встречалось много строк, идентифицируемых, как «Немного текста» ИЛИ «Не определено»)).

Обоснование: Такое правило позволяет идентифицировать как «Текст» сегменты, в которых подчеркивание текста сливается с самим текстом, например, на титульном листе в разделе с фамилиями студента и преподавателей.

«Рисунок»

Правило: Сегмент достаточно большой И (имеет цвет ИЛИ имеет много строк, идентифицируемых как «Черная линия средней длины»).

Обоснование: Рисунки в документах, как правило, занимают значительную часть страницы, поэтому небольшой сегмент, идентифицированный изначально, как «Черная линия средней длины», вероятно не является «Рисунком» или «Графиком».

Рисунки часто имеют цвет, а если они не цветные, и на них встретились «Черная линия средней длины», то вероятно строк, идентифицированных таким образом, на них будет достаточно много. Если таких строк мало, то скорее сегмент относится к классу «Схема алгоритма» нежели «Рисунок».

«График»

Правило: Сегмент имеет цвет И имеет не много строк, идентифицированных как «Черная линия средней длины» И количество вертикальных черных линий не меньше двух И количество белых пикселей преобладает над остальными.

Обоснование: Данное правило учитывает ситуации, когда прикрепленные в документе графики имеют небольшой масштаб (иначе они были бы изначально классифицированы как «Длинная черная линия»). А когда прикрепленные в документе графики имеют небольшой масштаб, они часто ограничены черным прямоугольником, откуда и появляется условие «количество вертикальных черных линий не меньше двух».

«Не определено»

Правило 1: Сегмент имеет единственную «Черную линию средней длины» ИЛИ (сегмент не очень высокий И имеет не много «Черных линий средней длины»).

Правило 2: Сегмент небольшой (меньше наперед заданной константы, например 20 пикселей).

Обоснование: Правило 1 позволяет классифицировать как «Не определено» формулы и уравнения, содержащие квадратные корни, дроби и знаки « \sum », ведь «Черная линия средней длины» часто встречается в единственном экземпляре именно в формулах и уравнениях.

Можно было бы классифицировать такие сегменты, как «Формула», а не «Не определено», однако разрабатываемый метод не позволяет отличать более простые формулы от текста, поэтому было принято решение не вводить для данной ситуации отдельный класс.

Если сегмент небольшой, то вероятно он не принадлежит к какому-либо другому классу, кроме как «Не определено».

«Схема алгоритма»

Также сегмент получает класс «Схема алгоритма», если ранее не удалось причислить его ни к одному другому классу.

Порядок проверки принадлежности к уточненному классу:

- 1) График,
- 2) Рисунок,

- 3) Схема алгоритма,
- 4) Текст,
- 5) Не определено,
- 6) Схема алгоритма.

Из «Длинная черная линия» в...

«Таблица»

Правило: Сегмент достаточно большой И содержит больше двух вертикальных линий И минимальное расстояние между двумя вертикальными линиями не слишком маленькое.

«Листинг»

Правило: Сегмент содержит ровно две вертикальные линии И не содержит черных линий среднего размера И (содержит хотя бы одну строку, классифицированную как «Много текста» ИЛИ не содержит строк, классифицированных как «Цвет»).

«Схема алгоритма»

Правило: Сегмент не имеет цвета и содержит не меньше двух черных линий среднего размера.

Обоснование: В схеме алгоритма часто встречается не менее двух черных линий среднего размера — начало блока и конец блока (например на IDEF0 диаграмме).

«График»

Правило: Сегмент имеет цвет И содержит немного длинных черных линий И количество вертикальных черных линий не меньше двух И количество белых пикселей преобладает над остальными.

Считается, что сегмент содержит несколько черных линий, если

$$\frac{count_long_black_line}{height} < C7, \quad (14)$$

где $C7$ — наперед заданная константа, соответствующая максимальному значению отношения количества строк, идентифицированных как «Длинная черная линия» к высоте сегмента, при котором можно считать, что сегмент содержит немного длинных черных линий.

«Не определено»

Правило: Сегмент небольшой.

«Рисунок»

Сегмент получает класс «Рисунок», если ранее не удалось причислить его ни к одному другому классу.

Порядок проверки принадлежности к уточненному классу:

- 1) Не определено,
- 2) График,
- 3) Таблица,
- 4) Листинг,
- 5) Схема алгоритма,
- 6) Рисунок.

2.2.3 Объединенная разметка

Исходные данные — уточненная разметка страницы. Получаемый результат — объединенная разметка страницы.

Объединенная разметка страницы формируется из уточненной разметки в несколько этапов:

- 1) Слияние небольших фоновых сегментов с ближайшим наибольшим сегментом;
- 2) Объединение сегментов после слияния;
- 3) Слияние фоновых сегментов с соседними, если соседние сегменты имеют один и тот же класс;
- 4) Объединение сегментов после слияния;
- 5) Смена класса небольших фоновых сегментов на «Не определено»;
- 6) Объединение сегментов после слияния;

- 7) Слияние небольших «Не определенных» сегментов с наибольшими соседними сегментами;
- 8) Объединение сегментов после слияния.

На рисунке 9 представлен алгоритм объединения сегментов после слияния.

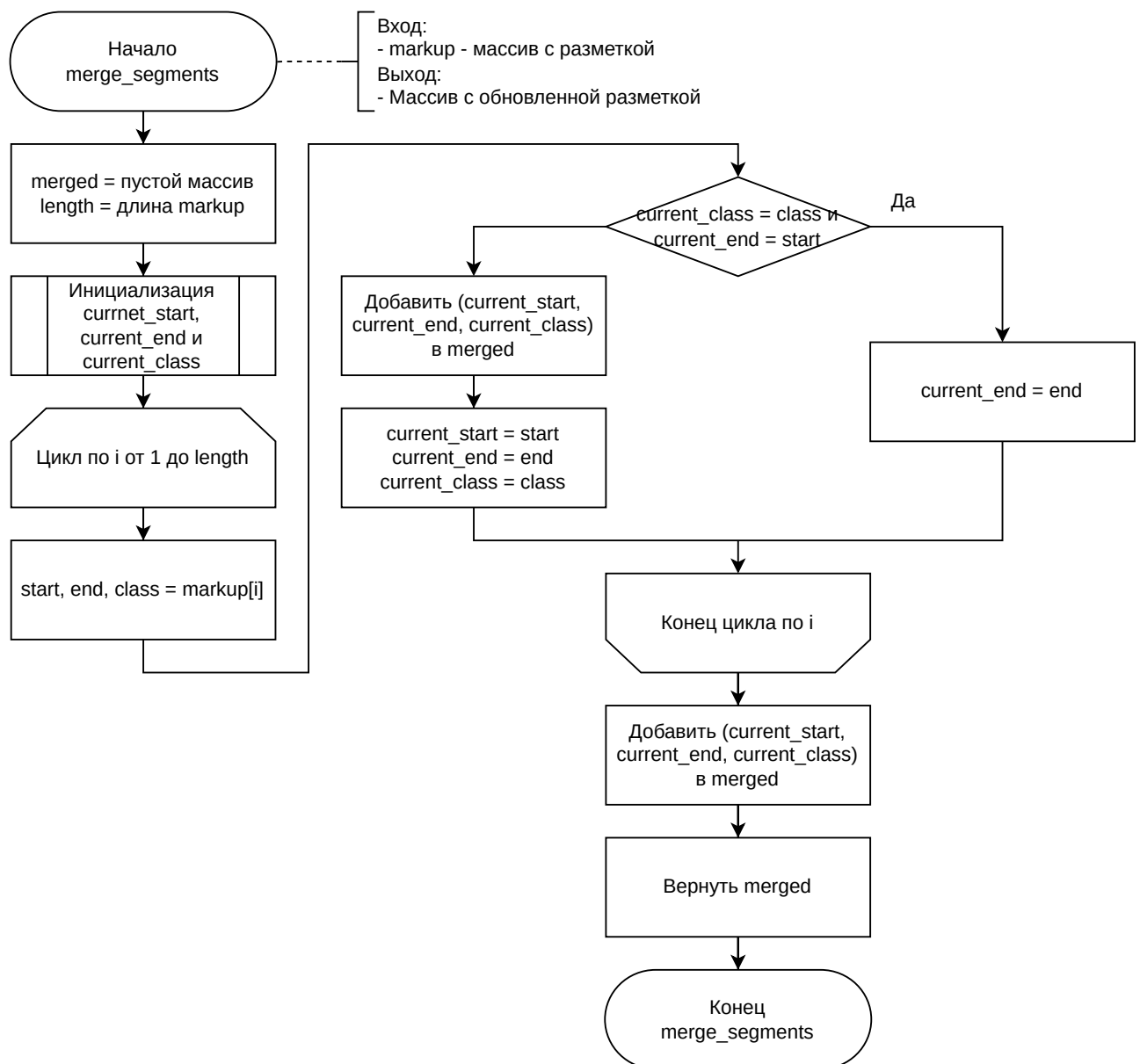


Рисунок 9 – Объединение сегментов после слияния

2.3 Тестирование и классы эквивалентности

Тестировать разрабатываемый метод имеет смысл на каждом этапе:

- 1) На этапе создания первичной разметки;

- 2) На этапе создания уточненной разметки;
- 3) На этапе создания объединенной разметки;

для лучшей локализации возможных ошибок.

Тестировать создаваемые разметки можно сравнивая с «эталонными» разметками — проходиться по документу построчно и сравнивать класс текущей строки с классом соответствующей строки в эталонной разметке.

Пример такого алгоритма тестирования приведен на рисунке ниже.

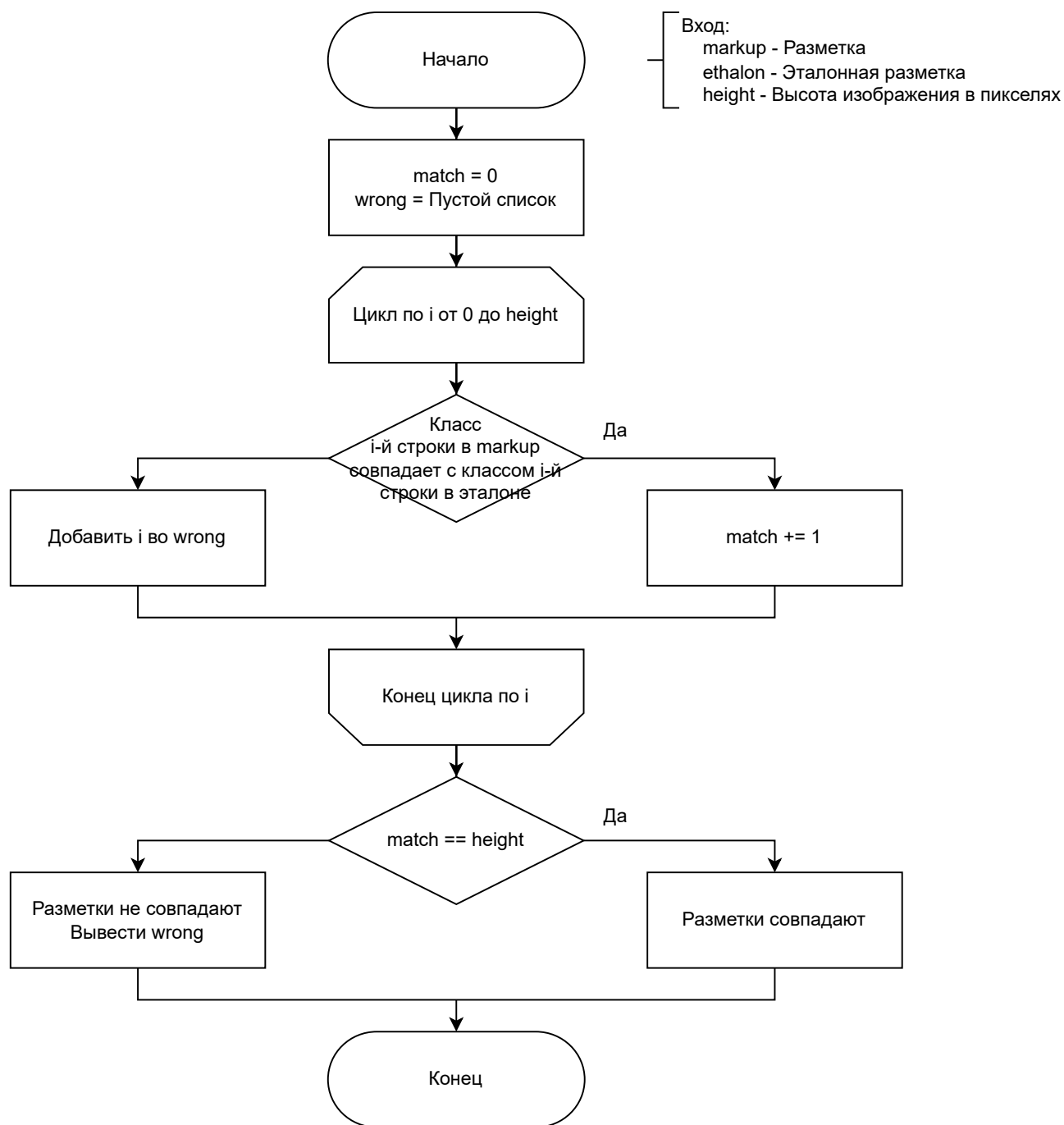


Рисунок 10 – Алгоритм сравнения разметки с эталонной

2.3.1 Тестирование первичной разметки

Классы эквивалентности при тестировании первичной разметки естественным образом соответствуют ее классам:

- Фон;
- Много текста;
- Немного текста;
- Длинная черная линия;
- Черная линия средней длины;
- Цвет;
- Не определено.

