

## ДЕТЕКТИРОВАНИЕ ТЕКСТОВЫХ ОБЛАСТЕЙ НА ИЗОБРАЖЕНИИ ДОКУМЕНТА МЕТОДОМ СЛИЯНИЯ

*В данной статье представлен метод автоматического детектирования текстовых областей на изображении документа с помощью слияния областей. Предложенный метод комбинирует простоту реализации алгоритма слияния с анализом и отсечением графических областей, которые не несут текстовой информации. Проведено моделирование работы предложенного метода, выделены классы практических задач и условия наиболее эффективного его применения.*

**Ключевые слова:** текстовая область, графическая область, слияние, изображение документа, пороговое значение, преобразование Хафа.

### Введение

Задача распознавания документов в современном мире постоянного конвертирования информации, хранящейся в традиционной форме, в соответствующий электронный вариант приобрела широкое распространение. При распознавании документа логичным этапом видится создание общей структуры документа (page layout analysis), распознавание текстовых элементов (OCR – optical character recognition), формирование таблиц, перенос изображений и формул и т.п.

Основную информационную составляющую при конвертировании информации в электронную форму образует текстовое наполнение. Реализацией задачи конвертирования текста на изображении в соответствующий электронный текст занимается OCR, в то же время, существующие методы решения данной проблемы ориентируются на распознавание именно текста, а не графической информации внутри документа – такой как схемы, диаграммы, чертежи. Распознавание же этих элементов с помощью существующих механизмов OCR чаще всего приводит к генерации лишнего текста, символов, которых в реальности в документе не существует.

Данная статья посвящена физическому анализу [1] структуры документа, т.е. такому анализу, который позволяет выделить такие физические элементы документа, как текст, графические объекты, таблицы. Существующие методы сегментации часто базируются на анализе проекций [2, 11], связанных компонент [3], фильтрации изображения, анализе границ и частот смены цвета [3, 11], слиянии-расщеплении элементов изображения и т.п.

Как логический, так и физический анализ документа может быть выполнен одним из двух

наиболее распространенных подходов [4–6, 12]: «снизу вверх» («bottom-up» – анализ мелких элементов изображения с последующим их слиянием, образуя при этом элементы более высокого логического уровня: блоки, параграфы, абзацы и т.п.) либо «сверху вниз» («top-down» – разбиение изображения на абзацы, блоки с их последующим анализом и разбиением) [4, 5, 7, 8]. Предложенный подход подразумевает использование стратегии «снизу-вверх», используя слияние элементов нижнего уровня.

Статья посвящена разработке метода, который комбинирует простоту реализации стратегии слияния с эффективным детектированием текстовых областей, текстовой информации внутри таблиц и графических областей, и является адаптируемым к структуре изображения документа. Большинство нюансов реализации данного метода являются универсальными и позволяют обрабатывать изображения документов различного типа.

### Предварительная обработка

Современные методы получения изображения документа, такие как сканирование или фотографирование, позволяют получить цветное или полутоновое изображение различного уровня четкости и качества. С практической точки зрения, целесообразно избегать обработки полноцветных и даже полутоновых изображений при решении таких проблем, в которых возможно пренебрежение подобной информацией. При обработке документа мы в дальнейшем будем оперировать понятиями «значимых» объектов (текста, таблиц, графических объектов) и «фона», разделить же сцену на изображении на фон и объект позволяет бинаризация. Следует отметить, что к проблеме выбора порога для бинаризации необходимо

подходить весьма серьезно, когда речь идет о практических задачах, требующих высокоточной обработки, таких, как собственно OCR, сегментация с целью измерения биологических объектов и т.п. При решении рассматриваемой проблемы определения текста мы не уделяем проблеме выбора порога особого значения в предположении, что изображение документа является высококонтрастным в том понимании, что текст от фона значительно отличается по яркости. В качестве значения порога бинаризации было выбрано толерантное значение в половину максимальной яркости, т.е. в порог был установлен в  $\delta = 128$ .

Кроме бинаризации часто также используют эквализацию гистограммы для улучшения контраста, фильтры шумоподавления (размытие, медианный фильтр), морфологические операции для устранения шумов и т.п.

В данной статье мы говорим об изображении документа, и о разбиении его на области. Под областью мы будем понимать атомарный объект, контент которого является однородным, т.е. содержит либо текст, либо изображение, либо таблицу и т.п. Область, содержащую какой-либо графический объект внутри изображения документа (схему, чертеж, изображение и т.п.) мы будем называть графической областью, область с контентом в виде текста – текстовой областью.

### Выделение потенциальных областей интереса

Метод анализа изображения документа «снизу-вверх» предполагает слияние мелких элементов с образованием более крупных [9]. Первоначальным этапом является сканирование изображения с помощью сканирующего окна, которые дает возможность определить минимально возможные области, содержащие значимые объекты, для последующего слияния. В качестве таких областей на первом этапе обработки предлагается использовать классическую окрестность каждого пикселя, который потенциально может быть частью текстовой области. Чаще всего при выполнении подобных операций принято использовать сканирующее окно размером 3x3 пикселей, но в нашей задаче было использовано окно размера 5x5, поскольку большее окно позволяет ускорить процесс обработки.

