

ISSN 1816-0301 (print)
УДК 004.912

Поступила в редакцию 21.12.2017
Received 21.12.2017

Ю. Б. Крапивин

Брестский государственный технический университет, Брест, Республика Беларусь

СИСТЕМА «ПлагиаТКонтроль» КАК ИНСТРУМЕНТ ЭКСПЕРТИЗЫ ТЕКСТОВЫХ ДОКУМЕНТОВ

Аннотация. Приводятся описание и анализ работы инструментально-программного комплекса «ПлагиаТКонтроль», позволяющего автоматизировать решение задачи распознавания в заданном текстовом документе фрагментов, заимствованных как из локальной полнотекстовой базы данных пользователя, так и из сети Интернет. Комплекс обеспечивает решение задачи с учетом не только явного, но и неявного заимствований с точностью до парадигм лексических единиц и отношений лексической и грамматической синонимии в соответствии с приведенной структурно-функциональной схемой системы автоматического распознавания заимствованных фрагментов. «ПлагиаТКонтроль» предоставляет возможность осуществлять работу в различных режимах, позволяя автоматизировать труд эксперта и существенно ускорить процедуру анализа документов на предмет наличия в них заимствований (плагиата) из других текстовых документов.

Ключевые слова: естественный язык, автоматическая обработка текстов, заимствованный фрагмент, cross-language-функциональность

Для цитирования. Крапивин, Ю. Б. Система «ПлагиаТКонтроль» как инструмент экспертизы текстовых документов / Ю. Б. Крапивин // Информатика. – 2018. – Т. 15, № 1. – С. 103–109.

Y. B. Krapivin

Brest State Technical University, Brest, Republic of Belarus

SYSTEM «PlagiarismControl» AS THE TOOL FOR THE EXPERTISE OF THE TEXT DOCUMENTS

Abstract. The description and the operability analysis of the implemented instrumental software system «PlagiarismControl» has been done. The system affords to automatize solving the task of the identification of the adopted fragments in the given text document both from the local full-text user's database and from the Internet. The system affords solving the task taking in account explicit as well as implicit adoptions with precision up to lexical units paradigms and both lexical and grammatical synonymy relations, according to the structural-functional schematic diagram of the system of the automatic recognition of reproduced fragments of the text documents. «PlagiarismControl» is able to work in different modes, to automatize the work of the expert and to speed up significantly the procedure of the analysis of the documents, with the purpose of recognition of the adoptions (plagiarism) from other text documents.

Keywords: natural language, automatic text processing, adopted fragment, cross-language functionality

For citation. Krapivin Y. B. System «PlagiarismControl» as the Tool for the Expertise of the Text Documents. *Informatics*, 2018, vol. 15, no. 1, pp. 103–109 (in Russian).

Введение. Быстрое развитие и проникновение информационных технологий во все сферы жизни человека сделали сегодня вполне естественным то, о чем еще несколько десятилетий назад можно было только мечтать: в несколько кликов мыши найти и обработать практически любую интересующую информацию, не обладая при этом специальными знаниями, а опираясь лишь на базовые принципы компьютерной грамотности, приобретаемые зачастую самостоятельно. Существующие информационно-поисковые системы и многочисленные специализированные информационные ресурсы в сети Интернет позволяют это сделать за считанные секунды. Естественно, подготовка практически любой квалификационной или исследовательской работы, начиная от школьного реферата и заканчивая диссертацией, также так или иначе основывается на таких возможностях. В связи с этим весьма актуальными являются проблемы обнару-

жения в текстовых документах фрагментов, заимствованных из других источников, каковыми все чаще становятся именно интернет-доступные документы, и анализа заимствованных фрагментов (ЗФ) с целью установления (неустановления) фактов цитирования, т. е. плагиата. При этом особенно трудоемкой является первая задача, которая прежде всего требует разработки эффективных средств автоматизации.

