

АЛЕКСАНДР ВАЛЕРЬЕВИЧ КУЧУГАНОВ, ДЕНИС РАШИДОВИЧ  
КАСИМОВ, ВАЛЕРИЙ НИКОНОВИЧ КУЧУГАНОВ,  
ПАВЕЛ ПЕТРОВИЧ ОСКОЛКОВ

*Ижевский Государственный Технический Университет им. М.Т. Калашникова (Россия)*

## **Автоматизация распознавания старославянских рукописей с помощью дескрипционных логик**

В статье предлагается подход к распознаванию старославянских символов, включающий этапы тернарной сегментации, выделения границ цветовых областей, синтеза двух вариантов скелетона, формирования нечеткого графа, создания логического описания в нечеткой дескрипционной логике и распознавания с помощью подсистемы формальных автоматических рассуждений. Приводятся результаты предварительных экспериментов, которые подтверждают перспективность предложенного подхода<sup>1</sup>.

**Ключевые слова:** распознавание, скелетон, структурные элементы, нечеткий граф, дескрипционная логика, старопечатные тексты, старославянский язык, древнерусский язык.

### **1. Введение**

Автоматизация процесса перевода старославянских текстов в электронную форму, позволяющую применять компьютерные методы обработки текстовой информации, существенно упрощает проведение исторических и лингвистических исследований. Традиционные OCR-системы не справляются с распознаванием древних текстов. Специализированные методы и системы демонстрируют недостаточную надежность распознавания — в среднем не выше 80% на манускриптах неплохого качества.

Авторы работы *Feature Selection for Classification of Old Slavic Letters* [Bande et al. 2014] предлагают два способа распознавания старославянских кириллических символов: на основе “деревьев решений” и путем нечеткой классификации. Оба метода оперируют статистическими и простыми структурными признаками символов (свыше 20 признаков).

Экспериментальная система распознает символы каждым из способов со средней точностью и полнотой 70–80%. Следует, однако, отметить сложность настройки системы простым пользователем (не экспертом). При этом возможность простой и быстрой настройки системы (в частности, под конкретного писца) является весьма эффективным способом повышения качества распознавания.

В другой работе — *Word-Based Adaptive OCR for Historical Books* [Kluzner et al. 2009] — описана методика распознавания целых слов, а не отдельных символов исторических текстов. Сначала производится выделение отдельных слов на изображении, затем выполняется кластеризация полученного множества фрагментов-слов таким образом, чтобы в каждом кластере были представлены экземпляры одного и того же слова. После кластеризации производится распознавание изображений слов с помощью традиционного OCR-средства. К примеру, книгу 18 века, написанную немецким готическим шрифтом, разработанная программная система корректно распознала на 86,6%, что на 4,1% выше результатов традиционной OCR-системы. Однако, подходы пословного распознавания плохо применимы к кириллическим текстам, т.к. слова в них не отделены друг от друга и между буквами имеется относительно большое расстояние. Выделение отдельных слов в таких текстах представляет непростую задачу даже для человека-неспециалиста.

## 2. Предлагаемый подход

Предлагаемый нами подход к распознаванию старославянских символов схематично представлен на рис. 1.

Сначала производится скелетизация изображений символов, состоящая из следующих этапов: 1) тернарная сегментация изображения; 2) выделение границ черных областей и границ объединения серых и