Для тестирования разметки Фона можно предложить следующие тестовые сценарии:

- Белая строка высотой в один пиксель;
- Больше одной фоновой строки подряд;
- Пустая страница.

Для тестирования вывода состояния «Много текста» можно предложить следующие тестовые сценарии:

- Строка текста на русском языке;
- Строка текста на английском языке;
- Строка текста на китайском языке;
- Строка текста длиной во всю ширину документа;
- Строка текста длиной больше половины документа.

Для тестирования вывода состояния «Немного текста» можно предложить следующие тестовые сценарии:

- Текст «Немного текста»;
- Текст «aa»;
- Строка текста длиной меньше половины документа.

Для тестирования вывода состояния «Длинная черная строка» можно предложить следующие тестовые сценарии:

- Длинная черная строка высотой в один пиксель;
- Длинная черная строка высотой в несколько пикселей;
- Листинг кода;
- Таблица;
- Черно-белое изображение в рамке.

Для тестирования вывода состояния «Черная строка средней длины» можно предложить следующие тестовые сценарии:

- Небольшая черная линия;
- Подчеркнутый текст;
- Блок схемы алгоритма;
- Схема алгоритма.

Для тестирования вывода состояния «Цвет» можно предложить следующие тестовые сценарии:

- Цветное изображение без рамки;
- Цветной текст;
- Цветной текст внутри не цветного.

Для тестирования вывода состояния «Не определено» можно предложить следующие тестовые сценарии:

- Граф сверточной нейронной сети без прямых линий;
- Черно-белое изображение без прямых черных линий.

2.3.2 Тестирование уточненной разметки

При тестировании уточненной разметки должна проверяться корректность вывода уточненного состояния для первичной разметки на основе описанных в данном разделе ранее правил.

Таким образом, классы эквивалентности при тестировании уточненной разметки следующие:

- «Не определено» → «Текст»,
- «Не определено» → «Таблица»,
- «Не определено» → «Листинг»,
- «Не определено» → «Схема алгоритма»,
- «Не определено» → «Рисунок»,
- «Не определено» → «График»,
- «Не определено» → «Не определено»,
- «Немного текста» → «Текст»,
- «Много текста» → «Таблица»,
- «Много текста» → «Листинг»,
- «Много текста» → «Текст»,
- «Цвет» → «График»,
- «Цвет» → «Не определено»,
- «Цвет» → «Рисунок»,
- «Черная линия средней длины» → «Текст»,
- «Черная линия средней длины» → «Рисунок»,
- «Черная линия средней длины» → «График»,
- «Черная линия средней длины» → «Не определено»,

- «Черная линия средней длины» → «Схема алгоритма»,
- «Длинная черная линия» → «Таблица»,
- «Длинная черная линия» → «Листинг»,
- «Длинная черная линия» → «Схема алгоритма»,
- «Длинная черная линия» → «График»,
- «Длинная черная линия» → «Не определено»,
- «Длинная черная линия» → «Рисунок».

Для каждого класса эквивалентности следует подобрать тестовые сценарии, покрывающие все правила, которые могут быть применены в конкретном классе эквивалентности для получения соответствующей уточненной разметки.

2.3.3 Тестирование объединенной разметки

При тестировании объединенной разметки должна проверяться корректность каждого этапа создания объединенной разметки.

Классы эквивалентности в таком случае соответствуют этапам создания объединенной разметки:

- Слияние небольших фоновых сегментов с ближайшим наибольшим сегментом;
- Слияние фоновых сегментов с соседними, если соседние сегменты имеют один и тот же класс;
- Смена класса небольших фоновых сегментов на «Не определено»;
- Слияние небольших «Не определенных» сегментов с наибольшими соседними сегментами;
- Объединение сегментов после слияния.

2.4 Структура разрабатываемого ПО

Структура разрабатываемого ПО представлена на рисунке 11 ниже.

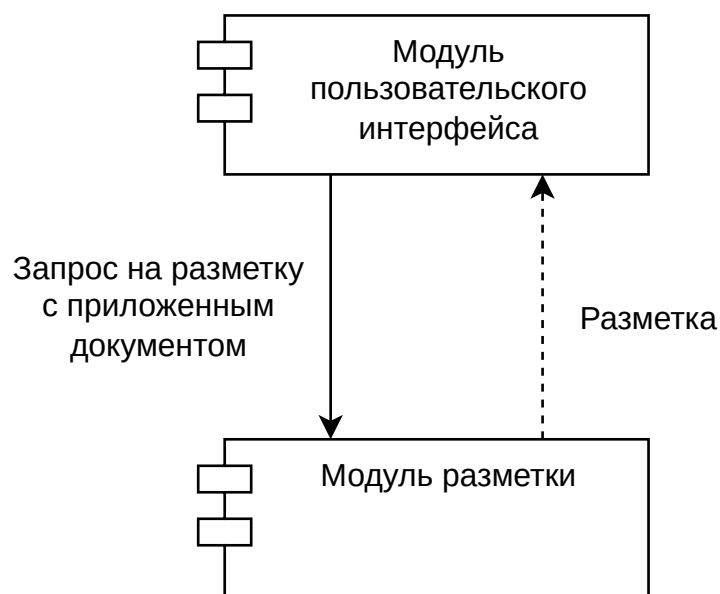


Рисунок 11 – Структура разрабатываемого ПО

Вывод

В данном разделе были приведены требования и ограничения разрабатываемого метода, были описаны основные этапы разрабатываемого метода, сценарии тестирования, классы эквивалентности тестов, а также структура разрабатываемого ПО.

3 Технологический раздел

В данном разделе будут выбраны средства реализации программного обеспечения, описаны основные функции разработанного программного обеспечения, приведены результаты тестирования, примеры пользовательского интерфейса, демонстрация работы программы, а также руководство пользователя.

3.1 Выбор средств реализации

Для реализации метода выделения составных частей научного текста на основе анализа распределения пикселей в сканирующей строке был выбран Python [12] по следующим причинам:

- язык программирования позволяет быстро создавать рабочие прототипы в связи с наличием большого количества библиотек;
- имеются библиотеки для работы с PDF документами и изображениями;
- есть опыт работы с данным языком программирования.

Для управления зависимостями Python-проекта и организации изолированных виртуальных окружений, обеспечивая воспроизводимость и независимость среды выполнения, был выбран `uv` [13].

Для работы с массивами данных (в том числе сканирующей строкой пикселей) была выбрана библиотека NumPy [14].

Для работы с изображениями и PDF были выбраны библиотеки PIL [15] и PyMuPDF [16] соответственно, так как в процессе работы программы используются как PDF документы, так и изображения, а PyMuPDF позволяет преобразовывать результат рендеринга PDF в формат, поддерживаемый PIL, который, в свою очередь, может быть преобразован в NumPy массив и использован для дальнейшей обработки.

Для создания веб-интерфейса была выбрана библиотека Gradio [17], так как данная библиотека позволяет быстро создать рабочий прототип графического интерфейса к программе.

В качестве среды разработки был выбран Neovim [18] по следующим причинам:

- данный текстовый редактор позволяет редактировать файлы с исходным кодом программы;
- быстрая и удобная навигация по файлам проекта;
- есть опыт использования данного текстового редактора.

3.2 Реализация программного обеспечения

В данном подразделе будут описаны входные и выходные данные основных функций разработанного программного обеспечения, а также приведены их листинги.

3.2.1 Точка входа

Входные данные:

- `pdf_path` — строка, путь к PDF документу,
- `json_path` — строка, имя создаваемого файла разметки;
- `markup_type` — целое число от 0 до 3, тип разметки; 0 — построчная разметка, 1 — первичная разметка, 2 — уточненная разметка, 3 — объединенная разметка.
- `num_workers` — целое число, число рабочих процессов, так как документ быстрее обрабатывать параллельно;
- `page_indices` — массив целых чисел, конкретные страницы для обработки, если они указаны; по умолчанию обрабатывается весь документ.

Выходные данные: Сохраненная в указанном JSON файле разметка. Исходный код данной функции приведен в листинге 1 ниже.

Листинг 1 – Точка входа, функция main

```

1 def main(pdf_path, json_path, markup_type, num_workers=8,
2         page_indices=None):
3
4     if page_indices is None:
5         doc = fitz.open(pdf_path)
6         total_pages = doc.page_count
7         page_indices = list(range(total_pages))
8         doc.close()
9
10    # Делит диапазон страниц на N примерно равных кусков.
11    chunks = chunk_indices(len(page_indices), num_workers)
12    pages_chunks = [ [page_indices[i] for i in chunk] for chunk
13                    in chunks ]
14
15    # .. (Создание шкалы прогресса и запуск ее в отдельном потоке)
16
17    try:
18        with ProcessPoolExecutor(max_workers=num_workers) as
19            executor:
20                futures = [executor.submit(process_pages, pdf_path,
21                                         chunk, markup_type,
22                                         queue) for chunk in
23                        pages_chunks]
24
25                for future in futures:
26                    results.extend(future.result())
27
28    finally:
29        stop_signal.set()
30        thread.join()
31
32    results.sort(key=lambda x: x["page"])
33
34    with open(json_path, "w", encoding="utf-8") as f:
35        json.dump(results, f, ensure_ascii=False, indent=4)
36
37    return results

```

3.2.2 Функция обработки страниц

Данная функция запускается для каждого процесса, после чего результаты обработки страниц объединяются.

Входные данные:

- `pdf_path` — строка, путь к PDF документу,
- `page_range` — целочисленный массив, страницы, которые требуется обработать конкретному процессу;
- `markup_type` — целое число от 0 до 3, тип разметки;
- `queue` — объект типа очередь, нужен исключительно для обновления строки прогресса в терминале.

Выходные данные: Частичная разметка в виде Python-словаря указанного в листинге 2 формата.

Исходный код данной функции приведен в листинге 2 ниже.

Листинг 2 – Функция обработки страниц

```
1 from fast import segdoc as sd
2 def process_pages(pdf_path, page_range, markup_type, queue=None):
3     """Функция для обработки поддиапазона страниц."""
4     partial_results = []
5     for i in page_range:
6         image = page_to_image(pdf_path, i)
7         image_np = np.array(image)
8         markup = sd(image_np, markup_type)
9         assert markup is not None
10        del image, image_np
11        partial_results.append({
12            "page": i + 1,
13            "segments": [
14                {"y_start": s[0], "y_end": s[1], "label": s[2]}
15                for s in markup
16            ]
17        })
18        if queue: # Нужно для шкалы прогресса
19            queue.put(1)
20    return partial_results
```

3.2.3 Интерфейсная функция для разметки

Входные данные:

- `image` — массив NumPy, представляющий изображение страницы;
- `v` — выбираемый тип разметки.

Выходные данные: Python-словарь, представляющий разметку.

Исходный код данной функции приведен в листинге 3 ниже.

Листинг 3 – Интерфейсная функция для разметки

```
1 def fn(sl):
2     return extract_line_features(sl, 0, None, WHITE_THRESH,
3         GRAY_TOL)
4
5 def segdoc(image, v):
6     sd = segment_document
7     if v == 0:
8         markup = segment_document_raw(image, fn)
9         return merge_segments(markup)
10
11     if v == 1:
12         markup = sd(image, fn, True)
13         return markup
14
15     if v == 2:
16         markup = sd(image, fn, False)
17         return markup
18
19     if v == 3:
20         markup = sd(image, fn, False)
21         return merge(markup)
```

3.2.4 Выделение характеристик строки пикселей

Входные данные:

- `sl` — массив NumPy, представляющий сканирующую строку;
- `left_margin` — целое число, индекс окончания левого поля документа;

- `right_margin` — целое число, индекс начала правого поля документа;
- `white_thresh` — целое число, пороговое значение для белого цвета, пиксели со значением выше него считаются белыми;
- `gray_tol` — целое число, пиксель считается цветным, если размах значений его каналов меньше `gray_tol`.