Для решения указанных задач автором был разработан инструментально-программный комплекс (ИПК) «ПлагиатКонтроль», ориентированный на обработку текстовых документов, в первую очередь диссертационных и дипломных работ, представленных как на русском, так и белорусском языках, с целью автоматического распознавания в них на лексико-грамматическом уровне заимствований из других текстовых документов, представленных в поисковом пространстве (ПП) на том же языке, что и входной документ, либо на обоих указанных языках (cross-language-функциональность).

Описание системы. Поисковое пространство, с которым работает ИПК, может включать в себя текстовые интернет-документы и документы из поисковой БД пользователя, например уже существующие у него или доступные ему коллекции диссертационных или дипломных работ. ИПК обеспечивает решение задачи на лексико-грамматическом уровне, т. е. с учетом не только явного, но и неявного заимствования с точностью до парадигм лексических единиц и отношений лексической и грамматической синонимии, а также с точностью до указанной пары языков, в соответствии с приведенной ниже структурно-функциональной схемой системы автоматического распознавания ЗФ (рис. 1), которая была разработана на случай решения задачи в n -язычной информационной среде, $n \geq 1$.

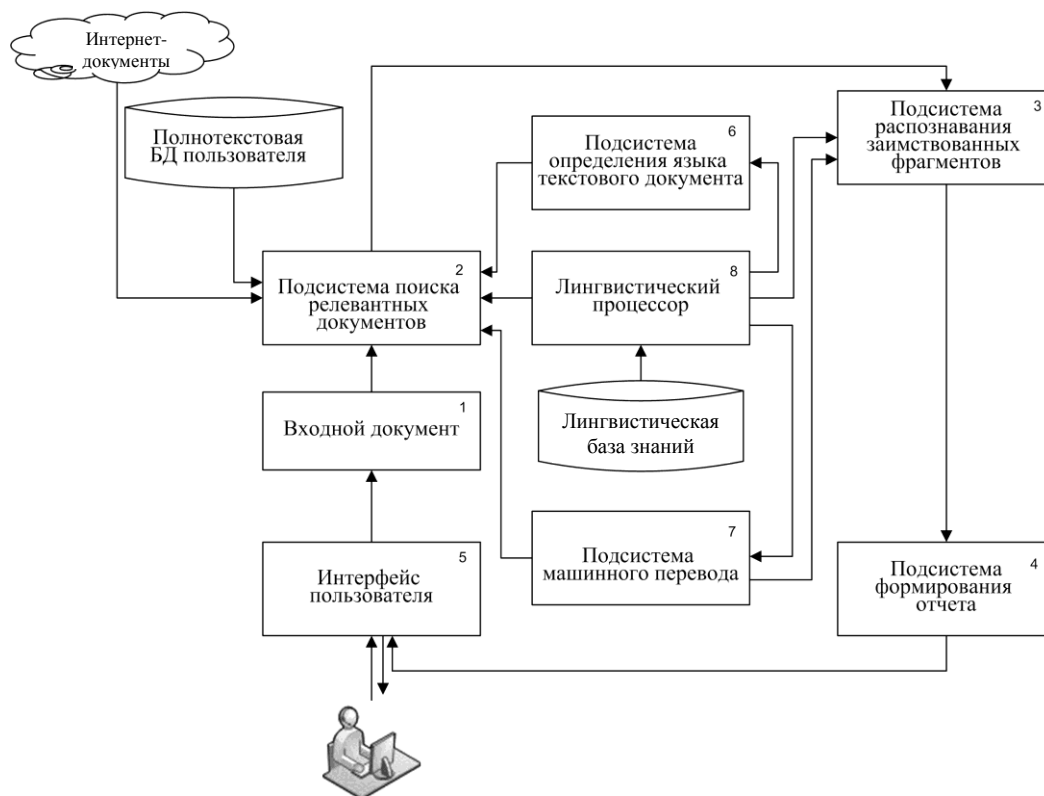


Рис. 1. Структурно-функциональная схема системы автоматического распознавания ЗФ

Входной информацией для системы является подлежащий анализу на предмет заимствования входной документ (блок 1), представленный в одном из наиболее распространенных форматов: TXT, RTF, DOC, DOCX, PDF, HTML. Ориентируясь на огромные объемы ПП и учитывая жесткие ограничения, накладываемые на время реакции системы автоматического распознавания ЗФ, этот документ в общем случае сначала поступает в подсистему поиска релевант-