Результатом обработки на этом этапе является множество областей, каждая из которых имеет площадь, равную сканирующему окну и расположена вокруг значимых пикселей, яркость которых удовлетворяет условию бинаризации (рис.1).

Переход на более высокий иерархический уровень выполняется путем слияния

пересекающихся квадратных областей с образованием произвольных прямоугольных регионов (рис.2).



Рис. 1 Результат сканирования изображения плавающим окном размера 5x5

*Статья посвящена разработке метод предложено использование эмпирических признаков в комбинации с применением проекционног расположения заголовка. Экспериментальные метода и позволили определить области его при*

Рис. 2 Результат слияния пересекающихся областей

Результатом подобного технического слияния без анализа контента являются области, которые могут содержать контент типа, не являющегося потенциально интересным, например, линии разметки документа. Подобные области могут быть идентифицированы по нестандартному отношению высоты области к ширине либо с использованием других шумоподавляющих фильтров. Например, фильтрация областей по отношению  $height_{r_1} / width_{r_1} > 0.02$ , где  $r_1$  – анализируемая область,  $height, width$  – ее высота и ширина соответственно, позволяет получить результат, представленный на рис. 3.

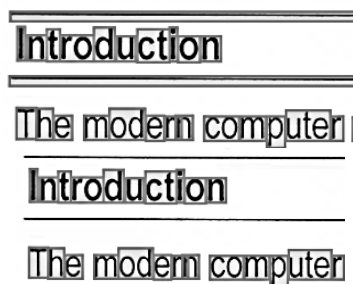


Рис. 3 Результат фильтрации потенциальных текстовых областей по отношению высоты к ширине: вверху – графическое отображение областей до фильтрации, внизу – графическое отображение областей после фильтрации.

Слияние квадратных пересекающихся областей позволяет сформировать совокупность прямоугольных областей, которые, в зависимости от выбранного ранее размера плавающего окна сканирования и особенностей шрифта могут содержать одну или несколько букв текста. Для получения области с контентом более высокого иерархического уровня (например, слова)

выполняется слияние тех областей, которые пересекаются в рамках допустимого радиуса  $\delta_1$ . Для этого каждая из соседних областей увеличивается по горизонтали и вертикали на значение радиуса  $\delta_1$ , после чего пересекающиеся области объединяются. Необходимость увеличения области на значение  $\delta_1$  по вертикали объясняется возможным соседством двух областей, которые содержат символы разного размера и расположены выше либо ниже базовой линии текста. Графическая интерпретация данного критерия приведена на рис. 4.



Рис.4 Графическая интерпретация критерия пересечения областей по радиусу  $\delta_1$

На рис.5 показан результат слияния набора текстовых областей с использованием различных пороговых значений.

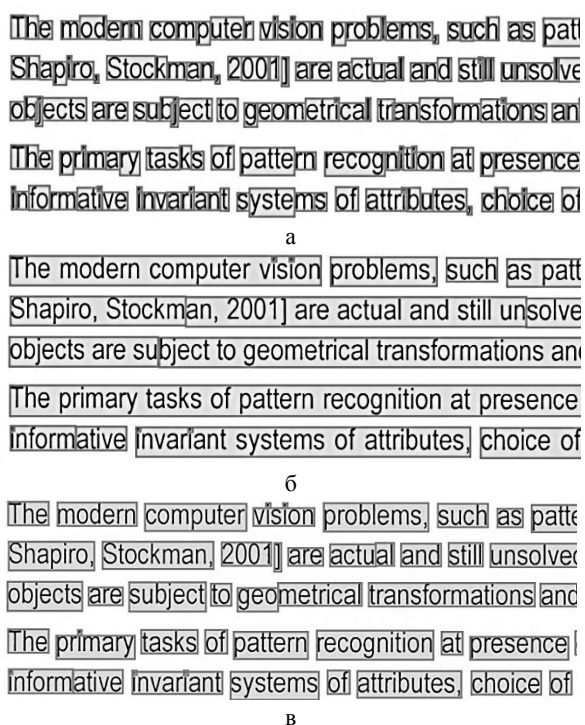


Рис.5 Результат слияния текстовых областей по горизонтали: а – графическое отображение областей перед слиянием, б – графическое отображение областей после слияния при  $\delta_1 = 15$ , в – графическое отображение областей после слияния при  $\delta_1 = 5$ .

Как видно, выбор значения  $\delta_1$  определяет строгость правила формирования единственной текстовой области в виде строки. При выборе значения порогового ограничения важную роль может иметь априорная информация о структуре

изображения документа, например,  $\delta_1$  непосредственно указывает, могут ли быть объединены в единственную область такие области с текстом, которые находятся в одной горизонтальной физической строке, но в разных столбцах документа.