Выходные данные: Структура `LineFeatures` (см. Листинг 4), описывающая характеристики распределения пикселей в сканирующей строке.

Исходный код данной функции приведен в листинге 4 ниже.

Листинг 4 – Структура `LineFeatures`

```

1 @dataclass
2 class LineFeatures:
3     count_white: int
4     count_color: int
5     count_gray: int
6     comp_lengths: List[int]
7     gap_lengths: List[int]
8     gray_comp_lengths: List[int]
9     color_comp_lengths: List[int]
10    first_nonwhite_index: int | None

```

Исходный код данной функции приведен в листинге 18 приложения А.

3.2.5 Функция создания построчной разметки

Входные данные:

- `image` — массив `NumPy`, представляющий изображение страницы;
- `line_feature_func` — функция для выделения характеристик распределения пикселей в сканирующей строке, возвращает `LineFeatures`.

Выходные данные: Построчная разметка в виде Python-словаря (использует названия состояний конечного автомата для классификации).

Исходный код данной функции приведен в листинге 5 ниже.

Листинг 5 – Функция создания построчной разметки (используется для отладки)

```
1 def segment_document_raw(  
2     image: np.ndarray ,  
3     line_feature_func: Callable[[np.ndarray], LineFeatures] ,  
4 ):  
5     results = []  
6     height = image.shape[0]  
7     for y in range(1, height):  
8         line = image[y:y+1]  
9         feat = line_feature_func(line)  
10        state = classify_line(feat)  
11        result = (y, y+1, StateNames[state])  
12        results.append(result)  
13    result = (height-1, height ,  
14              StateNames[classify_line(  
15                  line_feature_func(  
16                      image[height-1:height]))])  
17    results.append(result)  
18  
19    return results
```

3.2.6 Функция классификации строки

Входные данные: `feat` — структура `LineFeatures` с характеристиками обрабатываемой строки.

Выходные данные: Целое число, класс данной строки в терминах состояний конечного автомата.

Исходный код данной функции приведен в листинге 20 приложения А.

3.2.7 Первичная и уточненная разметки

Входные данные:

- `image` — массив `NumPy`, представляющий изображение страницы;
- `line_feature_func` — функция для выделения характеристик распределения пикселей в сканирующей строке, возвращает `LineFeatures`;

- raw — флаг создания первичной разметки, если установлен, будет создана первичная разметка, иначе — уточненная.

Выходные данные: Python-словарь, представляющий первичную или уточненную разметку.

Исходный код данной функции приведен в листинге 23 приложения А.

3.2.8 Функция обновления состояния КА

Входные данные:

- state — целое число, текущее состояние конечного автомата,
- feat — структура LineFeatures с характеристиками текущей строки.

Выходные данные: Обновленное состояние конечного автомата.

Листинг 6 – Функция обновления состояния конечного автомата

```
1 def update_state(state: int, feat: LineFeatures):  
2     inferred_state = classify_line(feat)  
3     return FSM[state][inferred_state]
```

Конечный автомат

Исходный код конечного автомата и связанной с ним структуры State приведен в листингах 25 и 27 приложения А соответственно.

3.2.9 Функция классификации сегмента

Входные данные:

- state — целое число, текущее состояние конечного автомата;
- sd — структура SegmentData, содержащая накопленную информацию о текущем сегменте;
- raw — флаг использования первичной разметки.

Выходные данные: Целое число, уточненный класс сегмента в терминах уточненной разметки.

Исходный код данной функции приведен в листинге 7 ниже.

Листинг 7 – Функция классификации сегмента

```
1 def classify_segment(state: int, sd: SegmentData, raw: bool =  
    False):  
2     handlers = {  
3         State.UNDEFINED: handle_undefined,  
4         State.BACKGROUND: handle_background,  
5         State.FEW_TEXT: handle_few_text,  
6         State.MANY_TEXT: handle_many_text,  
7         State.COLOR: handle_color,  
8         State.MEDIUM_BLACK_LINE: handle_medium_black_line,  
9         State.LONG_BLACK_LINE: handle_long_black_line,  
10    }  
11  
12    handler = handlers.get(state)  
13    assert handler is not None  
14  
15    if raw:  
16        return StateNames[state]  
17  
18    return handler(sd)
```

Исходный код для функций уточнения класса сегмента `handle_undefined`, `handle_background`, `handle_few_text`, `handle_many_text`, `handle_color`, `handle_medium_black_line`, `handle_long_black_line` приведен в листингах 28 – 39 приложения А.

3.2.10 Обновление информации о сегменте

Входные данные:

- `sd` — структура `SegmentData` информации о сегменте;
- `prev_feat`, `feat` — структуры `LineFeatures` информации о предыдущей и текущей строках соответственно;
- `line` — массив `NumPy`, представляющий сканирующую строку.

Выходные данные: Обновленные значения полей `sd`.

Исходный код для функции обновления информации о сегменте приведен в листинге 42 приложения А.

3.2.11 Слияние смежных сегментов одного класса

Входные данные: `arr` — массив разметки, состоящий из кортежей типа (начало сегмента, конец сегмента, класс сегмента).

Выходные данные: Массив разметки, состоящий из кортежей типа (начало сегмента, конец сегмента, класс сегмента) такой, что несколько подряд идущих сегментов одного класса в массиве `arr` соответствуют одному сегменту в выходном массиве.

Исходный код данной функции приведен в листинге 8 ниже.

Листинг 8 – Слияние смежных сегментов одного класса в один

```
1 def merge_segments(arr):
2     merged = []
3     (current_start, current_end, current_class) = arr[0]
4     for i in range(1, len(arr)):
5         start, end, cls = arr[i]
6         if current_class == cls and current_end == start:
7             current_end = end
8         else:
9             merged.append([current_start, current_end,
10                            current_class])
11             current_start, current_end, current_class = start,
12                            end, cls
13     merged.append([current_start, current_end, current_class])
14     return merged
```

3.2.12 Объединенная разметка

Входные данные: `markup` — массив разметки, состоящий из кортежей типа (начало сегмента, конец сегмента, класс сегмента).

Выходные данные: Обновленная разметка в виде массива, состоящего из кортежей типа (начало сегмента, конец сегмента, класс сегмента).

Исходный код функции создания объединенной разметки приведен в листинге 43 приложения А.

3.3 Результаты тестирования

После разработки программного обеспечения было проведено тестирование модуля разметки. Тестирование проводилось в соответствии с описанными в конструкторском разделе классами эквивалентности. При тестировании уточненной разметки классы эквивалентности были объединены в более крупные — на основе получаемого класса разметки на выходе. Результаты тестирования представлены в листинге 9 ниже.

Листинг 9 – Результаты тестирования

```
1 ===== short test summary info =====
2 PASSED test.py::TestPrimaryMarkup::test_background
3 PASSED test.py::TestPrimaryMarkup::test_few_text
4 PASSED test.py::TestPrimaryMarkup::test_undefined
5 PASSED test.py::TestPrimaryMarkup::test_many_text
6 PASSED test.py::TestPrimaryMarkup::test_long_black_line
7 PASSED test.py::TestPrimaryMarkup::test_medium_black_line
8 PASSED test.py::TestPrimaryMarkup::test_color
9 PASSED test.py::TestSpecifiedMarkup::test_to_text
10 PASSED test.py::TestSpecifiedMarkup::test_to_table
11 PASSED test.py::TestSpecifiedMarkup::test_to_code
12 PASSED test.py::TestSpecifiedMarkup::test_to_diagram
13 PASSED test.py::TestSpecifiedMarkup::test_to_figure
14 PASSED test.py::TestSpecifiedMarkup::test_to_plot
15 PASSED test.py::TestSpecifiedMarkup::test_to_undefined
16 PASSED test.py::TestMergedMarkup::test_merge_little_bg
17 PASSED test.py::TestMergedMarkup::test_merge_bg_interpolate
18 PASSED test.py::TestMergedMarkup::test_swap_bg_to_undefined
19 PASSED test.py::TestMergedMarkup::test_merge_undefined
20 PASSED test.py::TestMergedMarkup::test_merge_segments
```

3.4 Пользовательский интерфейс

На рисунках 12 – 14 ниже представлен порядок взаимодействия пользователя с графическим веб-интерфейсом.

Разметка PDF-документа

Загрузите PDF-документ

Drop File Here

- or -

Click to Upload

Тип разметки

Построчная

✓ Построчная

Первичная

Уточненная

Объединенная

Разметить

JSON разметка

Размеченный PDF




Use via API  · Built with Gradio  · Settings 

Рисунок 12 – Выбор типа разметки

Разметка PDF-документа

Загрузите PDF-документ

index.pdf


2.4 MB

Тип разметки


Уточненная

Разметить

JSON разметка


processing | 4.0s

Размеченный PDF


processing | 4.0s




Use via API  · Built with Gradio  · Settings 

Рисунок 13 – Разметка в процессе

Разметка PDF-документа

Загрузите PDF-документ

index.pdf2.4 MB

Тип разметки

Уточненная

Разметить

JSON разметка

index.markup.json103.5 KB

Размеченный PDF

index.annotated.pdf2.9 MB

Use via API · Built with Gradio · Settings

Рисунок 14 – Разметка завершена — пользователь может скачать разметку в формате JSON и/или PDF документ для ее визуализации

3.5 Демонстрация работы программы

На рисунках 15 – 18 ниже приведена демонстрация работы программы. Представлены результаты построчной (строке пикселей назначается класс на основе распределения пикселей в ней, без использования конечного автомата), первичной, уточненной и объединенной разметок.

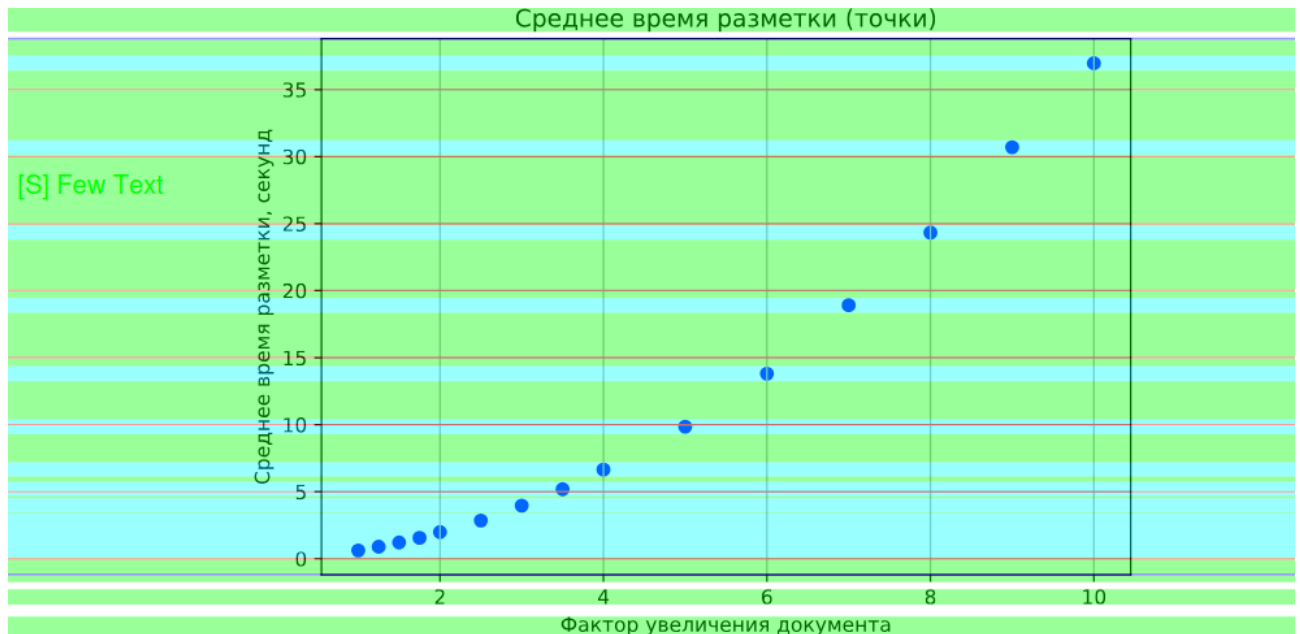


Рисунок 5 – Зависимость времени работы алгоритма от коэффициента увеличения документа при его преобразовании в изображение, исходные данные

Для определения зависимости времени разметки от фактора увеличения документа проведем степенную аппроксимацию на полученных данных.

```
[S] Few Text 1 import numpy as np
2 x = np.array([ ... ])
3 y = np.array([ ... ])
4 log_x = np.log(x)
5 log_y = np.log(y)
6 b, log_c = np.polyfit(log_x, log_y, 1)
7 c = np.exp(log_c)
8 print(f'y={c:.3f}*x^{b:.3f}') # y = 0.583 * x^1.782
```