ных документов (блок 2). Ее задача – за минимальное время найти в используемом ПП множество наиболее релевантных входному документу текстовых документов, представленных как на русском, так и на белорусском языках (т. е. построить так называемое минимизированное поисковое пространство (МПП)). Эти документы далее поступают в подсистему, собственно, распознавания заимствованных фрагментов (блок 3) и затем – в подсистему формирования отчета (блок 4) с последующим предоставлением его пользователю. Взаимодействие с системой осуществляется посредством интерфейса пользователя (блок 5), который поддерживает ввод документов и просмотр результатов поиска заимствований.

Очевидно, что функционирование системы в режиме cross-language требует функциональности подсистем определения языка текстового документа (блок 6) и машинного перевода (блок 7). При этом функциональность последней подсистемы в зависимости от количества используемых языков и параметров распределения на них текстовых документов в ПП ориентирована на перевод текстов самих документов или их поисковых образов. Первый случай может иметь место как при поиске релевантных документов, так и на этапе распознавания ЗФ. Например, найденные на основе перевода поискового образа документа релевантные интернет-документы будут переводиться в блоке 3, а для документов из полнотекстовой БД пользователя с целью минимизации общего времени решения задачи в ряде случаев целесообразно заранее получить эквиваленты на других естественных языках и хранить их в этой же БД. Отметим, что в том частном случае, когда имеет место одноязычная информационная среда, подсистема поиска релевантных документов не использует функциональность подсистем блоков 6 и 7. Более того, если ПП ограничивается только документами полнотекстовой БД пользователя, которая является относительно постоянной и небольшой (от нескольких единиц до нескольких сотен документов), то эта подсистема, «не включаясь», сразу передаст входной документ в подсистему распознавания ЗФ. Что касается лингвистического процессора (блок 8) и используемой им лингвистической базы знаний, то они обеспечивают автоматический лингвистический, в том числе и лингвостатистический, анализ текстовых документов в той мере, в какой это необходимо для подсистем блоков 2, 3, 6, 7, что, в свою очередь, зависит от используемых методов решения соответствующих задач. Можно говорить, что задачи автоматического лингвистического анализа текстовых документов (блок 8), а также распознавания их языка (блок 6) и машинного перевода (блок 7) составляют базовую лингвистическую обработку текстовых документов при решении задачи распознавания ЗФ.

Таким образом, в соответствии с приведенной схемой задача автоматического распознавания ЗФ в целом включает следующие подзадачи: поиска релевантных входному документу текстовых документов в полнотекстовой БД и в сети Интернет, автоматического распознавания языка текстового документа, машинного перевода текстовых документов и их поисковых образов в множестве заданных естественных языков, автоматического распознавания заимствованных фрагментов текстовых документов, автоматического построения отчета, автоматического лингвистического анализа текстовых документов.

Требуемая базовая автоматическая лингвистическая обработка текстовых документов в ИПК «ПлагиатКонтроль» осуществляется с помощью лингвистического процессора известной системы машинного перевода с белорусского языка на русский и обратно (СМП Б/Р) [1]. Его функциональность, как показал анализ, полностью соответствует (с точки зрения постановки задачи ИПК) требованиям, предъявляемым к модулю автоматического лингвистического анализа текстового документа и модулю машинного перевода. Учитывая, что речь идет только о двух и довольно близких языках, а также тот факт, что СМП Б/Р имеет высокие качественные показатели, а поисковое пространство в основном является русскоязычным, целесообразно было исследовать и уточнить принципиальную схему решения задачи на этапе поиска релевантных документов (блок 2), что привело к разработке структурно-функциональной схемы соответствующей подсистемы. В соответствии с ней распознавание ключевых слов во входном документе основывается на использовании метода $TF*IDF$ [2], получаемые при этом белорусскоязычные документы из МПП, а также входной текстовый документ переводятся с помощью СМП Б/Р на русский язык непосредственно перед началом распознавания ЗФ (блок 3). В случае если поисковая БД пользователя является относительно постоянной, то такой перевод рекомендуется осуществить заранее для всех белорусскоязычных документов этой БД, что существенно