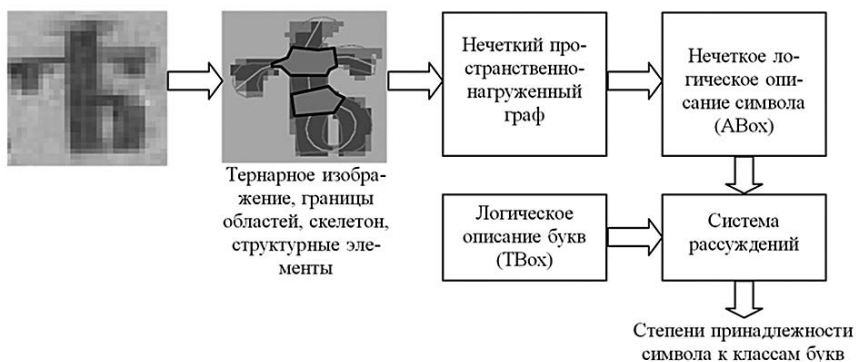


Рис. 1. Процесс описания, анализа и распознавания символов

черных областей; 3) синтез двух вариантов скелетона: черного и серо-черного; 4) аппроксимация скелетона; 5) идентификация областей перекрестков [Кучуганов 2012].

Затем выделяются конструктивные элементы (мачты, перекладины и т.п.) символа путем анализа формы отдельных цепочек элементов скелетона, а также соединительные элементы (области перекрестков).

Формируется нечеткий атрибутивный пространственно нагруженный граф изображения символа, вершины и ребра которого содержат количественные и качественные значения пространственных атрибутов, характеризующих конструктивные элементы символа.

Система имеет базу знаний о символах в форме ТВох-описания в нечеткой дескрипционной логике f-SHIN [Straccia 2001]. Используемые элементы:

1. **Атомарные концепты:** КонЭл (конструктивные элементы), СоедЭл (соединительные элементы), ПрямойФормы (объекты с прямой формой), КруглойФормы, ОвальнойФормы, ИзвилистойФормы, ДугообразнойФормы, ВобразнойФормы, УобразнойФормы, Замкнутый (замкнутые объекты), ВертОриентации (объекты, имеющие ориентацию “Север” или “Юг”), ГоризОриентации, НаклВлОриентации, НаклВпОриентации, Длинный (объекты большой длины), СреднейДлины, Короткий, Толстый (объекты большой ширины на протяжении всей длины), СреднейТолщины, Тонкий, ТолстыйВСередине (объекты большой ширины в середине) и др.
2. **Роли:** соед1L (соединительный элемент точкой 1 примыкает к точке L конструктивного элемента), соедL1  $\equiv$  соед1L– (роль соед1L является обратной по отношению к роли соедL1), соедDL (конструктивный элемент точкой D примыкает к точке L другого конструктивного элемента), соед7B1E (соединительный элемент точками 7 и 1 примыкает соответственно к точкам B и E конструктивного элемента), соед7B1E  $\equiv$  соед1E7B, соедB7E1  $\equiv$  соед7B1E– и т.п.

На рис. 2 представлены коды точек примыкания элементов. Точки примыкания соединительных элементов кодируются значениями их ориентации относительно центроидов элементов. Точки примыкания конструктивных элементов кодируются в зависимости от взаимного расположения по вертикали и горизонтали (верхняя или нижняя, левая или правая), а также по принципу выделения начальной и конечной точек в соответствии с направлением движения по элементу против часовой стрелки.

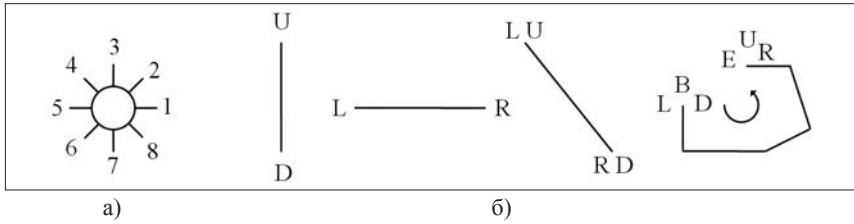


Рис. 2. Коды точек примыкания: а) у соединительных элементов; б) у конструктивных элементов

