

ҚАЗАҚСТАН РЕСПУБЛИКАСЫ
БІЛІМ ЖӘНЕ ҒЫЛЫМ МИНИСТРЛІГІ

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ
РЕСПУБЛИКИ КАЗАХСТАН



ҚазҰТЗУ ХАБАРШЫСЫ _____

_____ **ВЕСТНИК КазНУ**

VESTNIK KazNRTU _____

№3 (127)

Главный редактор
И. К. Бейсембетов – ректор

Зам. главного редактора
Б.К. Кенжалиев – проректор по науке

Отв. секретарь
Н.Ф. Федосенко

Редакционная коллегия:

С.Б. Абдыгаппарова, Б.С. Ахметов, З.С. Абишева- акад. НАНРК, Л.Б. Атымтаева, Ж.Ж. Байгунчечков- акад. НАНРК, А.Б. Байбатша, А.О. Байконурова, В.И. Волчихин (Россия), К. Дребенштед (Германия), Г.Ж. Жолтаев, Р.М. Искаков, С.Е. Кудайбергенов, С.Е. Кумеков, В.А. Луганов, С.С. Набойченко – член-корр. РАН, И.Г. Милев (Германия), С. Пежовник (Словения), Б.Р. Ракишев – акад. НАН РК, М.Б. Панфилов (Франция), Н.Т. Сайлаубеков, А.Р. Сейткулов, Фатхи Хабаши (Канада), Бражендра Мишра (США), Корби Андерсон (США), В.А. Гольцев (Россия), В. Ю. Коровин (Украина), М.Г. Мустафин (Россия), Фан Хуаан (Швеция), Х.П. Цинке (Германия), Т.А. Чепуштанова, Г.Ж. Елигбаева, Б.У. Куспангалиев

Учредитель:

Казахский национальный исследовательский технический университет
имени К.И. Сатпаева

Регистрация:

Министерство культуры, информации и общественного согласия
Республики Казахстан № 951 – Ж “25” 11. 1999 г.

Основан в августе 1994 г. Выходит 6 раз в год

Адрес редакции:

г. Алматы, ул. Сатпаева, 22,
каб. 616, тел. 292-63-46
Nina. Fedorovna. 52 @ mail.ru

- [13] Петрова Е.В. Взгляд в будущее – обзор новинок, представленных на выставке ConExpo-Con / Agg. // Строительная техника и технологии. – 2014. – №3 (103). – С. 30-48.
- [14] Клушанцев Б.В., Косарев А.И., Муйземнек Ю.А. Дробилки. Конструкция, расчет, особенности эксплуатации. – М.: Машиностроение, 1990. – С. 320.
- [15] Предварительный патент РК №19801, В02С 4/30, бюл. №8, 15.08.2008г.

Гурьянов Г.А., Ким В.А., Васильева О.Ю.

Қатты бөлшектерді бұзу процесінің механика-математикалық моделін әзірлеу және кедір-бұдырлы біліктердің геометриялық параметрлерін анықтау

Аңдатпа: Мақалада кедір-бұдырлы біліктердің конструкциясын біліктердің геометриялық профилі пішінін өзгерту арқылы жетілдіру, сондай-ақ кесек материалды жаңа профилді біліктің қармау процесінің компьютерлік моделі нәтижесі көрсетілген.

Түйінді сөздер: кесек материал, қармау бұрышы, ұнтақтау процесі, білікті уатқыш, компьютерлік модель.

Guryanov G.A., Kim V.A., Vasileva O.Yu.

Development of the mechanical-mathematical model of the process of the destruction of solid particles and definition of geometric parameters of convex-concave rollers

Summary: The article considers a variant of improving the design of a roller crusher by changing the geometric shape of the profile of rolls, as well as the results of computer simulation of the process of gripping a piece of material in a roller crusher with a new roll profile.

Keywords: lump material, gripping angle, grinding process, roller crusher, computer model.