## Выделение изображений

Разделение текстовых и графических областей при выполнении анализа документа является нетривиальной задачей, требующей индивидуального подхода к каждому типу изображения документа. Существует два основных класса сложности методов решения подобной проблемы, первый из которых оперирует количественными характеристиками области, такими, как плотность пикселей интереса, геометрические пропорции размеров [10, 15] и т.п., второй – использует более трудоемкие методы с большей точностью, которые опираются на глубокий анализ области и использование априорной информации о возможных различиях [13, 14].

В [15] предложено считать текстовой областью такую область, для которой выполнено условие:

$$0.045 < ratio_{r_1} < 0.444,$$

где  $ratio_{r_1} = (Total\ black\ pixels) / (Total\ pixels)$

для области  $r_1$ , при этом попадание за пределы минимальной границы соответствует пустой области, попадание за пределы максимальной – графической области.

Тем не менее, подобное условие может быть выполнено для большинства изображений документов лишь в тех случаях, если графический элемент внутри области  $r_1$  является достаточно контрастным и при этом его предварительная обработка (особенно бинаризация) сохраняет этот контраст.

В качестве обобщения предложенного метода может быть рассмотрено использование предварительной информации об области, содержащей изображение. Одной из самых легкодоступных в вычислительном плане характеристик является анализ физических размеров области, поскольку для множества изображений документов графические области значительно отличаются по размеру от текстовых областей.

Эмпирические исследования показали, что целесообразным может быть использование следующих условий для идентификации области  $r_1$ , содержащей изображение:

$$(ratio > 0.4 \vee ratio < 0.1) \wedge width_{r_1} < 2.5height_{r_1},$$

где  $width_{r_1}$ ,  $height_{r_1}$  – ширина и высота прямоугольной области соответственно. Таким

образом, к графической области на изображении документа предъявляются требования повышенной или пониженной плотности точек, а также несоответствие размерам типичной текстовой области, для которой ширина значительно больше высоты.

### Слияние областей по горизонтали

Получив множество текстовых областей после отсеивания графических, рассмотрим процедуру слияния тех текстовых областей, которые пересекаются в горизонтальном направлении. На рис. 6 приведена графическая интерпретация оценивания параметров  $d_1, d_2$ , по которым принимается решение об объединении. Значения определяются по двум прямоугольным областям  $r_1, r_2$ , каждая из которых имеет параметр  $top$ , определяющий значение ординаты верхней стороны, и параметр  $bottom$ , соответствующий ординате нижней стороны области:  $d_1 = |top_{r_1} - top_{r_2}|$ ,  $d_2 = |bottom_{r_1} - bottom_{r_2}|$ . В качестве критерия принятия положительного решения объединения предложено использовать дизъюнкцию  $(d_1 \leq \delta_2) \vee (d_2 \leq \delta_2)$ , где  $\delta_2$  – некоторое пороговое значение, указывающее на возможное отклонение слов текста от горизонтальной оси.

Введем также в рассмотрение дополнительное пороговое ограничение  $\delta_3$ , указывающее на максимально возможное горизонтальное расстояние между двумя текстовыми областями:  $\delta_3 = left_{r_1} + width_{r_1} - left_{r_2}$ , где  $r_1, r_2$  – анализируемые области,  $left$  – абсцисса левой границы области,  $width$  – ширина области. Параметр  $\delta_3$  дает возможность варьировать расстояние между текстовыми областями, определяющими слова.

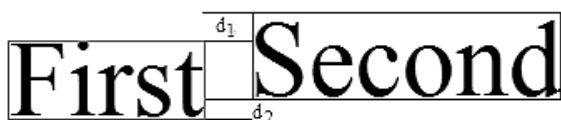


Рис.6 Графическая интерпретация значений  $d_1, d_2$

### Удаление таблиц

После выполнения всех предыдущих этапов выбор порога  $\delta_3$  с последующим слиянием может привести к тому, что таблицы на изображении документа будут определены как единый элемент (рис. 7, вверху). Для решения подобной проблемы целесообразно использовать метод для удаления линий, после чего выполнить повторное разбиение областей внутри исходной области с таблицей.

Известным методом детектирования линий на изображении является использование преобразования Хафа [16, 17]. Идея преобразования состоит в том, что каждый из значимых пикселей на изображении «голосует» за всевозможные прямые линии, частью которых он может являться. Соответственно, накопление «голосов» всех пикселей в массиве-накопителе с последующим поиском в нем глобальных и локальных максимумов позволяет определить прямые линии на изображении, за которые было отдано наибольшее число «голосов».

С практической точки зрения использование преобразования Хафа эффективно лишь для поиска линий на относительно небольших изображениях в силу его значительной трудоемкости (близка к полному перебору). Кроме того, классическая реализация преобразования Хафа позволяет определять идеальные линии, которые в результате оцифровки документа практически никогда не сохраняются.