3. **Производные концепты:** Линия  $\equiv$  КонЭл  $\cap$   $\neg$ Замкнутый  $\cap$  ПрямойФормы; Мачта  $\equiv$  Линия  $\cap$  ВертОриентации; Перекладина  $\equiv$  Линия  $\cap$  ГоризОриентации; НаклВлЛиния  $\equiv$  Линия  $\cap$  НаклВлОриентации; НаклВпЛиния  $\equiv$  Линия  $\cap$  НаклВпОриентации; Петля  $\equiv$  КонЭл  $\cap$   $\neg$ Замкнутый  $\cap$  ДугообразнойФормы; Круг  $\equiv$  КонЭл  $\cap$  Замкнутый  $\cap$  КруглойФормы; Овал  $\equiv$  КонЭл  $\cap$  Замкнутый  $\cap$  ОвальнойФормы и др. В таблице 1 представлен фрагмент описания букв и вспомогательных структур.

Таблица 1. Описание конфигураций (букв и вспомогательных структур)

Производный концепт	Примеры экземпляров	Определение
ВерхняяМачта-БуквыЪ		$\equiv$ Мачта $\cap$ СреднейДлины $\cap$ Толстый $\cap$ ЭсоедD3.(СоедЭл $\cap$ Эсоед7U.(Мачта $\cap$ (СреднейДлины $\cup$ Короткий) $\cap$ Толстый $\cap$ ЭсоедDL.(Перекладина $\cap$ СреднейДлины $\cap$ ТонкийВСередине)) $\cap$ Эсоед8E.(Петля $\cap$ СреднейДлины $\cap$ (Толстый $\cup$ СреднейТолщины)))
ВерхняяПравая-Перекладина		$\equiv$ (Перекладина $\cup$ НаклВпЛиния) $\cap$ (СреднейДлины $\cup$ Короткий) $\cap$ ТолстыйСлева $\cap$ ТонкийВСередине $\cap$ ТолстыйСправа
ВерхняяВер-тМачтаБуквыЪ		$\equiv$ ВерхняяМачтаБуквыЪ $\cap$ ЭсоедUL.ВерхняяПраваяПерекладина
ВерхняяМачта-БуквыЪ		$\equiv$ ВерхняяМачтаБуквыЪ $\cap$ ЭсоедUR.ВерхняяЛеваяПерекладина

Для распознаваемого символа строится нечеткое АВох-описание. В примере на рис. 3 сказано, что  $v1$  является экземпляром концепта КонЭл с полной уверенностью, экземпляром концепта ВертОриентации – со степенью принадлежности 0.3; между  $v1$  и  $v2$  существует отношение  $\text{соед}RU$  с полной уверенностью, между  $v2$  и  $v3$  – отношение  $\text{соед}D3$  со степенью уверенности 0.8 и т.д.

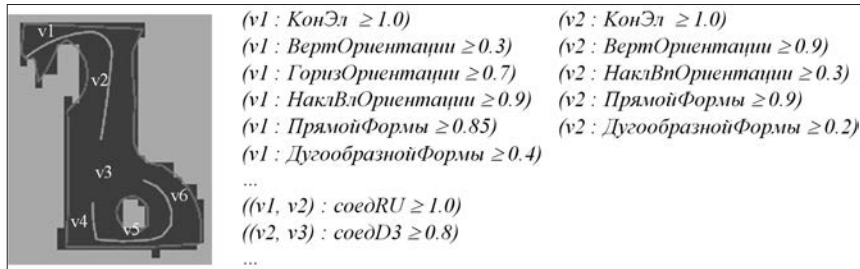


Рис. 3. Пример распознаваемого символа и его нечеткого логического описания

Распознавание символа мы сводим к выполнению в подсистеме автоматических рассуждений (reasoner) запросов следующего содержания: “определить элемент символа, который с высокой степенью принадлежности является экземпляром концепта ВерхняяМачтаБуквыГ”. В примере, представленном на рис. 3, подсистема рассуждений определит, что элемент  $v2$  является экземпляром концепта ВерхняяМачтаБуквыГ со степенью принадлежности, равной 0,93.