M. Mansurova, M. Kaipoldayev

(Al-Farabi Kazakh National University, Almaty, Kazakhstan
E-mail: mansurova.madina@gmail.com, m.kaipoldayev@gmail.com)

DEVELOPMENT OF A MODULE FOR DETECTING DUPLICATE TEXTS

Abstract: This paper is devoted to solve the problem of eliminating duplicating articles in the news web-portal. A model of the system for detecting and eliminating such articles is described. Also, two algorithms for comparing texts are given and an analysis of their effectiveness is made. The main task was to develop a module for eliminating duplicate articles for the news web-portal.

Keywords: w-shingling algorithm, Hirschberg's algorithm, duplication of texts, text processing.

М.Е. Мансурова, М.Е.Қайполдаев

(Казахский Национальный Университет им. Аль-Фараби, Алматы, Республика Казахстан.
E-mail: mansurova.madina@gmail.com, m.kaipoldayev@gmail.com)

РАЗРАБОТКА МОДУЛЯ ОБНАРУЖЕНИЯ ДУБЛИРУЮЩИХСЯ ТЕКСТОВ

Аннотация: Данная работа посвящена решению задачи устранения дублирующихся статей в новостном портале. Описана модель системы обнаружения и устранения таких статей. Также приведены два алгоритма сравнения текстов и сделан анализ их эффективности. Основной задачей являлась разработка модуля устранения дублирующихся статей для новостного портала.

Ключевые слова: алгоритм Шинглов, алгоритм Хиршберга, дублирование текстов, обработка текста.

1. Введение

Проблема повторяющегося содержимого текста очень становится очень актуальной в современных реалиях. Если содержимое различных текстов описывает одно и то же событие или объект, то это может привести к различным негативным последствиям. К примеру в поисковых системах необходимо выводить данные содержащие различный контент, так-как пользователю не желательно просматривать множество схожих сайтов. Также любой документ претендующий на уникальность должен пройти процедуру выявления плагиата. В таких ситуациях необходима разработка систем и алгоритмов, которые выявляют сходства содержимого текстов и устраняет дублирование документов.

Значительной проблемой анализа данных на дублирование является то, что с развитием интернет технологий количество веб-документов очень сильно возросло, что в свою очередь привело к увеличению числа схожих документов. Данный фактор привел к необходимости разработки более эффективных методов выявления дублирования.

2 Постановка задачи

На данный момент в лаборатории НИИ ММ КазНУ им. аль-Фараби разработан веб-портал который занимается сбором и анализом информации о чрезвычайных ситуациях на территории Казахстана. Данный сайт собирает и публикует новостные данные с официальных источников МВД РК. Однако, вследствие того что многие новостные сайты публикуют схожие новости, возникает проблема дублирующихся новостей.

Для обеспечения оригинальности публикуемых новостей необходимо разработать и добавить модуль фильтрации новых новостей и оставлять лишь один из многих дублирующихся новостей (рис. 1). Предварительно все собранные за небольшой промежуток времени новости будут храниться в кэш-памяти системы. Затем модуль выявления дублирующихся записей будет обрабатывать данные кэша и удалять записи которые удовлетворяют условиям дублирования статей. После чего на сайте публикуются лишь оригинальные новости. Эти же новости отправляются на хранение в базу данных новостей в качестве архива. На портале будет публиковаться лишь самая первая по дате публикации новость из всех одинаковых записей.