Применение преобразования Хафа сопряжено с некоторыми нюансами рассматриваемой предметной области, которые связаны с обработкой текста. С практической точки зрения массив-накопитель содержит указатели на параметры линий, которые могут не являться частью таблицы, а проходить, например, сквозь текст. Отслеживание и отсечение подобных ситуаций можно реализовать с помощью дополнительных условий, выдвигаемых к найденной линии, подлежащей удалению, а именно: значение длины линии (его эквивалентом является значение текущего максимума в массиве-накопителе) должно превышать половину максимально возможной длины линии в этом направлении.

Практические исследования показали, что целесообразным может быть удаление всех линий в прямоугольной области, которые превышают значение в 60% от глобального максимума, но не более 20 таких линий. Поскольку алгоритмически преобразование Хафа определяет линии на изображении без анализа его контекста, то в каждом из прямоугольников также могут быть удалены линии, являющиеся частью текста, например, четко выраженные вертикальные линии в буквах «I», «Г», «K», «H» и т.п. латинского алфавита. Для устранения последствий подобного удаления предлагается использовать слияние областей по горизонтали (аналогично ранее рассмотренному этапу при поиске графических областей).

Результаты удаления вертикальных и горизонтальных линий с помощью преобразования Хафа с последующим повторным разбиением изменившейся области представлены на рис. 7 внизу.

В качестве последнего этапа обработки выполняется объединение текстовых областей, которые расположены друг над другом и имеют одинаковую ширину (в рамках погрешности  $\delta_4$ ), что соответствует объединению физических строк текста в вертикальном направлении в единственную область.

Type	Description	Example
object	The ultimate base type of all other types	object o = null;
string	String type; a string is a sequence of Unicode characters	string s = "hello";
sbyte	8-bit signed integral type	sbyte val = 12;
short	16-bit signed integral type	short val = 12;
int	32-bit signed integral type	int val = 12;
long	64-bit signed integral type	long val1 = 12; long val2 = 34L;
byte	8-bit unsigned integral type	byte val1 = 12;
ushort	16-bit unsigned integral type	ushort val1 = 12;
uint	32-bit unsigned integral type	uint val1 = 12; uint val2 = 34U;
ulong	64-bit unsigned integral type	ulong val1 = 12; ulong val2 = 34U; ulong val3 = 56U; ulong val4 = 78UL;
float	Single-precision floating point type	float val = 1.23F;
double	Double-precision floating point type	double val1 = 1.23; double val2 = 4.56D;
bool	Boolean type; a bool value is either true or false	bool val1 = true; bool val2 = false;
char	Character type; a char value is a Unicode character	char val = 'h';
decimal	Precise decimal type with 28 significant digits	decimal val = 1.23M;

Type	Description	Example
object	The ultimate base type of all other types	object o = null;
string	String type; a string is a sequence of Unicode characters	string s = "hello";
sbyte	8-bit signed integral type	sbyte val = 12;
short	16-bit signed integral type	short val = 12;
int	32-bit signed integral type	int val = 12;
long	64-bit signed integral type	long val1 = 12; long val2 = 34L;
byte	8-bit unsigned integral type	byte val1 = 12;
ushort	16-bit unsigned integral type	ushort val1 = 12;
uint	32-bit unsigned integral type	uint val1 = 12; uint val2 = 34U;
ulong	64-bit unsigned integral type	ulong val1 = 12; ulong val2 = 34U; ulong val3 = 56U; ulong val4 = 78UL;
float	Single-precision floating point type	float val = 1.23F;
double	Double-precision floating point type	double val1 = 1.23; double val2 = 4.56D;
bool	Boolean type; a bool value is either true or false	bool val1 = true; bool val2 = false;
char	Character type; a char value is a Unicode character	char val = 'h';
decimal	Precise decimal type with 28 significant digits	decimal val = 1.23M;

Рис. 7 Блок, содержащий таблицу: сверху – до удаления прямых линий, внизу – после удаления линий.

## Эксперименты

Таким образом, все параметры предложенного метода обработки исходного изображения документа могут быть представлены в виде кортежа  $M = \langle \delta_1, \delta_2, \delta_3, \delta_4 \rangle$ , где  $\delta_1$  – радиус касания областей для первичного объединения;  $\delta_2$  – максимально возможное отклонение текстовых областей вдоль вертикальной оси для объединения областей в строке текста;  $\delta_3$  – максимально возможный интервал между текстовыми областями вдоль горизонтальной оси;  $\delta_4$  – максимальное значение отклонения областей для конечного слияния текстовых областей по вертикали.

Выбор параметров  $\delta_1 - \delta_4$  определяется многими факторами, такими как разрешение и размер изображения, параметры форматирования текста на документе и т.п. Решающее значение имеет размер используемого шрифта, на который могут оказывать влияние также все перечисленные факторы. Таким образом, предварительная оценка размеров текста на

изображении позволяет выполнить грубую оценку пороговых ограничений  $\delta_1 - \delta_4$ .