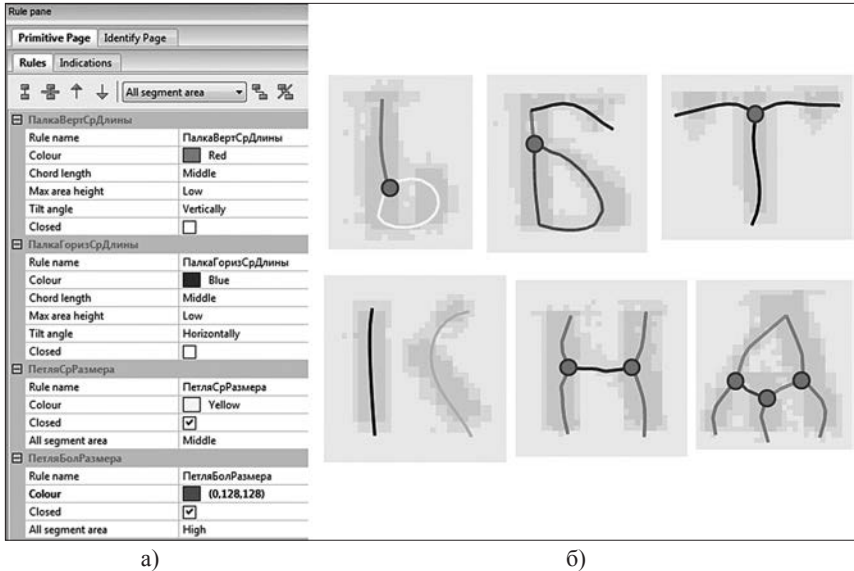
### 3. Описание реализации

В базе знаний системы имеются логические определения (DL-описания) конструктивных элементов, из которых состоят старославянские буквы. Интерфейс ввода определений конструктивных элементов представлен на рис. 4(а). Определением элемента является набор отличительных признаков и их возможных значений. Значения задаются как в четкой, так и нечеткой форме.

Конструктивные элементы описываются следующими атрибутами: длина хорды, количество сегментов слева и справа от хорды, площадь каждого сегмента, площадь левых/правых сегментов, площадь всех сегментов, максимальная высота сегмента, угол наклона элемента, замкнутость.

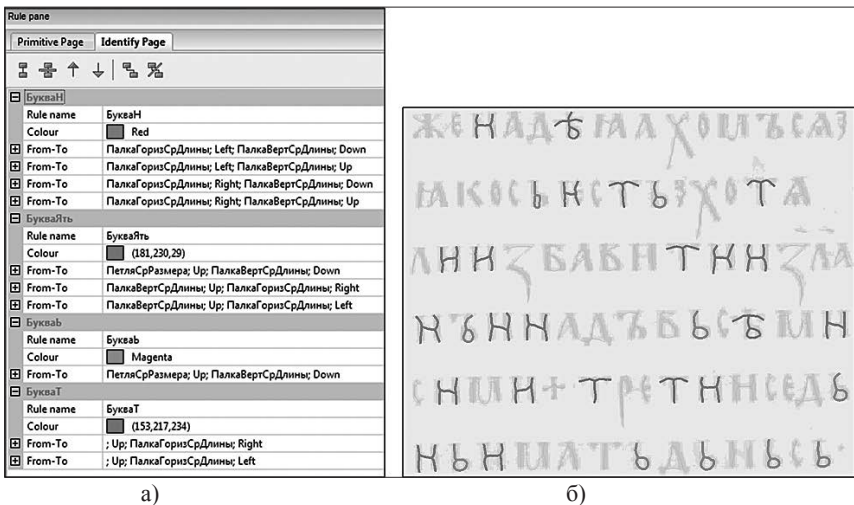
Один и тот же элемент может подходить под несколько определений с той или иной степенью уверенности (принадлежности).

Примеры конкретных элементов, соответствующих введенным определениям, показаны на рис. 4(б).