снижает общее время проведения экспертизы. Дополнительно в подсистеме поиска релевантных документов (блок 2) пользователю, т. е. эксперту, наряду с автоматическим построением поискового образа запроса в виде списка ключевых слов предоставляется возможность самому задать такой список. При реализации подсистемы используется поисковый механизм информационно-поисковой системы Google. Экспериментально установлено, что наиболее приемлемыми являются длина поискового образа запроса в три-пять ключевых слов и количество возвращаемых поисковой машиной релевантных документов не более 50. Для получения необходимых методом TF*IDF частот лексических единиц был построен корпус текстов объемом 976 382 слова. Он включает тексты различных предметных областей, полученных как из сети Интернет, так и из локальных коллекций рефератов, курсовых, дипломных, диссертационных работ и законодательных актов.

Проектирование и реализация подсистемы определения языка текстового документа (блок 6) осуществлялись в точном соответствии с методом, лингвистическими ресурсами и алгоритмом, представленными в [3], причем длина требуемого для этой цели фрагмента анализируемого текстового документа составляла не более одной страницы.

Проектирование и реализация подсистемы распознавания заимствованных фрагментов (блок 3) выполнены в полном соответствии с методом и алгоритмом, представленными в [4]. При этом было взято пороговое значение $\mu = 4$, определяющее максимально допустимое количество слов предложения из входного текста, не входящих в сравниваемое предложение из БД.

Что касается подсистемы формирования отчета (блок 4), то независимо от режима использования ИПК она формирует и предоставляет пользователю отчет о результатах автоматической экспертизы в виде сравнительной таблицы, в которой представлены эквивалентные фрагменты входного текстового документа (слева) и найденного в локальной БД или среди интернет-доступных документов (справа). Представление эквивалентных фрагментов осуществляется на языках оригинальных документов, а для уточнения контекста существует возможность перехода к текстам самих документов. Результаты работы сохраняются в соответствующих каталогах системы.

ИПК «ПлагиятКонтроль» предоставляет пользователю возможность в зависимости от типа поискового пространства, способов формирования поискового образа запроса и личных предпочтений осуществлять работу в различных режимах:

Интернет-авто. Система сравнивает входной документ с релевантными документами, которые она автоматически находит в сети Интернет.

Сайт-авто. Система сравнивает входной документ с релевантными документами, которые она автоматически находит в определяемом самим экспертом конкретном сайте сети Интернет.

Локальный-авто. Система сравнивает входной документ с релевантными документами, задаваемыми самим экспертом.

Интернет-интер. Система сравнивает входной документ с релевантными документами, которые она находит в сети Интернет на основании списка ключевых слов, задаваемого экспертом.

Сайт-интер. Система сравнивает входной документ с релевантными документами, которые она находит в определяемом самим экспертом конкретном сайте в сети Интернет на основании задаваемого им же списка ключевых слов.

Анализ результатов работы. ИПК «ПлагиятКонтроль» тестировался, например, на коллекции пояснительных записок к дипломным проектам студентов как гуманитарного, так и технического профилей в количестве 94 текстовых документов [5]. На рис. 2 показан фрагмент сравнительной таблицы совпадающих текстовых фрагментов документов: слева – фрагменты анализируемого документа «Товарная политика ОАО «ПродСервис» (на примере магазина «Корзинка»)», справа – фрагменты найденного релевантного ему документа «Совершенствование товарной стратегии магазина, дипломная работа», загруженного из коллекции рефератов на сайте <http://bibliofond.ru>.