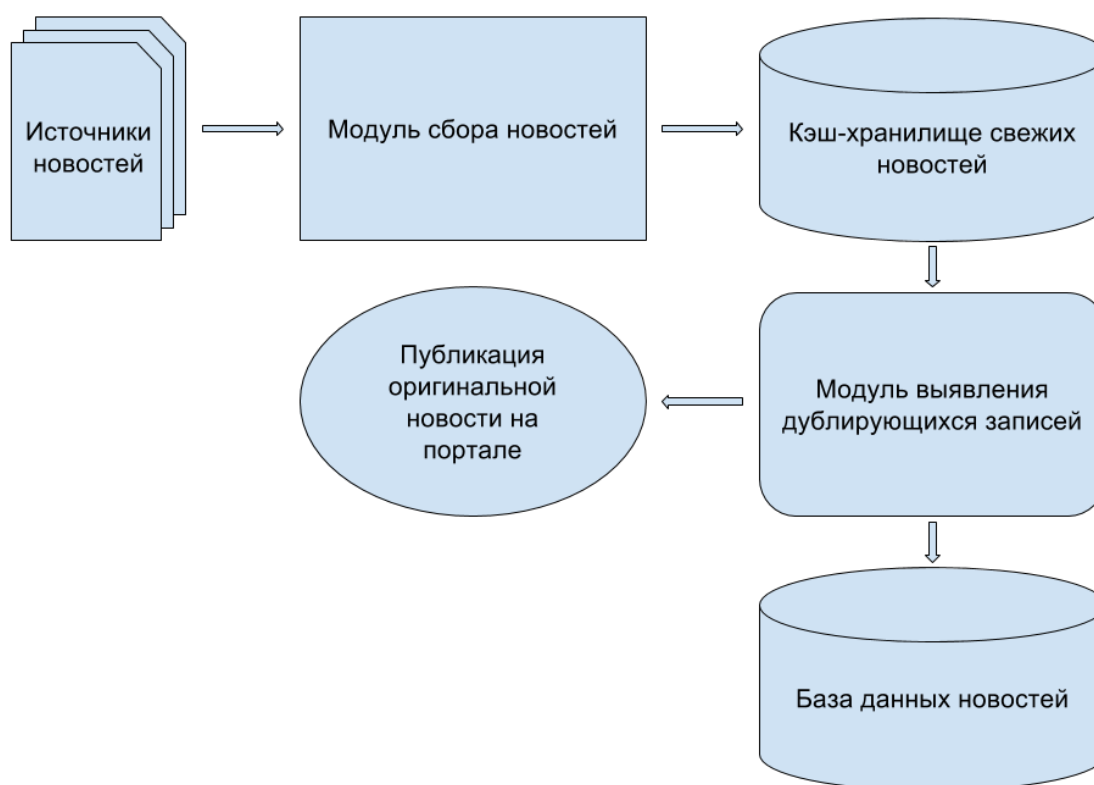


Рис. 1. Модель системы сбора и публикации новостей

3 Существующие решения

Для определения схожести текста вычисляются различные метрики. Одним из часто используемых мер схожести является расстояние Левинштейна[1]. Данное расстояние определяется как наименьшее количество операции удаления и вставки символа с помощью которых из одного текста можно получить второе. Имеются также модификации данного алгоритма которые учитывают вес операции. К примеру, операции удаления и вставки символа имеют вес 1, а операция замены символа имеют вес 2. Данная метрика называется “расстояние редактирования”.

Выше приведенные метрики применимы для сравнения строк, однако для сравнения больших текстов они не подходят, так-как эти метрики вычисляются с помощью сравнения символов. Для сравнения текстов объектом сравнения необходимо брать целые слова. В связи с чем широко распространен алгоритм шинглов разработанный А. Broder et al[2]. В данном методе сравниваемые тексты делятся на множество частей (шинглов), которые попарно сравниваются друг с другом. Дальнейшим развитием данного метода является такой алгоритм, при котором вычисляются 84 контрольных сумм каждого документа. Далее эти 84 шингла образуют 6 супершинглов. Затем каждый документ представляется всеми парными сочетаниями из 6 супершинглов, эти сочетания называются мегашинглами. Два текста будут считаться схожими если у них совпадает не менее одного мегашингла.

Другой подход основанный на лексической характеристике текста создан А. Chowdhury et al. Этот алгоритм основывается на вычислении дактилограммы I-Match. Для этого необходимо составить словарь со средними значениями IDF (*inverse document frequency*). После этого находится пересечение множеств слов из составленного словаря и всех слов текста. Затем, если значение этого пересечения больше заданного порога, то вычисляется дактилограмма I-Match. А два документа можно считать схожими по содержанию если данные дактилограммы идентичны.