**Рис. 4.** База знаний (TBox) о конструктивных элементах символов:  
 а) пример DL-описаний конструктивных элементов; б) пример распознавания конструктивных элементов по логическим описаниям

База знаний системы содержит также логические определения букв. На рис. 5(а) приведен интерфейс ввода определений букв. Буква рассматривается как конструкция из элементов, определенных ранее. В определении буквы указываются входящие в ее состав элементы и способы



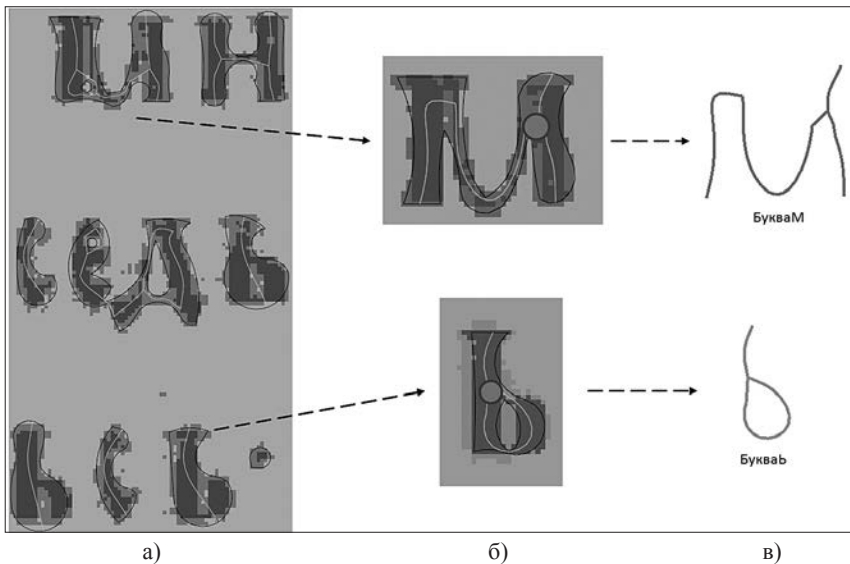
**Рис. 5.** База знаний (TBox) о символах: а) пример DL-описаний символов; б) пример распознавания символов по логическим описаниям

примыкания их друг к другу, т.е. какими точками они друг с другом соединяются.

На рис. 5 (б) приведены примеры конкретных букв, соответствующих введенным определениям. Распознавание осуществляет подсистема логических рассуждений.

Один и тот же символ подсистема рассуждений может отнести сразу к нескольким определениям букв. Каждый вариант распознавания обладает определенной степенью уверенности (принадлежности).

Для нераспознанных символов осуществляется попытка построения альтернативных версий скелетона. Для этого производится адаптивная цветовая сегментация изображения символа, которая позволяет получить более качественный скелетон и успешно распознать символ (рис. 6).



**Рис. 6.** Адаптивная цветовая сегментация нераспознанных символов:  
 а) исходный скелетон; б) результат адаптивной сегментации;  
 в) результат распознавания

Результатом работы модуля графического распознавания являются варианты значений каждого символа (рис. 7).

Эти данные далее передаются в модуль уточнения по грамматическому словарю древнерусского языка. Словарь создан в рамках проекта “Манускрипт” (manuscripts.ru) под руководством В.А. Баранова.