C:\PigInternet\Волж.rtf	c:\PigInternet\Saved\13032013_123710\qpy4eudk.cwf.html
0.1 Даже хорошо продуманные планы сбыта и рекламы не смогут исправить ошибки, допущенные при планировании ассортимента.	1.1 Даже хорошо продуманные планы сбыта и рекламы не смогут нейтрализовать последствия ошибок, допущенных ранее при установлении цены.
0.1 Важным элементом товарной политики является служба сервиса для клиентов.	1.1 Одним из элементов товарной политики является создание службы сервиса для клиентов.
0.1 Успех на рынке – главный критерий оценки деятельности предприятия	1.1 Поскольку рыночный (конечный) успех отныне является главным критерием оценки деятельности предприятий, а их рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой, то именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием, планированием ассортимента и его совершенствованием [20, с 2].
0.1 Рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой.	1.1 Поскольку рыночный (конечный) успех отныне является главным критерием оценки деятельности предприятий, а их рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой, то именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием, планированием ассортимента и его совершенствованием [20, с 2].
0.1 Именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием ассортимента и управлением процессом реализации продукции.	1.1 Поскольку рыночный (конечный) успех отныне является главным критерием оценки деятельности предприятий, а их рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой, то именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием, планированием ассортимента и его совершенствованием [20, с 2].
0.1 Обеспечение необходимого уровня обслуживания покупателей и роста основных экономических показателей деятельности товарного предприятия в значительной степени зависит от правильного формирования ассортимента товаров.	1.1 Обеспечение необходимого уровня обслуживания покупателей и роста основных экономических показателей деятельности товарного предприятия в значительной степени зависит от правильного формирования ассортимента товаров в его магазинах.
0.1 Процесс формирования товарного ассортимента в магазине должен исходить из формы его товарной специализации и размера торговой	1.1 Процесс формирования товарного ассортимента в магазине должен исходить из формы его товарной

Рис. 2. Фрагмент сравнительной таблицы текстов анализируемых документов

В первой строке имеет место значительное сходство формулировок, смысл предложений достаточно близок, особенно их первых частей:

«*Даже хорошо продуманные планы сбыта и рекламы не смогут исправить ошибки, допущенные при планировании ассортимента.*» и

«*Даже хорошо продуманные планы сбыта и рекламы не смогут нейтрализовать последствия ошибок, допущенных ранее при установлении цены.*»

Здесь и далее в примерах полужирным шрифтом выделены различия в существующих парах лексем с точностью до символа, а прямое полужирное начертание указывает на лексемы, не имеющие пары в сравниваемых фрагментах, т. е. добавленные и (или) удаленные.

Вторая пара найденных фрагментов еще более близка по смыслу:

«**Важным** элементом товарной политики является служба сервиса для клиентов.» и

«**Одним из элементов** товарной политики является **создание службы сервиса** для клиентов.»

Кроме того, в предложениях первых двух пар анализируемых документов присутствуют морфологические преобразования и синонимические конструкции.

Третья, четвертая и пятая пары фрагментов демонстрируют разделение сложносочиненного предложения релевантного документа на три простых предложения входного документа с последующим добавлением (удалением) лексических единиц и выполнением процедур их грамматического согласования:

«*Успех на рынке – главный критерий оценки деятельности предприятия.*»

«*Рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой.*»

«Именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием ассортимента и управлением процессом реализации продукции.» и

«Поскольку рыночный (конечный) успех отныне является главным критерием оценки деятельности предприятий, а их рыночные возможности предопределяются правильно разработанной и последовательно осуществляемой товарной политикой, то именно на основе изучения рынка и перспектив его развития предприятие получает исходную информацию для решения вопросов, связанных с формированием, планированием ассортимента и его совершенствованием [20, с 2].»