На рис. 8 представлены результаты работы предложенного метода по детектированию текстовых блоков для двух изображений документов, первый из которых (рис. 8, сверху) имеет разрешение 96dpi и размер 1981x1389 пикселей, второй (рис. 8, внизу) – разрешение 300dpi при размере 1984x2928. Для обработки был использован кортеж с параметрами  $M_1 = \langle 50, 50, 50, 25 \rangle$ . Как видно из рис.8, текстовые области, параметры которых не соответствуют параметрам кортежа, образуют отдельные области.

Харьковский национальный экономический университет

### АВТОМАТИЗАЦИЯ ВЫДЕЛЕНИЯ ЗАГОЛОВКА ПО ИЗОБРАЖЕНИЮ ДОКУМЕНТА

Статья посвящена разработке метода локализации заголовка на изображении документа. Предложено использование эмпирических признаков, которые выделяются из области среднего текста, и в комбинации с применением проекционного преобразования дана возможность выделить область расположения заголовка. Экспериментальные исследования подтвердили эффективность предложенного метода и позволили определить области его применимости.

Ключевые слова: эмпирические признаки, проекция, документ, изображение, блок, проекционное преобразование

#### Введение

Задача локализации заголовка на изображении документа часто является главной для анализа структуры документа и его содержимого с последующим распознаванием [1-3], например, с целью автоматического индексирования документов.

Выделяют два направления методов анализа структуры документов: физический и логический [4], первый из которых решает задачу разбиения

Заголовок  **Выводы** документа характеризуется некоторыми общими признаками, к которым можно отнести следующие:

- заголовок чаще может иметь больший размер шрифта, чем основной текст;
- заголовок может быть написан прописными символами;

**Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**

UDC 004.932

А.В. Трохимовский, Е.О. Перезин

Харьковский национальный экономический университет

### АВТОМАТИЗАЦИЯ ВЫДЕЛЕНИЯ ЗАГОЛОВКА ПО ИЗОБРАЖЕНИЮ ДОКУМЕНТА

Статья посвящена разработке метода локализации заголовка на изображении документа. Предложено использование эмпирических признаков, которые выделяются из области среднего текста, и в комбинации с применением проекционного преобразования дана возможность выделить область расположения заголовка. Экспериментальные исследования подтвердили эффективность предложенного метода и позволили определить области его применимости.

Ключевые слова: эмпирические признаки, проекция, документ, изображение, блок, проекционное преобразование

#### Введение

Задача локализации заголовка на изображении документа часто является главной для анализа структуры документа и его содержимого с последующим распознаванием [1-3], например, с целью автоматического индексирования документов.

Выделяют два направления методов анализа структуры документов: физический и логический [4], первый из которых решает задачу разбиения изображения на текстовые или графические блоки (страницы, параграфы, абзацы и т.д.) с целью автоматического индексирования документов.

Выделяют два направления методов анализа структуры документов: физический и логический [4], первый из которых решает задачу разбиения

Заголовок  **Выводы** документа характеризуется некоторыми общими признаками, к которым можно отнести следующие:

- заголовок чаще может иметь больший размер шрифта, чем основной текст;
- заголовок может быть написан прописными символами;

**Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**  
 **Выводы**  **Выводы**  **Выводы**  **Выводы**

Рис. 8 Примеры обработки изображений текстовых документов.

На рис.9 показаны примеры обработки двух изображений документов, которые содержат графические области. Первое изображение (рис. 9, слева) имеет разрешение 600dpi, размер 4008x6167 пикселей и было обработано с помощью кортежа с параметрами  $M_2 = \langle 70, 150, 57, 165 \rangle$ , второе (рис. 9, справа, [18]) – имеет разрешение 96dpi, размер 685x920 пикселей, для его обработки использовался

кортеж  $M_3 = \langle 5, 5, 5, 5 \rangle$ . Как видно, области, содержащие графические объекты (выделены на рис. 9 полупрозрачным заполнением), были успешно сегментированы и отделены от текстовых областей. Выбор пороговых значений для кортежа  $M_3$  обоснован в первую очередь мелким размером шрифта на изображении документа.

## Выводы

В данной статье предложен метод определения местоположения текстовых областей на изображении документа с использованием слияния областей минимального размера с последующим образованием прямоугольных областей вокруг текстовых объектов. Предложенный метод прост в реализации, а также поддерживает возможность сепарации изображений и удаления таблиц.

Наилучший результат применения рассмотренного метода может быть получен с учетом имеющейся априорной информации о структуре документа и размерах используемых элементов, для большинства классов практических задач можно выделить набор пороговых ограничений непосредственно в процессе обработки.

Научной новизной статьи является предложенный метод детектирования текстовых областей на изображении документа с помощью процедуры слияния «снизу вверх», который дает возможность отделить текстовые фрагменты на изображении документов от графических объектов, выделить текст внутри таблиц, сохранив при этом достаточный для решения практических задач уровень быстродействия.