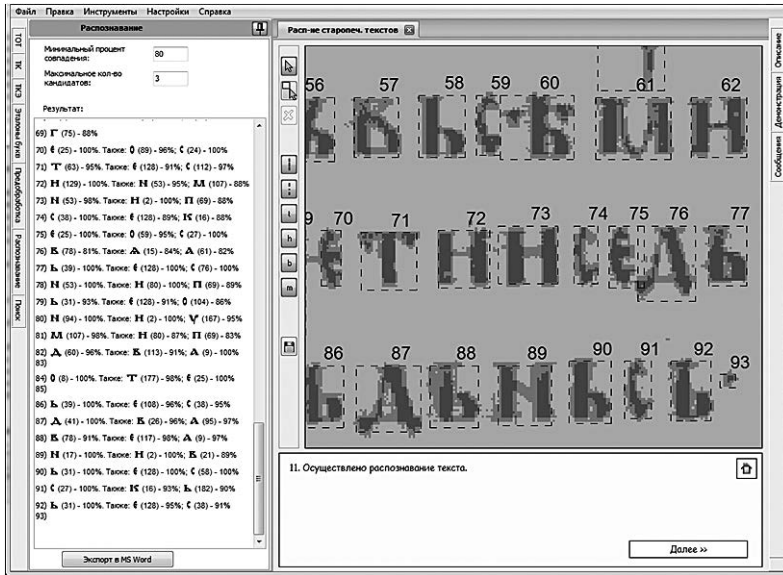


Рис. 7. Пример результатов графического распознавания символов

#### 4. Эксперименты

В экспериментальном исследовании использовался текст факсимильного издания Остромирова Евангелия 1056–1057 гг. Тестовый набор содержал 10 страниц данного текста.

В таблице 2 представлены значения показателей полноты и точности распознавания по отдельным буквам, полученные в нашей системе, а также в системе исследователей из Македонии [Bande et al. 2014] (по данным авторов).

Показатели вычислялись следующим образом: пусть  $S$  — множество изображений отдельных символов на старопечатных текстах из тестового набора;  $b$  — некоторая буква алфавита;  $S_b \subset S$  — подмножество символов, являющихся буквой  $b$ ;  $R_b \subset S$  — подмножество символов, которые система отнесла к букве  $b$ . Тогда:

$$\text{Полнота}_b = |R_b \cap S_b| / |S_b|,$$

$$\text{Точность}_b = |R_b \cap S_b| / |R_b|.$$

Для нашей системы в таблице 2 дополнительно приведены значения полноты с учетом кандидатов на результат распознавания, которые дают оценку того, какое качество распознавания можно достичь за счет последующего уточнения с помощью словаря древнерусского языка.

Полученные результаты показывают, что предложенный подход распознает символы со средней полнотой 86%, точностью 85%, что не-



Таблица 2. Сравнительные оценки качества распознавания символов

Буква	Наш экспериментальный модуль распознавания			Классификатор на основе дерева решений (Македония)	
	Полнота	Точность	Полнота с учетом потенциальных (для последующего уточнения с помощью словаря)	Полнота	Точность
А	0.75	0.87	1	0.5	0.63
Б	0.89	0.89	1	0.67	1
В	0.88	0.60	1	0.88	0.64
Г	0.81	0.59	1	0.56	1
Д	0.85	0.95	1	0.89	1
Е	0.80	0.92	0.90	0.8	1
Ж	0.83	0.89	1	1	0.82
З	0.86	0.74	1	1	0.82
И	0.92	0.90	1	0.86	0.86
Й	0.91	0.87	1	1	0.82
К	0.82	0.95	1	0.86	0.58
Л	0.79	0.70	0.85	0.44	0.67
М	0.96	0.92	1	0.71	0.71
Н	1	1	1	0.79	0.65
О	0.86	0.87	1	нет данных	
П	0.83	1	1	0.4	1
Ср.:	0.86	0.85	0.98	0.76	0.81

сколько выше, чем у известных аналогов. Здесь следует учитывать, что эксперименты в нашей системе и в аналогичных системах проводились на разных тестовых наборах изображений.

Также был произведен эксперимент по распознаванию текста Постной Триоди (РНБ, Погод. 41) (рис. 8).

База знаний была предварительно пополнена дополнительными вариантами описаний букв, поскольку в Постной Триоди (РНБ, Погод. 41)



Рис. 8. Фрагмент текста Постной Триоди (РНБ, Погод. 41)

написание некоторых букв (таких как “веди”, “земля”, “иже”) немного отличается от написаний в Остромировом Евангелии.