Листинг 3 – Степенная аппроксимация

Рисунок 15 – Пример построчной разметки

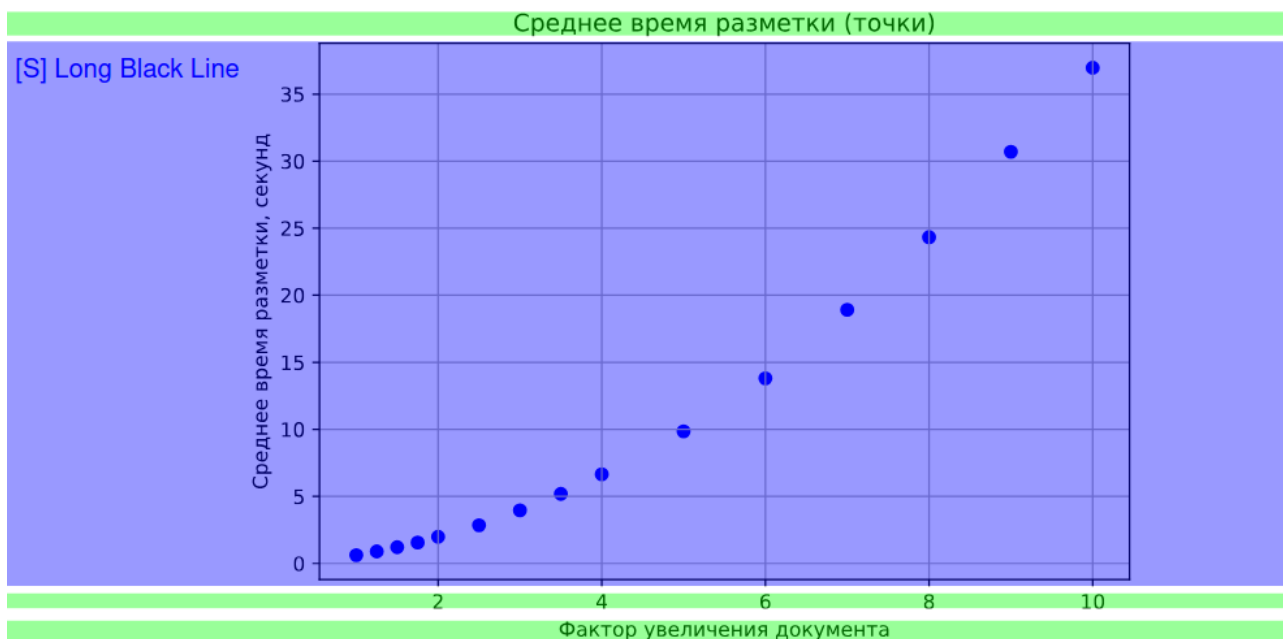


Рисунок 5 – Зависимость времени работы алгоритма от коэффициента увеличения документа при его преобразовании в изображение, исходные данные

Для определения зависимости времени разметки от фактора увеличения документа проведем степенную аппроксимацию на полученных данных.

```
[S] Long Black Line
import numpy as np
2 x = np.array([ ... ])
3 y = np.array([ ... ])
4 log_x = np.log(x)
5 log_y = np.log(y)
6 b, log_c = np.polyfit(log_x, log_y, 1)
7 c = np.exp(log_c)
8 print(f'y={c:.3f}*x^{b:.3f}') # y = 0.583 * x^1.782
```

Листинг 3 – Степенная аппроксимация

Рисунок 16 – Пример первичной разметки

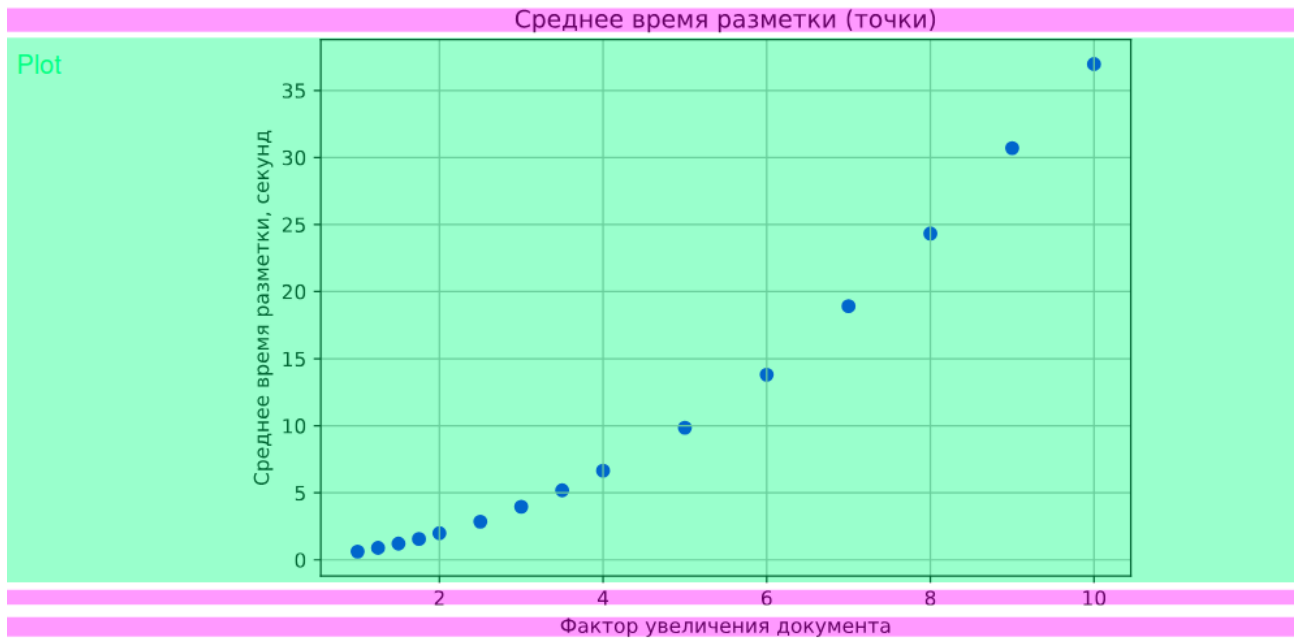


Рисунок 5 – Зависимость времени работы алгоритма от коэффициента увеличения документа при его преобразовании в изображение, исходные данные

Для определения зависимости времени разметки от фактора увеличения документа проведем степенную аппроксимацию на полученных данных.

Code

```

1 import numpy as np
2 x = np.array([ ... ])
3 y = np.array([ ... ])
4 log_x = np.log(x)
5 log_y = np.log(y)
6 b, log_c = np.polyfit(log_x, log_y, 1)
7 c = np.exp(log_c)
8 print(f'y={c:.3f}*x^{b:.3f}') # y = 0.583 * x^1.782

```

Листинг 3 – Степенная аппроксимация

Рисунок 17 – Пример уточненной разметки

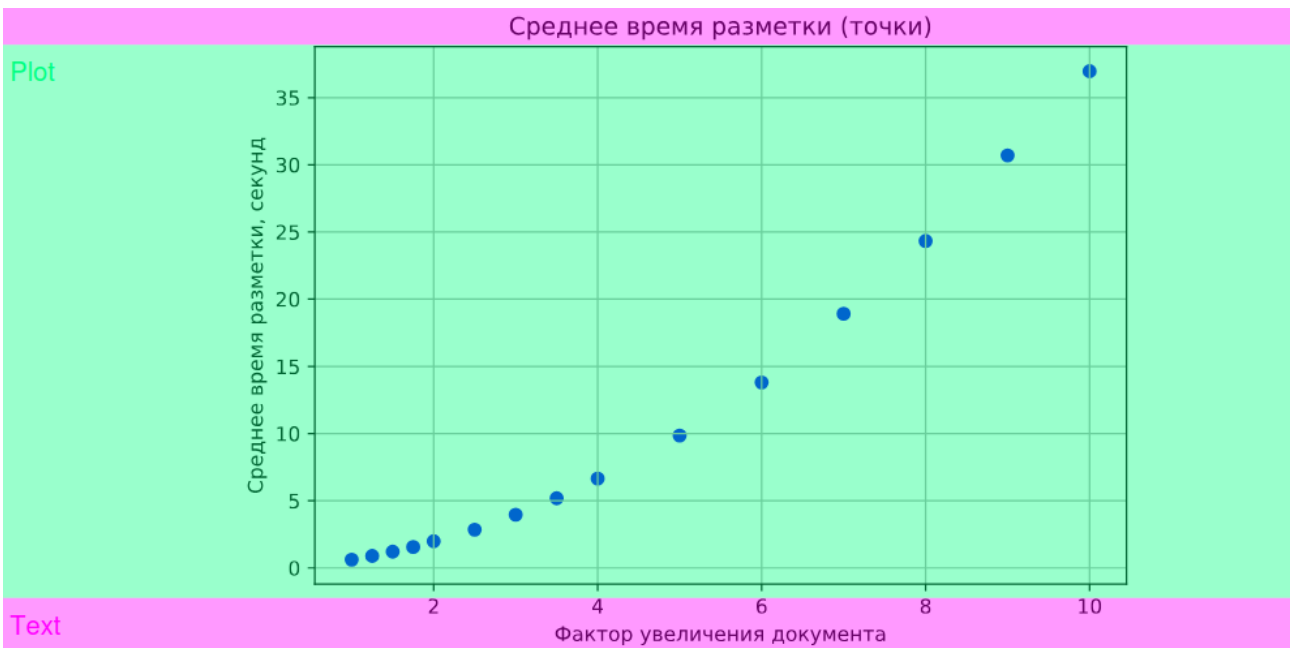


Рисунок 5 – Зависимость времени работы алгоритма от коэффициента увеличения документа при его преобразовании в изображение, исходные данные

Для определения зависимости времени разметки от фактора увеличения документа проведем степенную аппроксимацию на полученных данных.

Code

```
1 import numpy as np
2 x = np.array([ ... ])
3 y = np.array([ ... ])
4 log_x = np.log(x)
5 log_y = np.log(y)
6 b, log_c = np.polyfit(log_x, log_y, 1)
7 c = np.exp(log_c)
8 print(f'y={c:.3f}*x^{b:.3f}') # y = 0.583 * x^1.782
```

Листинг 3 – Степенная аппроксимация

Рисунок 18 – Пример объединенной разметки

В листингах 10 и 11 ниже представлен результат работы программы при запуске из терминала.

Листинг 10 – Запуск и результат работы разметки из терминала

```
1 uv run main.py ~/index.pdf x.json 3
2 Pages processed: 31/31 [00:06<00:00, 5.11 it/s]
```

Листинг 11 – Запуск и результат применения разметки к PDF из терминала

```
1 uv run apply.py ~/index.pdf x.json output.pdf
2 Annotating pages: 31/31 [00:00<00:00, 383.41 it/s]
3 Сохранено в: output.pdf
```

3.6 Руководство пользователя

Для запуска программы требуется установить uv командой, показанной в листинге 12 ниже.

Листинг 12 – Установка uv

```
1 curl -LsSf https://astral.sh/uv/install.sh | sh
```

После установки uv, следует установить зависимости приложения. Для этого, находясь в директории с программой, следует прописать команду, показанную в листинге 13 ниже.

Листинг 13 – Установка зависимостей

```
1 uv sync
```

Графический интерфейс запускается локально командой, показанной в листинге 14 ниже.

Листинг 14 – Запуск графического веб-интерфейса

```
1 uv run webgui.py
```

Также можно создать разметку без использования графического интерфейса. Для этого используется скрипт main.py, инструкция к использованию которого приведена в листинге 15 ниже.

Листинг 15 – Запуск скрипта для создания разметки

```
1 Использование: main.py [-w КОЛИЧЕСТВО_ПРОЦЕССОВ]
2                       [-p СТРАНИЦЫ]
3                       pdf_path json_path {0,1,2,3}
4
5 Сегментировать PDF страницы и экспортировать разметку в JSON.
6
7 Позиционные аргументы:
8   pdf_path           Путь ко входному PDF файлу
9   json_path          Путь к выходному JSON файлу
10  {0,1,2,3}          Тип разметки: 0 - построчная разметка ,
11                                     1 - первичная разметка ,
12                                     2 - уточненная разметка ,
13                                     3 - объединенная разметка
14
15 Опции:
16   -w, --workers КОЛИЧЕСТВО_ПРОЦЕССОВ
17                                   Количество рабочих процессов
18                                   (по умолчанию: 8)
19   -p, --pages СТРАНИЦЫ Диапазоны страниц для обработки ,
20                                   например, '1-3,5,7-9'
```

Для применения разметки к документу без графического интерфейса используется скрипт `apply.py`, инструкция к использованию которого приведена в листинге 16 ниже.