Одним из недостатков предложенного метода является необходимость подбора каждого из пороговых значений  $\delta_1 - \delta_4$ . Тем не менее, часто начальные их значения можно генерировать в автоматическом режиме на основе анализа физических параметров областей в процессе обработки. Например, значение  $\delta_1$  может быть грубо оценено как средняя ширина области на соответствующем этапе. Автоматический выбор эффективных пороговых значений, как и улучшение стабильности метода отделения графической области от текстовой может быть предметом для дальнейшей работы в этой области.

## Список литературы

1. Priti P. Rege. Text-Image Separation in Document Images Using Boundary/Perimeter Detection / Priti P. Rege, Chanchal A. Chandrakar // ACEEE International Journal on Signal & Image Processing. – Vol. 03, No. 01. – 2012. – Режим доступа: <http://hal.archives-ouvertes.fr/docs/00/74/79/44/PDF/70.pdf>.
2. Geometric Layout Analysis Techniques for Document Image Understanding: a Review [Электронный ресурс] / R. Cattoni, T. Coianiz, S. Messelodi, C. M. Modena // ITC-irst Technical Report TR#9703-09.– 1998.– Режим доступа: [http://www.ee.bgu.ac.il/~dinstein/stip2002/Seminar\\_papers/David\\_Cahana\\_Geometric%20Layout%20Analysis%20Techniques%20-%20a%20Review.pdf](http://www.ee.bgu.ac.il/~dinstein/stip2002/Seminar_papers/David_Cahana_Geometric%20Layout%20Analysis%20Techniques%20-%20a%20Review.pdf). – Название с экрана.
3. Negi, A. Localization, Extraction and Recognition of Text in Telugu Document Images / Negi A., Shanker K.N.,

**РАСПОЗНАВАНИЕ ИЗОБРАЖЕНИЙ СИМВОЛОВ ПО ИХ ЛИНЕЙНОМУ ОПИСАНИЮ**

А.В. Горюховатский, Е.О. Передрий  
Харьковский национальный экономический университет  
Харьков, пр. Ленина, 5а, goryukhovatskiy@rambler.ru

Проблема распознавания текста остается одной из наиболее актуальных в области систем компьютерного зрения ввиду ее необычайной сложности. Наиболее качественные методы распознавания символов связаны обычно с использованием нейронных сетей при значительных временных затратах на их обучение. Использование же структурных и численных признаков, которые характеризуют простотой построения в реальных задачах остается недостаточно исследованным.

Исследования посвящены разработке метода построения линейного описания изображения символа, которое является инвариантно-устойчивой характеристикой его структурных свойств. Идея метода заключается в приближенном к числовому распознавании, а именно – распознавании по форме.

В рамках предлагаемого подхода каждый символ представляется в виде иерархии переходов между некоторыми состояниями образующих его структурных элементов. Особенности построения таких описаний и их сравнения должны обеспечить устойчивость линейного описания к изменению масштаба и начертания одного класса.

Рассмотрим изображение  $I(x, y)$  как совокупность горизонтальных линий  $L = \{L_1, L_2, \dots, L_n\}$  с некоторым шагом  $\Delta y$ . Блоком назовем последовательность пикселей одинаковой яркости. В рамках каждой из линий может быть выделено  $m$  блоков  $B = \{B_1, B_2, \dots, B_m\}$ , чередование которых и образует каждую из линий  $L$ . В дальнейшем под блоком будем подразумевать лишь те их них, которые несут в рамках решаемой задачи смысловую нагрузку, в нашем случае – блоки черных пикселей. Например, рис. 1 иллюстрирует линию, проводящую сквозь символ изображения, которая состоит из четырех блоков.




Рисунок 1. Графическая интерпретация блока

Значение  $m$  формируется исходя из практических соображений после анализа множества эталонных изображений. Например, для символов латинского алфавита количество подобных блоков в линиях колеблется от 1 до 4. Необходимо учитывать, что  $m$  может быть различным даже для одного и того же набора символов различных шрифтов из-за различного качества.

Пусть  $S^i = \{s_1, s_2, \dots, s_k, i, i = 1, \dots, m$  – множество состояний, которые могут образовывать комбинации блоков в каждой линии,  $k_i$  – количество возможных состояний в множестве  $S^i$ . Формирование блоками определенного состояния может быть определено как на основании анализа только текущего блока, так и с учетом предыдущих. Например, единственный блок в линии может указывать на одно из состояний  $S^i = \{ \text{vertical line}, \text{horizontal line}, \text{diagonal line} \}$ . Возможный набор состояний может быть сгенерирован и скомпонован в группы после автоматического анализа всех эталонных изображений.

Таким образом, каждая линия  $L$  изображения отображается в одно из всех возможных состояний:  $l \Rightarrow s$ , соответственно, множество всех линий формирует набор состояний  $L \Rightarrow C = \{s\}$ , который и образует его линейное описание.