Анализировались две версии скелетона: 1) скелетон черного и 2) скелетон серо-черного. Обе версии скелетона представлены на рис. 9. Черным цветом отображены правильно распознанные символы; белым цветом показаны нераспознанные фрагменты скелетонов.



*Рис. 9. Результаты распознавания символов: а) символы, правильно распознанные по скелетону черного; б) символы, правильно распознанные по скелетону серо-черного*

В общей сложности на данном изображении программа правильно распознала 287 из 354 символов (81%). Распознавание надстрочных знаков в данной версии системы не реализовано.

### **Заключение**

Таким образом, применение дескрипционной логики позволяет непротиворечиво и структурировано представить синтаксис изображения старославянского символа в удобочитаемом формате (как для человека, так и для компьютера) и производить анализ средствами формальных логических рассуждений.

Новизна метода распознавания состоит в адаптивной тернарной сегментации изображения и синтезе описания на основе аппарата дескрипционной логики.

Дальнейшее совершенствование подхода предполагает: усиление алгоритма адаптивной цветовой сегментации, усиление выразительности дескрипционного описания букв, автоматический синтез эталонных DL-описаний новых вариантов букв, разработку модуля уточнения результатов распознавания на основе словарей древнерусского языка.

ПРИМЕЧАНИЕ

<sup>1</sup> Работа выполнена при частичной поддержке РФФИ (№ 15-07-08077 А) и Государственного задания МОиН РФ (Проект № 625).

ЛИТЕРАТУРА

- Кучуганов А. В., 2012.: Когнитивный алгоритм построения геометрического остова невыпуклых фигур, *Приволжский научный журнал*, № 3 (23). Н. Новгород, 84–89.
- Bande С.М., Klekovska М., Nedelkovski I. et al., 2014: Feature Selection for Classification of Old Slavic Letters, *Control Engineering and Applied Informatics*, 16.4. 81–90.
- Kluzner, V., Tzadok A., Shimony Y. et al., 2009: Word-Based Adaptive OCR for Historical Books, *Proceedings of the 10th International Conference on Document Analysis and Recognition (ICDAR2009)*. 501–505.
- Straccia U., 2001: Reasoning within fuzzy description logics, *Journal of Artificial Intelligence Research*, 14. 137–166.

BIBLIOGRAPHY (TRANSLITERATION)

- Bande С.М., Klekovska М., Nedelkovski I. et al., 2014: Feature Selection for Classification of Old Slavic Letters, *Control Engineering and Applied Informatics*, 16.4. 81–90.
- Kluzner, V., Tzadok A., Shimony Y. et al., 2009: Word-Based Adaptive OCR for Historical Books, *Proceedings of the 10th International Conference on Document Analysis and Recognition (ICDAR2009)*. 501–505.
- Kuchuganov A. V., 2012.: Kognitivnyy algoritm postroeniya geometricheskogo ostova nevyuklykh figur, *Privolzhskiy nauchnyy zhurnal*, № 3 (23). N. Novgorod, 84–89.
- Straccia U., 2001: Reasoning within fuzzy description logics, *Journal of Artificial Intelligence Research* 14. 137–166.

ALEKSANDER KUCHUGANOV, DENIS KASIMOV,  
VALERIY KUCHUGANOV, PAVEL OSKOLKOV

**Automation of Recognizing Old Slavonic Manuscripts  
with help of Description Logics**

In the paper an approach to recognition of old Slavonic symbols is proposed, which includes the stages of ternary segmentation, detecting edges of areas, synthesis of two skeleton variants, forming a fuzzy graph, creating a logical description in a fuzzy description logic, and recognition with help of a subsystem of formal automatic reasoning. Results of preliminary experiments are provided, which confirm the perspective of the proposed approach.

**Keywords:** recognition, skeleton, structural elements, fuzzy graph, description logic, old printed texts, Old Church Slavonic language, Old Russian language.