Листинг 16 – Запуск скрипта для применения разметки к PDF документу

```
1 Использование: apply.py [-h] pdf_path json_path output_path
2
3 Разметить PDF цветными сегментами из JSON файла разметки.
4
5 Позиционные аргументы:
6   pdf_path          Путь ко входному PDF файлу
7   json_path         Путь к JSON файлу с разметкой
8   output_path       Путь для сохранения размеченного PDF
```

Вывод

В данном разделе были выбраны средства реализации программного обеспечения, описаны основные функции разработанного программного обеспечения, приведены результаты тестирования, примеры пользовательского интерфейса, демонстрация работы программы, а также руководство пользователя.

4 Исследовательский раздел

В данном разделе будет приведено описание исследования, технические характеристики устройства, на котором выполнялись замеры времени, будет представлены и проанализированы результаты исследования.

Целью данного исследования является определение зависимостей времени выполнения разметки и максимального объема используемой памяти в зависимости от количества процессов, выполняющих разметку.

Предположительно, разметка будет производиться наиболее эффективно в случае, когда количество рабочих процессов совпадает с количеством физических ядер процессора. Также тип разметки не будет влиять на максимальный объем используемой памяти.

4.1 Описание исследования

В листинге 17 ниже и листингах 46 – 50 приложения А представлен исходный код скриптов, запускаемых для выполнения замеров времени и максимального объема используемой во время разметки памяти.

Листинг 17 – bench.py — Скрипт для запуска теста производительности

```
1 import subprocess
2
3 PDF_PATH = "/home/rukost/report.pdf"
4 MARKUP_TYPE = "0"
5 PAGES = "1-50"
6
7 for workers in [1, 2, 4, 8, 16, 32, 64]:
8     print(f"\n---_Benchmark_with_{workers}_workers_---")
9     result = subprocess.run([
10         "python", "main-bench.py",
11         PDF_PATH,
12         MARKUP_TYPE,
13         "-w", str(workers),
14         "-p", PAGES
15     ], capture_output=True, text=True)
16     print(result.stdout)
```

Для получения данных по разным типам разметок, при запуске скрипта `bench.py` изменялся параметр `MARKUP_TYPE`.

Скрипт запускался десять раз для каждого типа разметки, после чего полученные значения усреднялись.

4.2 Технические характеристики

Технические характеристики устройства, на котором выполнялись замеры времени, представлены ниже.

- 1) Процессор: AMD Ryzen 7 4700U 2.0 ГГц [19], 8 физических ядер, 8 потоков;
- 2) Оперативная память: 8 ГБ, DDR4, 3200 МГц;
- 3) Операционная система: Arch [20];
- 4) Версия ядра: 6.14.9.

При выполнении замеров времени ноутбук был подключен к сети электропитания, был запущен терминал Alacritty [21], в котором выполнялся тест производительности.

4.3 Результаты исследования

В данном подразделе будут приведены и проанализированы результаты исследования времени выполнения и максимального объема памяти, затрачиваемой в процессе построения, первичной, уточненной и объединенной разметок.

4.3.1 Исследование времени выполнения

В таблице 9 и на рисунке 19 представлены полученные в ходе исследования данные. Данные, приведенные в таблице являются усредненными значениями из десяти замеров.

Таблица 9 – Время выполнения разметки в зависимости от ее типа и количества процессов, с

Кол-во процессов	Построчн.	Первичн.	Уточн.	Объед.
1	31.80	41.21	41.29	41.31
2	16.70	21.72	21.96	22.03
4	10.16	13.05	13.11	13.19
8	7.03	8.71	8.72	8.96
16	7.29	8.83	8.96	9.23
32	7.40	9.23	9.24	9.47
64	7.82	9.67	9.73	9.79

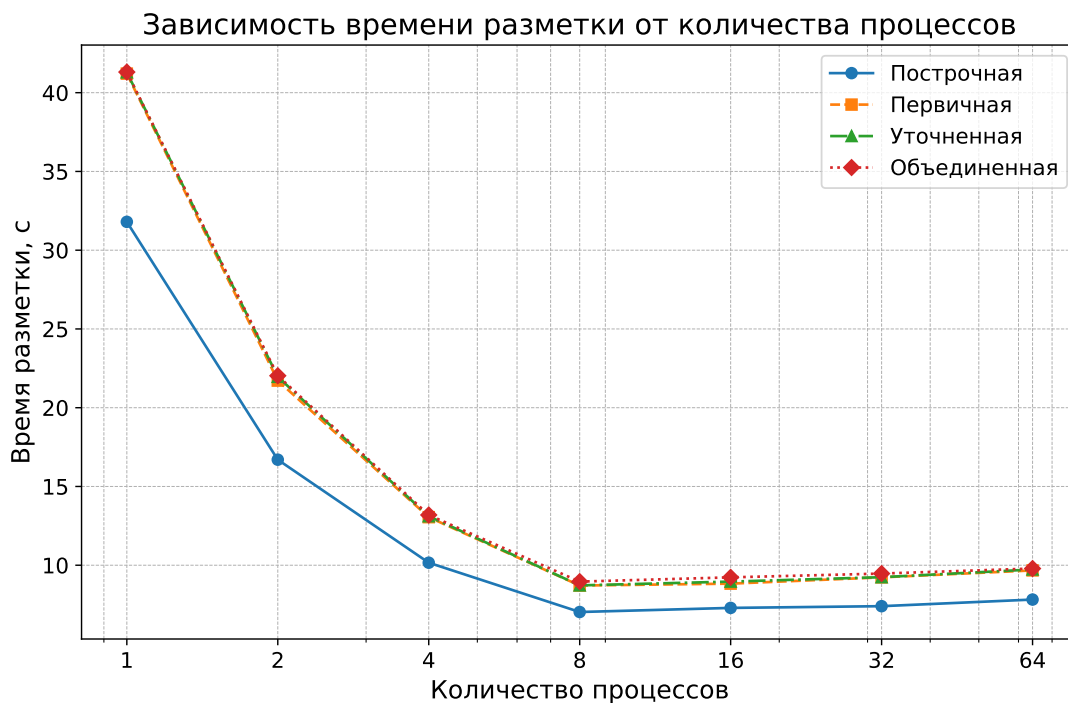


Рисунок 19 – Время выполнения разметки в зависимости от ее типа и количества процессов

По полученным данным можно сделать следующие выводы:

- Среди разметок наиболее быстрой является построчная, так как для нее нужно меньше вычислений;
- Увеличение числа процессов, используемых для разметки, позволяет ускорить процесс разметки;
- Наибольшая скорость разметки наблюдается, когда число процессов соответствует числу физических ядер процессора.

4.3.2 Исследование затрачиваемой памяти

В таблице 10 и на рисунке 20 представлены полученные в ходе исследования данные. Данные, приведенные в таблице являются усредненными значениями из десяти замеров.

Таблица 10 – Максимальное количество памяти, используемое в ходе разметки в зависимости от ее типа и количества процессов, МБ

Кол-во процессов	Построчн.	Первичн.	Уточн.	Объед.
1	247.39	246.58	246.56	245.09
2	374.05	373.12	373.26	373.16
4	627.94	625.74	627.79	627.11
8	1085.41	1096.95	1098.18	1122.85
16	1933.34	1934.43	1932.78	1935.12
32	3757.83	3760.95	3758.53	3758.38
64	6437.75	6444.76	6438.00	6440.77

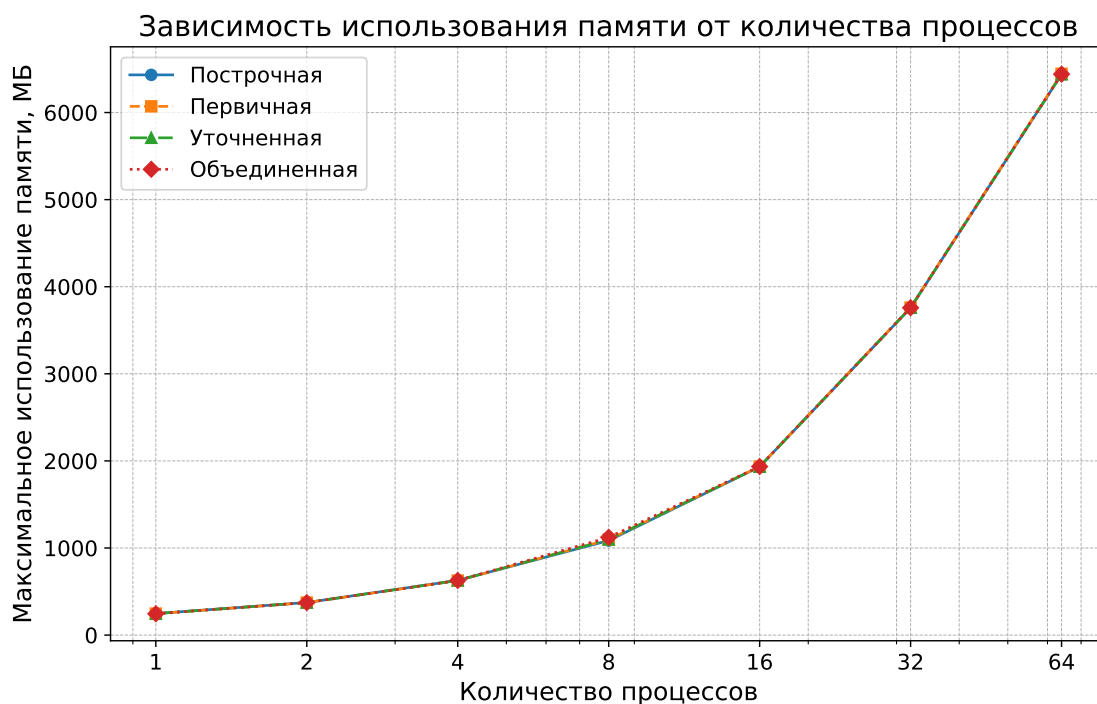


Рисунок 20 – Максимальное количество памяти, используемое в ходе разметки в зависимости от ее типа и количества процессов

По полученным данным можно сделать следующие выводы:

- Максимальное количество памяти, используемое при разметке, не зависит от типа разметки;
- Основным фактором, влияющим на максимальное количество памяти, используемое при разметке, является количество создаваемых в ходе разметки процессов;
- Объем используемой памяти увеличивается прямо пропорционально количеству рабочих процессов.

Вывод

Таким образом, в ходе исследования было выявлено, что разметка наиболее эффективна по времени, когда количество рабочих процессов соответствует количеству физических ядер процессора; самой быстрой является построчная разметка, и объем используемой памяти увеличивается пропорционально количеству рабочих процессов.

ЗАКЛЮЧЕНИЕ

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Clarivate Unveils Journal Citation Reports 2023 — A Trusted Resource to Support Research Integrity and Promote Accurate Journal Evaluation [Электронный ресурс]. — URL: <https://clarivate.com/news/clarivate-unveils-journal-citation-reports-2023-a-trusted-resource-to-support-research-integrity-and-promote-accurate-journal-evaluation/> (дата обращения: 03.06.2025).
2. Cui L. Document AI: Benchmarks, Models and Applications / Xu Y., Lv T., Wei F. 2021. — URL: <https://arxiv.org/abs/2111.08609>.
3. Shehzadi T. A Hybrid Approach for Document Layout Analysis in Document images / Stricker D., Afzal M.Z. 2024. — URL: <https://arxiv.org/abs/2404.17888>.
4. Bhowmik S. Document Layout Analysis. — Springer Singapore, 2023 — 86 с.
5. Binmakhashen G.M., Mahmoud S.A. Document Layout Analysis: A Comprehensive Survey // ACM Comput. Surv. 2019. Т. 52, № 6.
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. Document image analysis: A primer / L. O’Gorman, V. Govindaraju // 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.
12. Python — A programming language that lets you work quickly and integrate systems more effectively [Электронный ресурс]. – URL: <https://www.python.org/> (дата обращения: 26.05.2025).
13. uv — An extremely fast Python package and project manager, written in Rust [Электронный ресурс]. – URL: <https://docs.astral.sh/uv/> (дата обращения: 26.05.2025).
14. NumPy — The fundamental package for scientific computing with Python [Электронный ресурс]. – URL: <https://numpy.org/> (дата обращения: 26.05.2025).
15. Pillow — The Python Imaging Library adds image processing capabilities to your Python interpreter [Электронный ресурс]. – URL: <https://pypi.org/project/pillow/> (дата обращения: 26.05.2025).
16. PyMuPDF — A high-performance Python library for data extraction, analysis, conversion and manipulation of PDF (and other) documents [Электронный ресурс].
17. Gradio — Build and share delightful machine learning apps [Электронный ресурс]. – URL: <https://www.gradio.app/> (дата обращения: 26.05.2025).
18. Neovim — Hyperextensible Vim-based text editor [Электронный ресурс]. – URL: <https://neovim.io> (дата обращения: 26.05.2025).
19. AMD Ryzen 7 4700U [Электронный ресурс]. – URL: <https://www.amd.com/en/product/9096> (дата обращения: 26.05.2025).