Определим набор правил, которые будут формировать состояния. Блоки и линии характеризуются параметрами, определяющими их положение и состояние линии в целом. В качестве таких параметров были выбраны их центры тяжести  $c$  и ширина  $w$ :  $b = \{c, w\}$ . Собственно правила отображения совокупности блоков внутри конкретной линии в какое-либо из состояний представляют собой набор case-ситуаций. Формированных, исходя из анализа эталонного множества. Например, следующая ситуация для  $S^1$  отображает линию  $L$ , состоящую из единственного блока  $b$ , в одно из состояний  $s_1$ :

$$w_1 > \delta w, b = \{c_1, w_1, l = \{b_1\}\} \rightarrow l \Rightarrow s_1$$

где  $\delta \in [0, 1]$  – некоторое пороговое значение,  $w$  – ширина изображения,  $s_1 \in S^1$

**Visual Code Marker Detection**

Daniel Blatnik, Adhesh Banerjee      EE-588, Spring 2006

abstract—This report discusses the algorithm we implemented to identify and read data from visual code markers. In addition to outlining the algorithm steps, we summarize the results our algorithm achieved on a set of training images.

**1. Introduction**

The problem we address is as follows: Given a JPEG image scattered with visual code markers (Fig. 1), we wish to determine the coordinates of the center of the upper-left square of each marker, as well as the bits encoded by each marker. Our strategy was to identify possible guide bars through region labeling followed by a series of checks on each pair of regions. To keep the execution time as short as possible while limiting false positives and negatives, we tried to strike a balance between robustness and efficiency. After finding a pair of guide bars, we identify the four corners of the code marker using data extracted from the Radon transform of the guide bar regions. Finally, we apply a projective transform to map the four points, which form an arbitrary quadrilateral, to a square in the conventional  $x-y$  grid. The data bits can then be read by simple thresholding.

**2. Region Labeling and Region Assembly**

The next step is to detect all of the dark regions in the image. We use 4-connectivity rather than 8-connectivity because we expect all pixels in a guide bar to be well connected to others. Next, before entering our pairwise guidebar search, we attempt to remove as many regions as possible that clearly do not fit the characteristics of a guide bar. The pairwise search will consume almost all of the execution time of our program and will have  $O(n^2)$  complexity. To remove roughly 25% of the regions, for example, would take our algorithm twice as fast. We choose to eliminate based on size any region smaller than 50 pixels or larger than 8000 pixels. The lower bound is based on the size of guide bars in the 12 training images that we were provided, in which we found no guide bar regions containing less than 100 pixels. The upper bound is based on the largest possible code marker that would fit in a 640x480 image, which would contain roughly 480x480 pixels. Its short guide bar could consume no more than 5/121 of these pixels, leaving 8000 as a reasonable upper bound.

The next characteristic of guide bars that we take advantage of is the fact that they are black regions on white backgrounds. Looking now at the red, green, and blue color components of the pixels in the regions, we remove a region from consideration if the mean squared distance of its color components from the average of its color components exceeds threshold. This works since we expect the RGB components of a black pixel to be very near each other.

**3. Binary Thresholding**

We start by converting the image to grayscale by retaining only the luminance of the original image. Color is not an immediate concern since our first goal is to label the




Figure 1: Input image containing three markers.

Рис. 9 Примеры обработки изображений текстовых документов, которые содержат графические области.

Cherreddi C.K. // *Proceedings of the Seventh International Conference on Document Analysis and Recognition.*– 2003 .– P. 1193 – 1197.– Режим доступа: <http://www.hserus.net/~cck/pubs/icdar.pdf>

4. Smith, R. Hybrid Page Layout Analysis via Tab-Stop Detection [Электронный ресурс] / R. Smith // *10th International Conference on Document Analysis and Recognition (ICDAR), 2009, 26-29 July.*– P. 214-245.– Режим доступа: <http://dejanseo.com.au/research/google/35094.pdf>. – Название с экрана.

5. Nagy, G. Hierarchical representation of optically scanned documents / G. Nagy, S. Seth // *International Conference on Pattern Recognition (ICPR), Montreal, 1984.*– P. 347-349.

6. Faure, C. Extracting the tables of contents from the images of documents [Электронный ресурс] / C. Faure // Режим доступа: <http://perso.telecom-paristech.fr/~cfaure/articles/ria00.pdf>. – Название с экрана.