Существуют подходы основывающиеся на сходстве фонетического звучания слова. На данный момент большинство фонетических алгоритмов разработаны для английского языка, следовательно не могут корректно применяться для языков не схожих с данным. Наиболее популярными фонетическими алгоритмами являются Soundex, New York State Identification and Intelligence System, Cologne phonetics и т.д.

4 Выбор алгоритма и его реализация

Для решения задачи нахождения схожести текстов нами были выбраны два алгоритма: алгоритм w-shingling и алгоритм Хиршберга[3]. Выбор данных методов обусловлен их эффективностью и относительной простотой реализации.

4.1 Алгоритм w-shingling

W-shingling - это алгоритм оценки сходства двух документов. Суть данного алгоритма заключается в нахождении количества шинглов принадлежащих обоим документам. Затем находится отношение общих шинглов двух текстов к общему количеству шинглов в двух текстах.

Шингл - это последовательность, фиксированной длины, подряд идущих слов в тексте. Например, рассмотрим следующее выражение: “Мороз и солнце; день чудесный!”. Разделим данное предложение на шинглы с длиной 3. Для наглядности не будем учитывать знаки препинания. Получаем: “Мороз и солнце”, “и солнце день”, “солнце день чудесный”. Несложно догадаться, что количество шинглов равно

количество слов в тексте - длина шингла + 1.

Перед применением данного алгоритма необходимо провести предварительную обработку текста, т.е. удалить все знаки препинания, удалить стоп-слова (слова не имеющие смысловой нагрузки), привести слова к стандартной форме (именительный падеж, единственное число и т. п.).

Предобработка текста для применения алгоритма w-shingling

- Удаление всех лишних символов. Для осуществления данной операции была написана функция, которая удаляет все символы кроме А-Z, а-z, А-Я, а-я.

- Удаление стоп-слов. Для этого действия применялась функция которая сравнивает слова текста со стоп-словами, которые хранились в файле где записаны все слова такого типа. При совпадении из текста удалялось данное слово.

- Приведение слов к нормальной форме. Для этого применялся лемматизатор Apache Lucene.

Для нахождения дактилограммы шинглов была применена хэш-функция CRC32. Дактилограмма необходима для упрощения процесса сравнения шинглов, так как гораздо быстрее сравнивать числовые значения хэш-функций, нежели сравнивать шинглы в виде строк.

При написании программы пришлось совместить деление слов на шинглы и нахождение дактилограммы (т.е. контрольной суммы), так как имелась предстала следующая проблема. У шинглов, с одинаковым набором слов но расположенных в разном порядке, дактилограмма должна быть одинаковой. К примеру шинглы “солнце встает над Алматы” и “над Алматы встает солнце” должны иметь одинаковое значение дактилограммы.