Обнаружение указанных соответствий стало возможным благодаря тому, что разработанный алгоритм сравнения предложений [4] рассматривает их не как цепочки, а как множества слов, исключает высокочастотные общеупотребительные лексические единицы, которыми, как правило, являются предлоги, артикли, союзы и т. п., а само равенство предложений определяется в общем случае с точностью до канонических форм слов и отношений синонимии.

Шестая пара фрагментов отличается лишь отсутствием в предложении анализируемого документа трех лексем из релевантного ему предложения, причем их смысл почти одинаков:

«Обеспечение необходимого уровня обслуживания покупателей и роста основных экономических показателей деятельности товарного предприятия в значительной степени зависит от правильного формирования ассортимента товаров.» и

«Обеспечение необходимого уровня обслуживания покупателей и роста основных экономических показателей деятельности товарного предприятия в значительной степени зависит от правильного формирования ассортимента товаров в его магазинах.»

resumeBel.pdf	resumeRus.pdf
С. 1. Ключавыя словы: натуральная мова, лінгвістычны аналіз тэксту, алгарытм, канечны аўтамат, лінгвістычны працэсар, праблемаарыентаваная мова, корпус тэкстаў, лінгвістычная база ведаў, слоўнік, лінгвістычныя правілы, рэгулярны выраз.	С. 1. Ключевые слова: естественный язык, лингвистический анализ текста, алгоритм, конечный автомат, лингвистический процессор, проблемно-ориентированный язык, корпус текстов, лингвистическая база знаний, словарь, лингвистическое правило, регулярное выражение.
С. 1. Методы даследавання: метады камп'ютарнай лінгвістыкі, тэорыя формальных моў, аўтаматныя граматыкі, тэорыя канечных аўтаматаў і алгарытмаў, аб'ектнаарыентаванае праграмаванне.	С. 1. Методы исследования: методы компьютерной лингвистики, теория формальных языков, автоматные грамматики, теория конечных автоматов и алгоритмов, объектно-ориентированное программирование.
С. 1. У рабоце сфармулявана канцэпцыя базавага лінгвістычнага працэсара і вызначана яго функцыянальнасць, якая з'яўляецца ўніверсальнай у дадзеным да розных моў і задач іх апрацоўкі.	С. 1. В работе сформулирована концепция базового лингвистического процессора и определена его функциональность, которая является универсальной по отношению к различным ЕЯ и задачам их обработки.
С. 1. Распрацавана мова пашыраных рэгулярных выразаў WRE для формальнага апісання лінгвістычных правілаў, якая максімальна судаднесена з патрабаваннем яе даступнасці для выкарыстання экспертамі магчымасцю абагульнення патэрнаў і пабудовы эфектыўнага алгарытмічнага забеспячэння базавага лінгвістычнага працэсара.	С. 1. Разработан язык расширенных регулярных выражений WRE для формального описания лингвистических правил, который максимально соотносен с требованием его доступности для использования экспертами, возможностью обобщения разрабатываемых правил и построения эффективного алгоритмического обеспечения базового лингвистического процессора.
С. 1. Распрацаваны метады, алгарытмы, тэхналогіі і праграмае забеспячэнне ўкаранены ў складзе прамысловых сістэм GoldFire, TechOptimizer, Knowlegist, Cobrain, карыстальнікамі якіх з'яўляюцца многія найбуйнейшыя кампаніі свету, а таксама ўкаранены ў навучальны працэс у Беларускам дзяржаўным універсітэце.	С. 1. Степень использования: разработанные методы, алгоритмы, технологии и программное обеспечение внедрены в составе промышленных систем GoldFire, TechOptimizer, Knowlegist, Cobrain, пользователями которых являются многие крупнейшие компании мира, а также внедрены в учебный процесс в Белорусском государственном университете.