20. Arch Linux — A simple, lightweight distribution [Электронный ресурс]. — URL: <https://archlinux.org/> (дата обращения: 26.05.2025).
21. Alacritty — A fast, cross-platform, OpenGL terminal emulator [Электронный ресурс]. — URL: <https://github.com/alacritty/alacritty> (дата обращения: 26.05.2025).

ПРИЛОЖЕНИЕ А

Листинг 18 – Функция для выделения характеристик строки пикселей
(часть 1)

```
1 def extract_line_features(  
2     scanline: np.ndarray,  
3     left_margin: int = 0,  
4     right_margin: Optional[int] = None,  
5     white_thresh: int = WHITE_THRESH,  
6     gray_tol: int = GRAY_TOL,  
7 ) -> LineFeatures:  
8     end = right_margin if right_margin is not None else  
9         scanline.shape[0]  
10    segment = scanline[left_margin:end]  
11    if segment.size == 0:  
12        raise ValueError("Segment_length_is_zero. Проверьте_  
13        left_margin/right_margin.")  
14  
15    (mask_white,  
16    mask_nonwhite,  
17    mask_color,  
18    mask_gray) = get_masks(segment, white_thresh, gray_tol)  
19  
20    # Счётчики  
21    count_white = int(np.sum(mask_white))  
22    count_color = int(np.sum(mask_color))  
23    count_gray = int(np.sum(mask_gray))  
24  
25    def rle_lengths(mask: np.ndarray) -> List[int]:  
26        if mask.ndim != 1:  
27            mask = mask.ravel()  
28        padded = np.pad(mask.astype(np.int8), (1, 1),  
29            mode='constant')  
30        diffs = np.diff(padded)  
31        run_starts = np.where(diffs == 1)[0]  
32        run_ends = np.where(diffs == -1)[0]  
33        return (run_ends - run_starts).tolist()
```

Листинг 19 – Функция для выделения характеристик строки пикселей
(часть 2)

```
1  comp_lengths = rle_lengths(mask_nonwhite)
2  gray_comp_lengths = rle_lengths(mask_gray)
3  color_comp_lengths = rle_lengths(mask_color)
4  # gap_lengths: разности между концами и началом следующих комп
   # онент
5  padded_nonwhite = np.pad(mask_nonwhite.astype(np.uint8), (1,
   1))
6  diffs = np.diff(padded_nonwhite)
7  starts = np.where(diffs == 1)[0]
8  ends = np.where(diffs == -1)[0]
9
10 # Защита от несогласованных размеров: отбрасываем лишние значе
   # ния
11 min_len = min(len(starts), len(ends))
12 starts = starts[:min_len]
13 ends = ends[:min_len]
14
15 # Если после обрезки осталось хотя бы 2 компонента, можно иска
   # ть промежутки
16 if len(starts) > 1:
17     gap_lengths = (starts[1:] - ends[:-1]).tolist()
18 else:
19     gap_lengths = []
20
21 # Первый не-белый пиксель
22 nonwhite_idx = np.flatnonzero(mask_nonwhite)
23 first_nonwhite_index = int(nonwhite_idx[0] + left_margin) if
   nonwhite_idx.size > 0 else None
24
25 return LineFeatures(
26     count_white,
27     count_color,
28     count_gray,
29     comp_lengths,
30     gap_lengths,
31     gray_comp_lengths,
32     color_comp_lengths,
33     first_nonwhite_index,
34 )
```

Листинг 20 – Функция классификации строки (часть 1)

```

1 def classify_line(feats: LineFeatures):
2     cond_background = (
3         feat.first_nonwhite_index is None
4     )
5     if cond_background:
6         return State.BACKGROUND
7
8     min_long_black_line_length =
9         get_min_long_black_line_length(feats)
10    has_single_gray_comp = (
11        len(feats.comp_lengths) == 1 and
12        len(feats.gap_lengths) == 0 and
13        len(feats.color_comp_lengths) == 0
14    )
15    pretty_long_gray_comp = (
16        feat.count_gray > min_long_black_line_length
17    )
18    cond_long_black_line = (
19        has_single_gray_comp and
20        pretty_long_gray_comp
21    )
22    if cond_long_black_line:
23        return State.LONG_BLACK_LINE
24
25    min_medium_black_line_length =
26        get_min_medium_black_line_length(feats)
27    has_medium_sized_gray_comp = (
28        any(i > min_medium_black_line_length
29            for i in feats.gray_comp_lengths)
30    )
31    cond_medium_black_line = (
32        has_medium_sized_gray_comp
33    )
34    if cond_medium_black_line:
35        return State.MEDIUM_BLACK_LINE
36
37    n = MIN_NUMBER_OF_COMPONENTS_TO_BE_CONSIDERED_A_LOT
38    has_a_f_lot_of_comps = (
39        len(feats.comp_lengths) >
40        MIN_NUMBER_OF_COMPONENTS_TO_BE_CONSIDERED_A_F_LOT

```

Листинг 21 – Функция классификации строки (часть 2)

```

1      )
2      has_a_lot_of_comps_and_no_color = (
3          len(feat.comp_lengths) > n and
4          feat.count_color == 0
5      )
6      cond_many_text = (
7          has_a_f_lot_of_comps or
8          has_a_lot_of_comps_and_no_color
9      )
10     if cond_many_text:
11         return State.MANY_TEXT
12
13     has_color = (
14         feat.count_color > 0
15     )
16     cond_color = (
17         has_color
18     )
19     if cond_color:
20         return State.COLOR
21
22     if len(feat.comp_lengths) == 0:
23         mean_comp = 0
24     else:
25         mean_comp = np.mean(feat.comp_lengths)
26
27     if len(feat.gap_lengths) == 0:
28         mean_gap = 0
29         std_gap = 0
30         z_scores = []
31     else:
32         mean_gap = np.mean(feat.gap_lengths)
33         std_gap = np.std(feat.gap_lengths)
34         z_scores = (np.array(feat.gap_lengths) - mean_gap) /
35                     std_gap
36
37     has_huge_gaps = any(abs(z) > HUGE_GAP_ZSCORE for z in
38                         z_scores)
39     has_a_few_comps = (
40         len(feat.comp_lengths) <= n

```


Листинг 22 – Функция классификации строки (часть 3)

```

1      )
2      comps_are_mostly_small = (
3          mean_comp < MAX_COMPONENT_LENGTH_TO_BE_CONSIDERED_SMALL
4      )
5      gaps_are_mostly_small = (
6          mean_gap < MAX_COMPONENT_LENGTH_TO_BE_CONSIDERED_SMALL
7      )
8      cond_few_text = (
9          has_a_few_comps and
10         comps_are_mostly_small and
11         gaps_are_mostly_small and
12         not has_huge_gaps
13     )
14     if cond_few_text:
15         return State.FEW_TEXT
16
17     return State.UNDEFINED

```

Листинг 23 – Функция для создания первичной или уточненной разметки (часть 1)

```

1 def segment_document(
2     image: np.ndarray,
3     line_feature_func: Callable[[np.ndarray], LineFeatures],
4     raw: bool = False,
5 ):
6     empty_line = np.zeros_like(image[0:1]).reshape(-1,
7         image[0:1].shape[-1]).min(axis=-1).astype(int)
8     def empty_segment_data():
9         return SegmentData(
10             0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
11             empty_line.copy(), empty_line.copy()
12         )
13     def reset_segment_data(sd: SegmentData):
14         sd.start = sd.end
15         sd.count_single_long_black_line = 0
16         sd.count_single_medium_black_line = 0
17         sd.count_long_black_line = 0
18         sd.count_medium_black_line = 0

```

Листинг 24 – Функция для создания первичной или уточненной разметки
(часть 2)

```
1         sd.count_total_medium_black_line = 0
2         sd.count_many_text = 0
3         sd.count_color = 0
4         sd.count_few_text = 0
5         sd.count_undefined = 0
6         sd.count_white_px = 0
7         sd.count_color_px = 0
8         sd.count_gray_px = 0
9         sd.heatmap_black = np.zeros_like(empty_line)
10        sd.heatmap_color = np.zeros_like(empty_line)
11
12        results = []
13        height = image.shape[0]
14        prev_state = State.BACKGROUND
15        prev_feat = line_feature_func(image[0:1])
16        sd = empty_segment_data()
17        for y in range(1, height):
18            line = image[y:y+1]
19            feat = line_feature_func(line)
20            state = update_state(prev_state, feat)
21
22            bg_started = state == State.BACKGROUND and prev_state !=
                State.BACKGROUND
23            bg_finished = state != State.BACKGROUND and prev_state ==
                State.BACKGROUND
24            if bg_started or bg_finished:
25                class_name = classify_segment(prev_state, sd, raw)
26                result = (sd.start, sd.end, class_name)
27                results.append(result)
28                reset_segment_data(sd)
29
30            update_segment_data(sd, prev_feat, feat, line)
31            prev_state = state
32            prev_feat = feat
33        class_name = classify_segment(prev_state, sd, raw)
34        result = (sd.start, sd.end, class_name)
35        results.append(result)
36        return results
```

Листинг 25 – Конечный автомат (часть 1)

```

1 FSM = {
2     State.BACKGROUND: {
3         State.UNDEFINED:      State.UNDEFINED,
4         State.MANY_TEXT:      State.MANY_TEXT,
5         State.FEW_TEXT:       State.FEW_TEXT,
6         State.LONG_BLACK_LINE: State.LONG_BLACK_LINE,
7         State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
8         State.COLOR:          State.COLOR,
9         State.BACKGROUND:     State.BACKGROUND,
10    },
11    State.UNDEFINED: {
12        State.UNDEFINED:      State.UNDEFINED,
13        State.MANY_TEXT:      State.MANY_TEXT,
14        State.FEW_TEXT:       State.UNDEFINED,
15        State.LONG_BLACK_LINE: State.UNDEFINED,
16        State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
17        State.COLOR:          State.COLOR,
18        State.BACKGROUND:     State.BACKGROUND,
19    },
20    State.MANY_TEXT: {
21        State.UNDEFINED:      State.MANY_TEXT,
22        State.MANY_TEXT:      State.MANY_TEXT,
23        State.FEW_TEXT:       State.MANY_TEXT,
24        State.LONG_BLACK_LINE: State.MANY_TEXT,
25        State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
26        State.COLOR:          State.MANY_TEXT,
27        State.BACKGROUND:     State.BACKGROUND,
28    },
29    State.FEW_TEXT: {
30        State.UNDEFINED:      State.UNDEFINED,
31        State.MANY_TEXT:      State.MANY_TEXT,
32        State.FEW_TEXT:       State.FEW_TEXT,
33        State.LONG_BLACK_LINE: State.LONG_BLACK_LINE,
34        State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
35        State.COLOR:          State.COLOR,
36        State.BACKGROUND:     State.BACKGROUND,
37    },

```

Листинг 26 – Конечный автомат (часть 2)

```

1  State.LONG_BLACK_LINE: {
2      State.UNDEFINED:      State.LONG_BLACK_LINE,
3      State.MANY_TEXT:      State.LONG_BLACK_LINE,
4      State.FEW_TEXT:       State.LONG_BLACK_LINE,
5      State.LONG_BLACK_LINE: State.LONG_BLACK_LINE,
6      State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
7      State.COLOR:          State.LONG_BLACK_LINE,
8      State.BACKGROUND:     State.BACKGROUND,
9  },
10 State.MEDIUM_BLACK_LINE: {
11     State.UNDEFINED:      State.MEDIUM_BLACK_LINE,
12     State.MANY_TEXT:      State.MEDIUM_BLACK_LINE,
13     State.FEW_TEXT:       State.MEDIUM_BLACK_LINE,
14     State.LONG_BLACK_LINE: State.LONG_BLACK_LINE,
15     State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
16     State.COLOR:          State.COLOR,
17     State.BACKGROUND:     State.BACKGROUND,
18 },
19 State.COLOR: {
20     State.UNDEFINED:      State.COLOR,
21     State.MANY_TEXT:      State.COLOR,
22     State.FEW_TEXT:       State.COLOR,
23     State.LONG_BLACK_LINE: State.LONG_BLACK_LINE,
24     State.MEDIUM_BLACK_LINE: State.MEDIUM_BLACK_LINE,
25     State.COLOR:          State.COLOR,
26     State.BACKGROUND:     State.BACKGROUND,
27 },
28 }
```