7. Namboodiri, A. Document Structure and Layout Analysis [Электронный ресурс] / Аноор М. Namboodiri, Anil K. Jain // *Digital Document Processing.*– 2007.– P. 29-48.– Режим доступа: [http://pdf.aminer.org/000/348/003/document\\_page\\_segmentation\\_and\\_layout\\_analysis\\_using\\_soft\\_ordering.pdf](http://pdf.aminer.org/000/348/003/document_page_segmentation_and_layout_analysis_using_soft_ordering.pdf)

8. Learning Logic Programs for Layout Analysis Correction [Электронный ресурс] / M. Berardi, M. Ceci, F. Esposito, D. Malerba // *Proceedings of the 4th International Workshop on Multilingual OCR (MOCR '13).*– 2013.– Режим доступа: <http://www.di.uniba.it/~malerba/publications/icml03.pdf>

9. Chang, F. Chinese Document Layout Analysis Using An Adaptive Regrouping Strategy [Электронный ресурс] / Fu Chang, Shih-Yu Chu, Chi-Yen Chen // *Pattern Recognition, 38(2).* – 2005.– P. 261-271.– Режим доступа: <http://www.iis.sinica.edu.tw/papers/fchang/1566-F.pdf>. – Название с экрана.

10. Text Detection and Translation from Natural Scenes [Электронный ресурс] / J. Gao, J. Yang, Y. Zhang, A. Waibel // *School of Computer Science Carnegie Mellon University.*– 2001. – Режим доступа: <http://www.dtic.mil/dtic/tr/fulltext/u2/a455563.pdf>. – Название с экрана.

11. Bukhari, S.S. High Performance Layout Analysis of Arabic and Urdu Document Images / S.S. Bukhari, F. Shafait, T.M. Breuel // *Proceedings of the 11th International Conference on Document Analysis and Recognition (ICDAR 2011).*– 2011.– P. 1275 - 1279.

12. Historical Document Layout Analysis Competition / A. Antonacopoulos, C. Clausner, C. Papadopoulos, S. Pletschacher // *Proceedings of the 11th International Conference on Document Analysis and Recognition (ICDAR 2011).*– 2011.– P. 1516 - 1520.

13. Bukhari, S.S. Improved Document Image Segmentation Algorithm using Multiresolution Morphology / S.S. Bukhari, F. Shafait, T.M. Breuel // *Proceedings of the Document Recognition and Retrieval XVIII - DRR 2011, 18th Document Recognition and Retrieval Conference, San Jose, CA, USA, January 24-29.*– 2011.– P. 1- 10.

14. Gupta, N. Image Segmentation for Text Extraction / N. Gupta, V.K. Bange // *Proceedings of the 2nd International Conference on Electrical, Electronics and Civil Engineering (ICEECE'2012), Singapore, April 28-29.*– 2012.

15. Le, D.X. Automated Borders Detection and Adaptive Segmentation for Binary Document Images / D. X. Le, G. R. Thoma, H. Wechsler // *Proceedings of the International Conference on Pattern Recognition (ICPR '96).*– Vol. 7276. – P. 737 - 741.

16. Hough P.V.C. Methods and Means for Recognizing Complex Patterns. U.S. Patent 069654, 1962.

17. Duda, R.D. Use of the Hough transform to detect lines and curves in pictures [Электронный ресурс] / R.D. Duda, P.E. Hart // *Communications of the ACM.* – 1972.– Vol. 15, No 1.– P. 11–15.– ISSN:0001-0782. – Режим доступа: [www.ai.sri.com/pubs/files/tm036-duda71.pdf](http://www.ai.sri.com/pubs/files/tm036-duda71.pdf)

18. Blatnik, D. Visual Code Marker Detection [Электронный ресурс] / D. Blatnik, A. Banerjee // Режим доступа: [http://www.stanford.edu/class/ee368/Project\\_06/Project/ee368\\_reports/ee368group12.pdf](http://www.stanford.edu/class/ee368/Project_06/Project/ee368_reports/ee368group12.pdf). – Название с экрана.

**Рецензент:** д-р техн. наук, проф. Пуятин Е.П., Харьковский национальный университет радиоэлектроники

**Авторы:**

**Гороховатский Алексей Владимирович**

Харьковский национальный экономический университет им. С. Кузнеця,

Харьков, кандидат технических наук, доцент кафедры информатики и компьютерной техники.

Тел. 702-06-74 (4-38), E-mail: [gorohovatsky@rambler.ru](mailto:gorohovatsky@rambler.ru)

## Детектування текстових областей на зображенні документа методом злиття

О.В. Гороховатський

В статті запропоновано метод автоматичного детектування текстових областей на зображенні за допомогою злиття областей. Запропонований підхід комбінє простоту реалізації алгоритму злиття із аналізом та відокремленням графічних областей, котрі не несуть текстової інформації. Проведено моделювання роботи запропонованого методу, визначено класи практичних задач та умови найбільш ефективного його застосування.

**Ключові слова:** текстова область, графічна область, злиття, зображення документа, порогове обмеження, перетворення Хафа.

## The detection of text regions on image of a document using merge method

O.V. Gorokhovatskyi

Paper is devoted to the construction of the method of automatical detection of text regions on image of a document using regions merging. Proposed approach combines the simplicity of merging algorithm implementation with an opportunity of removing image elements which do not contain text information. Experimental investigations were performed, classes of practical tasks to use proposed approach in the most effective way were distinguished.

**Keywords:** text region, graphic region, merging, image of a document, threshold limitation, Hough transform.