Рис. 3. Фрагмент сравнительной таблицы текстов анализируемых рефератов на белорусском и русском языках

Система «ПлагиаТКонтроль» тестировалась и успешно используется в Высшей аттестационной комиссии Республики Беларусь для экспертизы диссертационных работ с привлечением информационных ресурсов Национальной библиотеки Беларуси. На рис. 3 приведен еще один фрагмент сравнительной таблицы текстов анализируемых документов. Он интересен тем, что система в качестве документа, релевантного входному (диссертации), выбрала в том числе эту же диссертацию и, функционируя в режиме cross-language, распознала в качестве ЗФ для реферата этой диссертации на белорусском языке ее же реферат, но на русском языке. Очевид-

но, что приведенные здесь пары предложений очень близки по смыслу, а полученное решение стало возможным благодаря высокому качеству подсистемы МП.

В результате тестирования в 84 % работ был обнаружен хотя бы один эквивалентный фрагмент, а размер заимствованных фрагментов в среднем в расчете на один документ составил 6,15 % его объема, причем в 36,17 % работ размер заимствованных фрагментов был больше среднего. Детальный анализ сгенерированных системой отчетов о результатах автоматической экспертизы подтвердил факты заимствований без ссылок на авторов в объемах, превышающих, например, допустимые в БГУ (не более 10 % текста работы) в 25,53 % документов.

Что касается скорости работы системы, то в среднем время проведения процедуры автоматического анализа одного дипломного проекта (средний размер которого составлял 91 974 символа) с помощью ИПК на ЭВМ типа Intel Pentium Dual CPU T2370, 1,73 ГГц, ОЗУ 2 ГБ, под управлением ОС Windows XP варьировалось в диапазоне от 90 до 180 с, а время, необходимое эксперту для обработки полученных результатов с целью признания или непризнания распознанных системой заимствованных текстовых фрагментов плагиата, в среднем занимало от 20 до 60 мин.

Заключение. Полученные результаты позволяют считать, что разработанный ИПК «ПлагиатКонтроль» является эффективным средством автоматизации труда эксперта и существенно ускоряет процедуру анализа диссертационных, дипломных и курсовых работ на русском и белорусском языках на предмет наличия в них плагиата, совершенного путем заимствования фрагментов из других текстовых документов, представленных в локальных поисковых БД пользователя и в сети Интернет без указания на них ссылок.

Список использованных источников

1. Воронков, Н. В. Методы, алгоритмы и модели систем автоматического реферирования текстовых документов : автореф. дис. ... канд. техн. наук : 05.13.17 / Н. В. Воронков ; Бел. гос. ун-т. – Минск, 2007. – 22 с.
2. Robertson, S. Understanding Inverse Document Frequency: On Theoretical Arguments for IDF / S. Robertson // Journal of Documentation. – 2004. – №. 60(5). – P. 503–520.
3. Крапивин, Ю. Б. Автоматическое определение языка текстового документа для основных европейских языков / Ю. Б. Крапивин // Информатика. – 2011. – № 31. – С. 112–116.
4. Крапивин, Ю. Б. Автоматический поиск заимствованных из интернет-источников фрагментов / Ю. Б. Крапивин // Искусственный интеллект. – 2012. – № 4. – С. 183–189.
5. Крапивин, Ю. Б. Об использовании системы автоматического распознавания воспроизведенных фрагментов текстовых документов в учебном процессе / Ю. Б. Крапивин // Информационные технологии и системы 2013 (ИТС 2013) : материалы Междунар. науч. конф. – Минск, 2013. – С. 142–143.

Информация об авторе

Крапивин Юрий Борисович – старший преподаватель, Брестский государственный технический университет (ул. Московская, 267, Брест, Республика Беларусь). E-mail: ybox@list.ru

Information about the author

Yury B. Krapivin – Senior Lecturer, Brest State Technical University (267, Moscovskaya Ave., Brest, Republic of Belarus). E-mail: ybox@list.ru