Листинг 27 – Структура State

```

1 class State:
2     UNDEFINED      = 0
3     BACKGROUND     = 1
4     FEW_TEXT       = 2
5     MANY_TEXT      = 3
6     COLOR          = 4
7     MEDIUM_BLACK_LINE = 5
8     LONG_BLACK_LINE = 6
```

Листинг 28 – Функция `handle_undefined` (часть 1)

```

1 def handle_undefined(sd: SegmentData):
2     def text(sd: SegmentData):
3         height = sd.end - sd.start
4         not_very_high = height <
5             MAX_SEGMENT_HEIGHT_TO_BE_CONSIDERED_NOT_VERY_HIGH
6         had_a_lot_of_few_text = (sd.count_few_text / height) >
7             MIN_FEW_TEXT_RATIO_TO_BE_CONSIDERED_A_LOT
8         return (
9             not_very_high or
10            had_a_lot_of_few_text
11        )
12
13     def code(sd: SegmentData):
14         height = sd.end - sd.start
15
16         high_vbls = sd.heatmap_black == height
17         padded = np.concatenate(([False], high_vbls, [False]))
18         diff = np.diff(padded.astype(int))
19         n_vertical_black_lines = len(np.where(diff == 1)[0])
20
21         return n_vertical_black_lines == 2
22
23     def figure(sd: SegmentData):
24         height = sd.end - sd.start
25         pretty_high = height >
26             MIN_FIGURE_HEIGHT_TO_BE_CONSIDERED_HIGH
27         return pretty_high
28
29     def plot(sd: SegmentData):
30         height = sd.end - sd.start
31
32         high_vbls = sd.heatmap_black >=
33             PLOT_VERTICAL_LINE_HEIGHT_CORRECTION * height
34         padded = np.concatenate(([False], high_vbls, [False]))
35         diff = np.diff(padded.astype(int))
36         n_vertical_black_lines = len(np.where(diff == 1)[0])
37
38         return n_vertical_black_lines == 1

```

Листинг 29 – Функция `handle_undefined` (часть 2)

```
1     if text(sd):
2         return ClassNames[Class.TEXT]
3
4     if code(sd):
5         return ClassNames[Class.CODE]
6
7     if figure(sd):
8         return ClassNames[Class.FIGURE]
9
10    if plot(sd):
11        return ClassNames[Class.PLOT]
12
13    return ClassNames[Class.UNDEFINED]
```

Листинг 30 – Функция `handle_background`

```
1 def handle_background(_: SegmentData):
2     return ClassNames[Class.BACKGROUND]
```

Листинг 31 – Функция `handle_few_text`

```
1 def handle_few_text(_: SegmentData):
2     return ClassNames[Class.TEXT]
```

Листинг 32 – Функция `handle_many_text` (часть 1)

```
1 def handle_many_text(sd: SegmentData):
2     def table(sd: SegmentData):
3         height = sd.end - sd.start
4
5         high_vbls = sd.heatmap_black == height
6         padded = np.concatenate(([False], high_vbls, [False]))
7         diff = np.diff(padded.astype(int))
8         lbl_start_indices = np.where(diff == 1)[0]
9         n_vertical_black_lines = len(lbl_start_indices)
10        has_more_than_two_vertical_lines = n_vertical_black_lines
11        > 2
12
13        min_space_is_reasonably_small = True
```

Листинг 33 – Функция `handle_many_text` (часть 2)

```

1      if has_more_than_two_vertical_lines:
2          min_space_is_reasonably_small =
3              min(np.diff(lbl_start_indices)) >
4              MIN_REASONABLY_SMALL_SPACE_BETWEEN_TWO_COLUMNS_IN_TABLE
5
6      pretty_high = height >
7          MIN_SEGMENT_HEIGHT_TO_BE_CONSIDERED_HIGH
8
9      return (
10         pretty_high and
11         has_more_than_two_vertical_lines and
12         min_space_is_reasonably_small
13     )
14
15 def code(sd: SegmentData):
16     height = sd.end - sd.start
17
18     high_vbls = sd.heatmap_black == height
19     padded = np.concatenate(([False], high_vbls, [False]))
20     diff = np.diff(padded.astype(int))
21     n_vertical_black_lines = len(np.where(diff == 1)[0])
22
23     has_two_vertical_lines = n_vertical_black_lines == 2
24     has_no_mbls = sd.count_single_medium_black_line == 0
25     had_many_text = sd.count_many_text > 0
26     has_no_color = sd.count_color == 0
27
28     pretty_high = height >
29         MIN_SEGMENT_HEIGHT_TO_BE_CONSIDERED_HIGH
30
31     return (
32         pretty_high and
33         has_two_vertical_lines and
34         has_no_mbls and
35         (
36             had_many_text or
37             has_no_color
38         )
39     )

```

Листинг 34 – Функция `handle_many_text` (часть 3)

```

1  if table(sd):
2      return ClassNames[Class.TABLE]
3
4  if code(sd):
5      return ClassNames[Class.CODE]
6
7  return ClassNames[Class.TEXT]
```

Листинг 35 – Функция `handle_color`

```

1  def handle_color(sd: SegmentData):
2      def plot(sd: SegmentData):
3          height = sd.end - sd.start
4          high_vbls = sd.heatmap_black >=
5              PLOT_VERTICAL_LINE_HEIGHT_CORRECTION * height
6          padded = np.concatenate(([False], high_vbls, [False]))
7          diff = np.diff(padded.astype(int))
8          n_vertical_black_lines = len(np.where(diff == 1)[0])
9          has_single_vertical_axis = n_vertical_black_lines == 1
10         has_pretty_small_color_to_white_relation =
11             (sd.count_color_px /
12              sd.count_white_px) <
13             MIN_COLOR_TO_WHITE_RATIO_TO_BE_CONSIDERED_SMALL
14         return (
15             has_single_vertical_axis and
16             has_pretty_small_color_to_white_relation
17         )
18
19     def undefined(sd: SegmentData):
20         height = sd.end - sd.start
21         smol = height < MAX_SEGMENT_HEIGHT_TO_BE_CONSIDERED_SMALL
22         return smol
23
24     if plot(sd):
25         return ClassNames[Class.PLOT]
26
27     if undefined(sd):
28         return ClassNames[Class.UNDEFINED]
29
30     return ClassNames[Class.FIGURE]
```


Листинг 36 – Функция `handle_medium_black_line` (часть 1)

```

1 def handle_medium_black_line(sd: SegmentData):
2     def text(sd: SegmentData):
3         height = sd.end - sd.start
4         had_a_lot_of_a_lot_of_text = (sd.count_many_text /
5             height) > MIN_MANY_TEXT_RATIO_TO_BE_CONSIDERED_A_LOT
6         pretty_small = height <
7             MAX_SEGMENT_HEIGHT_FOR_A_TEXT_TO_BE_CONSIDERED_NOT_SMALL
8         had_a_lot_of_undefined = (sd.count_undefined / height) >
9             MIN_UNDEFINED_RATIO_TO_BE_CONSIDERED_A_LOT
10        had_a_lot_of_few_text = (sd.count_few_text / height) >
11            MIN_FEW_TEXT_RATIO_TO_BE_CONSIDERED_A_LOT
12        return (
13            had_a_lot_of_a_lot_of_text or
14            (
15                pretty_small and
16                (
17                    had_a_lot_of_undefined or
18                    had_a_lot_of_few_text
19                )
20            )
21        )
22
23    def figure(sd: SegmentData):
24        height = sd.end - sd.start
25        has_color = sd.count_color > 0
26        has_a_lot_of_mbls = (sd.count_medium_black_line / height)
27            > MAX_MBL_RATIO_TO_BE_CONSIDERED_FEW
28        pretty_high = height >
29            MIN_SEGMENT_HEIGHT_TO_BE_CONSIDERED_HIGH
30        return (
31            pretty_high and
32            (
33                has_color or
34                has_a_lot_of_mbls
35            )
36        )

```

Листинг 37 – Функция `handle_medium_black_line` (часть 2)

```

1  def plot(sd: SegmentData):
2      height = sd.end - sd.start
3      has_color = sd.count_color > 0
4      has_a_few_mbls = (sd.count_medium_black_line / height) <
        MAX_MBL_RATIO_TO_BE_CONSIDERED_FEW
5
6      high_vbls = sd.heatmap_black >=
        PLOT_VERTICAL_LINE_HEIGHT_CORRECTION * height
7      padded = np.concatenate(([False], high_vbls, [False]))
8      diff = np.diff(padded.astype(int))
9      n_vertical_black_lines = len(np.where(diff == 1)[0])
10
11     all_px = sd.count_white_px + sd.count_color_px +
        sd.count_gray_px
12     has_many_white_pixels = sd.count_white_px / all_px >
        MIN_WHITE_PIXELS_RATIO_TO_BE_CONSIDERED_MANY
13
14     return (
15         has_color and
16         has_a_few_mbls and
17         n_vertical_black_lines >= 2 and
18         has_many_white_pixels
19     )
20
21  def equation(sd: SegmentData):
22      height = sd.end - sd.start
23      not_very_high = height <
        MIN_SEGMENT_HEIGHT_TO_BE_CONSIDERED_HIGH
24      has_a_few_mbls = sd.count_single_medium_black_line <
        MAX_MBL_RATIO_TO_BE_CONSIDERED_FEW
25      has_single_mbl = sd.count_single_medium_black_line == 1
26      return (
27          (
28              not_very_high and
29              has_a_few_mbls
30          ) or
31          has_single_mbl
32      )

```

Листинг 38 – Функция `handle_medium_black_line` (часть 3)

```
1 def diagram(sd: SegmentData):
2     return sd.count_single_medium_black_line > 1
3
4 def undefined(sd: SegmentData):
5     height = sd.end - sd.start
6     smol = height < MAX_SEGMENT_HEIGHT_TO_BE_CONSIDERED_SMALL
7     return smol
8
9 if plot(sd):
10     return ClassNames[Class.PLOT]
11
12 if figure(sd):
13     return ClassNames[Class.FIGURE]
14
15 if diagram(sd):
16     return ClassNames[Class.DIAGRAM]
17
18 if text(sd):
19     return ClassNames[Class.TEXT]
20
21 if equation(sd):
22     return ClassNames[Class.UNDEFINED]
23
24 if undefined(sd):
25     return ClassNames[Class.UNDEFINED]
26
27 return ClassNames[Class.DIAGRAM]
```

Листинг 39 – Функция `handle_long_black_line` (часть 1)

```

1 def handle_long_black_line(sd: SegmentData):
2     def table(sd: SegmentData):
3         height = sd.end - sd.start
4
5         high_vbls = sd.heatmap_black == height
6         padded = np.concatenate(([False], high_vbls, [False]))
7         diff = np.diff(padded.astype(int))
8         lbl_start_indices = np.where(diff == 1)[0]
9         n_vertical_black_lines = len(lbl_start_indices)
10        has_more_than_two_vertical_lines = n_vertical_black_lines
11        > 2
12
13        min_space_is_reasonably_small = True
14        if has_more_than_two_vertical_lines:
15            min_space_is_reasonably_small =
16                min(np.diff(lbl_start_indices)) >
17                MIN_REASONABLY_SMALL_SPACE_BETWEEN_TWO_COLUMNS_IN_TABLE
18
19        return (
20            has_more_than_two_vertical_lines and
21            min_space_is_reasonably_small
22        )
23
24    def code(sd: SegmentData):
25        height = sd.end - sd.start
26
27        high_vbls = sd.heatmap_black == height
28        padded = np.concatenate(([False], high_vbls, [False]))
29        diff = np.diff(padded.astype(int))
30        n_vertical_black_lines = len(np.where(diff == 1)[0])
31
32        has_two_vertical_lines = n_vertical_black_lines == 2
33        has_no_mbls = sd.count_single_medium_black_line == 0
34        had_many_text = sd.count_many_text > 0
35        has_no_color = sd.count_color == 0

```

Листинг 40 – Функция `handle_long_black_line` (часть 2)

```

1      return (
2          has_two_vertical_lines and
3          has_no_mbls and
4          (
5              had_many_text or
6              has_no_color
7          )
8      )
9
10     def diagram(sd: SegmentData):
11         has_no_color = sd.count_color == 0
12         has_a_few_mbls = sd.count_single_medium_black_line >= 2
13         return (
14             has_no_color and
15             has_a_few_mbls
16         )
17
18     def plot(sd: SegmentData):
19         height = sd.end - sd.start
20         has_color = sd.count_color > 0
21         has_a_few_lbls = (sd.count_long_black_line / height) < 0.1
22
23         high_vbls = sd.heatmap_black >=
24             PLOT_VERTICAL_LINE_HEIGHT_CORRECTION * height
25         padded = np.concatenate(([False], high_vbls, [False]))
26         diff = np.diff(padded.astype(int))
27         n_vertical_black_lines = len(np.where(diff == 1)[0])
28
29         return (
30             has_color and
31             has_a_few_lbls and
32             n_vertical_black_lines >= 2
33         )
34
35     def undefined(sd: SegmentData):
36         height = sd.end - sd.start
37         smol = height < MAX_SEGMENT_HEIGHT_TO_BE_CONSIDERED_SMALL
38         return smol