<i>Amangeldiyev S., Maksimov V.</i>	
IMPLEMENTATION OF IT TECHNOLOGIES IN THE TPP OF KAZAKHSTAN.....	381
<i>Kulikov V.Yu., Kyon Sv.S., Issagulov A.Z., Chsherbakova Y.P., Kovaleva T.V.</i>	
DEVELOPMENT OF THE TECHNOLOGY OF MANUFACTURING SHELL MOLDS WITH UNIFORM HARDNESS THROUGHOUT THE VOLUME.....	388
<i>Ashirbayeva N., Duisebayeva P., Ashirbayeva Zh., Alibekova</i>	
DYNAMIC FIELDS OF STRESS IN ELASTIC BODY WITH Zh. FOREIGN INCLUSION.....	393
<i>Otynshiyeva AM, Musapirova G.D.</i>	
RESEARCH METHODS OF SITE OPTIMIZATION.....	399
<i>Kartbayev A.ZH.</i>	
A LEARNING OF THE KAZAKH LANGUAGE MODEL INTERPOLATION PARAMETERS	406
<i>Alzhanova A. Ye., Azmukhanov A.A.</i>	
NANOMECHANICAL PROPERTIES OF ION IRRADIATED SiO ₂ /Si SYSTEMS.	410
<i>Zharkevich O., Nurzhanova O., Mateshov A.</i>	
OPTIMIZATION OF CONSTRUCTIVE PARAMETERS FOR HYDRAULIC CYLINDERS OF POWERED SUPPORT.....	415
<i>Tergeussizova A.S.</i>	
AUTOMATIC MANAGEMENT OF TECHNOLOGICAL PRODUCTION AND MATHEMATICAL MODELING OF STABILITY OF THE EXHAUST OF THE OPTICAL FIBER EXHAUST.....	419
<i>Serikbayeva A.K., Sameshova A.K.</i>	
STUDY OF PHASE TRANSFORMATIONS IN THE "PbO-S" SYSTEM.....	424
<i>Almuratova N.K., Toigozhinova Zh. Zh.</i>	
SUSTAINABILITY OF THE DYNAMICS OF NONLINEAR SYSTEM «FREQUENCY CONVERTER - ASYNCHRONOUS ENGINE».....	430
<i>Sagitov P.I., Asanova K.S., Toigozhinova Zh. Zh.</i>	
OTIMIZATION OF ENERGY SAVING IN A REGULATED SYSTEM FREQUENCY CONVERTER ASYNCHRONOUS ENGINE.....	435
<i>Tergeussizova A.S.</i>	
MANAGEMENT SYSTEMS FOR THE OPTICAL FIBER EXHAUST AND INNOVATIVE TECHNOLOGIES FOR ITS PRODUCTION.....	439

Physico-mathematical sciences

<i>Nurseiit S.T., Sariyeva A.K., Danlybaeva A.K.</i>	
ANALYSIS OF DOMESTIK AND FOREIGN SYSTEMS OF ENVIROMENTAL CERTIFICATION OF CONSTRUCTION SITES.....	448
<i>Baimatova N</i>	
INVESTIGATION OF EFFICIENCY OF BENZENE, TOLUENE, ETHYLBENZENE AND O-XYLENE ADSORPTION FROM INDOOR AIR BY MODIFIED CARBON-BASED ADSORBENTS.....	453
<i>Mirzakhmedova G.A.</i>	
MANAGEMENT OF ECONOMIC CLUSTERS.....	460
<i>Omar A., Amantayeva A., Bissekenova A., Agishev A., Hohlov S.</i>	
INFORMATION ENTROPY SPECTRA OF HOT STARS.....	465
<i>Rogovoy A.V., Karasheva K.</i>	
PROPERTIES OF ONE CLASS SOLUTIONS OF TRICOMI PROBLEM FOR MIXED TYPE EQUATIONS	469
<i>Turdykhanova D.L., Berikov D.B., Zhumadilov K.Sh.</i>	
STATISTICAL ANALYSIS OF THE HOMOGENEITY OF THE NUCLIDE COMPOSITION IN DIFFERENT AREAS.....	478
<i>Sakypbekova M.Zh.</i>	
NUMERICAL SOLUTION OF THE BASIC EQUATIONS OF HYDRODYNAMICS IN THE SIMULATION OF A TWO-DIMENSIONAL FLOW.....	482
<i>Guryanov G.A., Kim V.A., Vasileva O.Yu.</i>	
DEVELOPMENT OF THE MECHANICAL-MATHEMATICAL MODEL OF THE PROCESS OF THE DESTRUCTION OF SOLID PARTICLES AND DEFINITION OF GEOMETRIC PARAMETERS OF CONVEX-CONCAVE ROLLERS.....	486
<i>Mansurova M., Kaipoldayev M.</i>	
DEVELOPMENT OF A MODULE FOR DETECTING DUPLICATE TEXTS.....	495