ALEKSANDR KUČUGANOV, DENIS KASIMOV,  
VALERIJ KUČUGANOV, PAVEL OSKOLKOV

### Senujų slavų rankraščių tekstų atpažinimo automatizavimas deskriptinės logikos pagalba

Straipsnio autoriai siūlo naują senosios slavų kalbos simbolių atpažinimo metodiką, kurios etapai: ternarinė segmentacija, spalvų sričių ribų išskyrimas, dviejų skeleto variantų sintezė, neryškaus grafo formavimas, loginio aprašo kūrimas netikslioje diskriptinėje logikoje ir senųjų simbolių atpažinimas automatinės sistemos pagalba. Straipsnyje pateikiami preliminarių bandymų rezultatai, patvirtinantys siūlomo atpažinimo metodo perspektyvumą.

**Reikšminiai žodžiai:** atpažinimas, struktūriniai elementai, neryškus grafas, deskriptinė logika, senųjų spausdinių tekstai, senoji slavų kalba, senoji rusų kalba.

Поступило в редакцию: 16 марта 2018 г.

Принято к печати: 15 мая 2018 г.

Александр Валерьевич Кучуганов, к.т.н., доцент кафедры Автоматизированные системы обработки информации и управления Ижевского государственного технического университета им. М.Т. Калашникова.

Aleksander Kuchuganov, Ph.D., Assoc. prof., Department of Automated Data Processing and Control Systems, Kalashnikov Izhevsk State Technical University.

Aleksandr Kučiuganov, technikos mokslų daktaras, Iževsko valstybinio M.T. Kalašnikovo technikos universiteto Automatizuotų informacijos ir valdymo apdorojimo sistemų katedros docentas.

E-mail: aleks\_kav@udm.ru.

Денис Рашидович Касимов, к.т.н., доцент кафедры Автоматизированные системы обработки информации и управления Ижевского государственного технического университета им. М. Т. Калашникова.

Denis Kasimov, Ph.D., Assoc. prof., Department of Automated Data Processing and Control Systems, Kalashnikov Izhevsk State Technical University.

Denis Kasimov, technikos mokslų daktaras, Iževsko valstybinio M. T. Kalašnikovo technikos universiteto Automatizuotų informacijos ir valdymo apdorojimo sistemų katedros docentas.

E-mail: kasden@mail.ru

Валерий Никонорович Кучуганов, д.т.н., профессор, заведующий кафедрой Автоматизированные системы обработки информации и управления Ижевского государственного технического университета им. М. Т. Калашникова.

Valerij Kuchuganov, Habil. prof., Head, Department of Automated Data Processing and Control Systems, Kalashnikov Izhevsk State Technical University.

Valerij Kučiuganov, habilituotas technikos mokslų daktaras, profesorius; Iževsko valstybinio M.T. Kalašnikovo technikos universiteto Automatizuotų informacijos ir valdymo apdorojimo sistemų katedros vedėjas.

E-mail: kuchuganov@istu.ru;

Павел Петрович Осолков, инженер кафедры Автоматизированные системы обработки информации и управления Ижевского государственного технического университета им. М. Т. Калашникова.

Pavel Oskolkov, engineer, Department of Automated Data Processing and Control Systems, Kalashnikov Izhevsk State Technical University.

Pavel Oskolkov, inžinierius, Iževsko valstybinio M.T. Kalašnikovo technikos universiteto Automatizuotų informacijos ir valdymo apdorojimo sistemų katedra.

E-mail: aleks\_kav@udm.ru

XIII в. кн. 1.

**МѢСЯЦЕВЪ АКАСТА**

✠ ПІСА ПІЕРІММА ТАА ПІРІА

**П**РІЖАУ ТАНИИ ТИ ПІРИ Б ТАЦО ПУТІАКА  
СТЕЖІА ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА

**П**ІМЦІ ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА  
ПІРІА ПІРІА ПІРІА ПІРІА ПІРІА

Тобинъ 1855  
Порфиръ 1855  
Мелгаровъ 1890  
Сидоръ 1890