```

Листинг 41 – Функция `handle_long_black_line` (часть 3)

```
1  if undefined(sd):
2      return ClassNames[Class.UNDEFINED]
3
4  if plot(sd):
5      return ClassNames[Class.PLOT]
6
7  if table(sd):
8      return ClassNames[Class.TABLE]
9
10 if code(sd):
11     return ClassNames[Class.CODE]
12
13 if diagram(sd):
14     return ClassNames[Class.DIAGRAM]
15
16 return ClassNames[Class.FIGURE]
```

Листинг 42 – Функция обновления информации о сегменте

```

1 def update_segment_data(sd: SegmentData, prev_feat, feat:
  LineFeatures, line: np.ndarray):
2     prev_state = classify_line(prev_feat)
3     state = classify_line(feat)
4     sd.end += 1
5     if prev_state != state and state == State.LONG_BLACK_LINE:
6         sd.count_single_long_black_line += 1
7     elif state == State.LONG_BLACK_LINE:
8         sd.count_long_black_line += 1
9     elif prev_state != state and state == State.MEDIUM_BLACK_LINE:
10        sd.count_single_medium_black_line += 1
11    elif state == State.MEDIUM_BLACK_LINE:
12        sd.count_medium_black_line += 1
13    elif state == State.MANY_TEXT:
14        sd.count_many_text += 1
15    elif state == State.COLOR:
16        sd.count_color += 1
17    elif state == State.FEW_TEXT:
18        sd.count_few_text += 1
19    elif state == State.UNDEFINED:
20        sd.count_undefined += 1
21    elif state == State.BACKGROUND:
22        return
23    sd.count_white_px += feat.count_white
24    sd.count_color_px += feat.count_color
25    sd.count_gray_px += feat.count_gray
26
27    min_medium_black_line_length =
        get_min_medium_black_line_length(feat)
28    # NOTE: Account for several MBLs on single line
29    count_total_medium_black_lines =
        sum(np.array(feat.gray_comp_lengths) >
30
31    sd.count_total_medium_black_line +=
        count_total_medium_black_lines
32    (_, _, mask_color, mask_gray) = get_masks(line, WHITE_THRESH,
        GRAY_TOL)
33    sd.heatmap_black += mask_gray.astype(int)
34    sd.heatmap_color += mask_color.astype(int)

```

Листинг 43 – Функция создания объединенной разметки (часть 1)

```

1 def merge(markup: List[Tuple[int, int, str]]):
2     tmp_markup = [markup[0]]
3     new_markup = [markup[0]]
4     prev_start = curr_start = next_start = prev_end = curr_end =
        next_end = 0
5     prev_height = curr_height = next_height = 0
6     prev_class = curr_class = next_class = "None"
7
8     # Merge small Background segment with nearest larger segment
9     for i in range(1, len(markup) - 1):
10         (prev_start, prev_end, prev_class) = markup[i - 1]
11         (curr_start, curr_end, curr_class) = markup[i]
12         (next_start, next_end, next_class) = markup[i + 1]
13
14         prev_height = prev_end - prev_start
15         curr_height = curr_end - curr_start
16         next_height = next_end - next_start
17
18         if curr_class == ClassNames[Class.BACKGROUND] and
            curr_height <
19             MAX_SEGMENT_HEIGHT_TO_BE_MERGED_WITH_NEAREST_LARGER_SEGMENT:
20             if prev_height > curr_height:
21                 curr_class = prev_class
22             elif next_height > curr_height:
23                 curr_class = next_class
24
25         segment = (curr_start, curr_end, curr_class)
26         tmp_markup.append(segment)
27
28         segment = (next_start, next_end, next_class)
29         tmp_markup.append(segment)
30
31     new_markup = merge_segments(tmp_markup)
32     tmp_markup = [new_markup[0]]
33
34     # Remove background between same-class segments
35     for i in range(1, len(new_markup) - 1):
36         (prev_start, prev_end, prev_class) = new_markup[i - 1]
37         (curr_start, curr_end, curr_class) = new_markup[i]
38         (next_start, next_end, next_class) = new_markup[i + 1]

```


Листинг 44 – Функция создания объединенной разметки (часть 2)

```

1      prev_height = prev_end - prev_start
2      curr_height = curr_end - curr_start
3      next_height = next_end - next_start
4
5      if (
6          prev_class == next_class and
7          curr_class == ClassNames[Class.BACKGROUND]
8      ):
9          curr_class = next_class
10
11     segment = (curr_start, curr_end, curr_class)
12     tmp_markup.append(segment)
13
14     segment = (next_start, next_end, next_class)
15     tmp_markup.append(segment)
16
17     new_markup = merge_segments(tmp_markup)
18     tmp_markup = [new_markup[0]]
19
20     # Make small backgrounds uncertain
21     for i in range(len(new_markup)):
22         (curr_start, curr_end, curr_class) = new_markup[i]
23         curr_height = curr_end - curr_start
24
25         if curr_class == ClassNames[Class.BACKGROUND] and
26            curr_height <
27            MAX_BACKGROUND_HEIGHT_TO_BECOME_UNDEFINED:
28             curr_class = ClassNames[Class.UNDEFINED]
29
30         segment = (curr_start, curr_end, curr_class)
31         tmp_markup.append(segment)
32
33     new_markup = merge_segments(tmp_markup)
34     tmp_markup = [new_markup[0]]
35
36     # Merge uncertainty with greatest neighbor
37     for i in range(1, len(new_markup) - 1):
38         (prev_start, prev_end, prev_class) = new_markup[i - 1]
39         (curr_start, curr_end, curr_class) = new_markup[i]
40         (next_start, next_end, next_class) = new_markup[i + 1]

```

Листинг 45 – Функция создания объединенной разметки (часть 3)

```

1      prev_height = prev_end - prev_start
2      curr_height = curr_end - curr_start
3      next_height = next_end - next_start
4
5      if curr_class == ClassNames[Class.UNDEFINED] and
        curr_height < MAX_UNDEFINED_HEIGHT_TO_BE_MERGED:
6          if prev_class != ClassNames[Class.BACKGROUND] and
            prev_height > curr_height:
7              curr_class = prev_class
8          elif next_class != ClassNames[Class.BACKGROUND] and
            next_height > curr_height:
9              curr_class = next_class
10
11         segment = (curr_start, curr_end, curr_class)
12         tmp_markup.append(segment)
13
14     if next_class == ClassNames[Class.UNDEFINED] and next_height
        < MAX_UNDEFINED_HEIGHT_TO_BE_MERGED:
15         next_class = curr_class
16
17     segment = (next_start, next_end, next_class)
18     tmp_markup.append(segment)
19
20     tmp_markup = tmp_markup[1:]
21     second_class = tmp_markup[1][2]
22     first_start = tmp_markup[0][0]
23     first_end = tmp_markup[0][1]
24     tmp_markup[0] = [first_start, first_end, second_class]
25     if tmp_markup[-1][2] == ClassNames[Class.BACKGROUND]:
26         s = tmp_markup[-1][0]
27         e = tmp_markup[-1][1]
28         c = tmp_markup[-2][2]
29         tmp_markup[-1] = (s, e, c)
30
31     new_markup = merge_segments(tmp_markup)
32
33     return new_markup

```

Листинг 46 – Скрипт, используемый в исследовании (часть 1)

```

1 from PIL import Image
2 import numpy as np
3 import fitz
4 from concurrent.futures import ProcessPoolExecutor
5 from memory_profiler import memory_usage
6
7 from fast import segdoc as sd
8 import time
9
10 def page_to_image(pdf_path: str, page_index: int, scale_factor:
    float = 3) -> Image.Image:
11     doc = fitz.open(pdf_path)
12     page = doc.load_page(page_index)
13     pix = page.get_displaylist()
14         .get_pixmap(matrix=fitz.Matrix(scale_factor,
15                                         scale_factor))
16     img = Image.frombytes("RGB", (pix.w, pix.h), pix.samples)
17     doc.close()
18     return img
19
20 def process_pages(pdf_path, page_range, markup_type):
21     """Функция для обработки поддиапазона страниц."""
22     partial_results = []
23     for i in page_range:
24         image = page_to_image(pdf_path, i)
25         image_np = np.array(image)
26         markup = sd(image_np, markup_type)
27         assert markup is not None
28         del image, image_np
29         partial_results.append({
30             "page": i + 1,
31             "segments": [
32                 {"y_start": s[0], "y_end": s[1], "label": s[2]}
33                 for s in markup
34             ]
35         })
36     return partial_results

```

Листинг 47 – Скрипт, используемый в исследовании (часть 2)

```
1 def chunk_indices(total, chunks):
2     """Делит_диапазон_страниц_на_N_примерно_равных_кусков."""
3     avg = total // chunks
4     remainder = total % chunks
5     indices = []
6     start = 0
7     for i in range(chunks):
8         end = start + avg + (1 if i < remainder else 0)
9         indices.append(list(range(start, end)))
10        start = end
11    return indices
12
13 def main(pdf_path, markup_type, num_workers=8, page_indices=None):
14     results = []
15
16     if page_indices is None:
17         doc = fitz.open(pdf_path)
18         total_pages = doc.page_count
19         page_indices = list(range(total_pages))
20         doc.close()
21
22     chunks = chunk_indices(len(page_indices), num_workers)
23     pages_chunks = [ [page_indices[i] for i in chunk] for chunk
24                      in chunks ]
25
26     with ProcessPoolExecutor(max_workers=num_workers) as executor:
27         futures = [executor.submit(process_pages, pdf_path,
28                                   chunk, markup_type) for chunk in pages_chunks]
29         for future in futures:
30             results.extend(future.result())
31
32     results.sort(key=lambda x: x["page"])
33
34     return results
```

Листинг 48 – Скрипт, используемый в исследовании (часть 3)

```

1 def parse_page_ranges(pages_str, max_page):
2     """Преобразует строку диапазонов страниц в список индексов_
        (0-based). """
3     pages = set()
4     for part in pages_str.split(','):
5         if '-' in part:
6             start, end = map(int, part.split('-'))
7             pages.update(range(start - 1, min(end, max_page))) #
                -1 т.к. 0-based
8         else:
9             page = int(part)
10            if 1 <= page <= max_page:
11                pages.add(page - 1)
12    return sorted(pages)
13
14 def benchmarked_main(pdf_path, markup_type, num_workers,
    page_indices):
15     def runner():
16         return main(pdf_path, markup_type, num_workers,
            page_indices)
17
18     start_time = time.time()
19
20     mem_usage, result = memory_usage(
21         (runner, (), {}),
22         interval=0.1,
23         timeout=None,
24         include_children=True,
25         max_usage=True,
26         retval=True
27     )
28
29     end_time = time.time()
30
31     print(f"Elapsed_time: {end_time - start_time:.2f} seconds")
32     print(f"Max_memory_usage: {mem_usage:.2f} MiB")
33
34     return result

```

Листинг 49 – Скрипт, используемый в исследовании (часть 4)

```

1 if __name__ == "__main__":
2     import argparse
3
4     parser = argparse.ArgumentParser(description="Segment_PDF_
5         pages_and_export_markup_to_JSON.")
6     parser.add_argument("pdf_path", type=str, help="Path_to_the_
7         input_PDF_file")
8     parser.add_argument(
9         "markup_type",
10        type=int,
11        choices=[0, 1, 2, 3],
12        help="Markup_type:_0_-raw,_1_-primary,_2_-specified,_3_-merged"
13    )
14    parser.add_argument(
15        "-w", "--workers",
16        type=int,
17        default=8,
18        help="Number_of_worker_processes_(default:_8)"
19    )
20    parser.add_argument(
21        "-p", "--pages",
22        type=str,
23        default=None,
24        help="Page_ranges_to_process,_e.g.,_'1-3,5,7-9'"
25    )
26
27    args = parser.parse_args()
28
29    # Открываем PDF, чтобы узнать число страниц
30    doc = fitz.open(args.pdf_path)
31    total_pages = doc.page_count
32    doc.close()
33
34    if args.pages:
35        pages_to_process = parse_page_ranges(args.pages,
36            total_pages)
37    else:
38        pages_to_process = list(range(total_pages))

```

Листинг 50 – Скрипт, используемый в исследовании (часть 5)

```
1 markup = benchmarked_main(  
2     pdf_path=args.pdf_path,  
3     markup_type=args.markup_type,  
4     num_workers=args.workers,  
5     page_indices=pages_to_process  
6 )
```