

ISSN 2708-2032
e-ISSN 2708-2040



**INTERNATIONAL
UNIVERSITY**

**INTERNATIONAL
JOURNAL OF INFORMATION
& COMMUNICATION TECHNOLOGIES**

**Volume 2, Issue 2
June, 2021**

ҚАЗАҚСТАН РЕСПУБЛИКАСЫНЫҢ БІЛІМ ЖӘНЕ ҒЫЛЫМ МИНИСТРЛІГІ
МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РЕСПУБЛИКИ КАЗАХСТАН
MINISTRY OF EDUCATION AND SCIENCE OF THE REPUBLIC OF KAZAKHSTAN



**INTERNATIONAL JOURNAL OF
INFORMATION AND COMMUNICATION
TECHNOLOGIES**

**МЕЖДУНАРОДНЫЙ ЖУРНАЛ
ИНФОРМАЦИОННЫХ И
КОММУНИКАЦИОННЫХ ТЕХНОЛОГИЙ**

**ХАЛЫҚАРАЛЫҚ АҚПАРАТТЫҚ ЖӘНЕ
КОММУНИКАЦИЯЛЫҚ
ТЕХНОЛОГИЯЛАР ЖУРНАЛЫ**

Том 2, Выпуск 2
Июнь, 2021

Главный редактор – Ректор АО МУИТ, профессор, д.т.н.
Ускенбаева Р.К.

Заместитель главного редактора – Проректор по НиМД, PhD, ассоц.профессор
Дайнеко Е.А.

Отв. секретарь – PhD, ассоц.профессор, директор департамента по науке
Кальпеева Ж.Б.

ЧЛЕНЫ РЕДКОЛЛЕГИИ:

Отельбаев М. д.т.н., профессор, АО «МУИТ», Рысбайулы Б., д.т.н., профессор, АО «МУИТ», Куандыков А.А., д.т.н., профессор, АО «МУИТ», Синчев Б.К., д.т.н., профессор, АО «МУИТ», Дузбаев Н.Т., PhD, проректор по ЦИИ, АО «МУИТ», Ыдырыс А., PhD, заведующая кафедрой «МКМ», АО «МУИТ», Касымова А.Б., PhD, заведующая кафедрой «ИС», АО «МУИТ», Шильдибеков Е.Ж., PhD, заведующий кафедрой «ЭиБ», АО «МУИТ», Ипалакова М.Т., к.т.н., ассоц. профессор, заведующая кафедрой «КИИБ», АО «МУИТ», Айтмагамбетов А.З., к.т.н., профессор, АО «МУИТ», Амиргалиева С.Н., д.т.н., профессор, АО «МУИТ», Ниязгулова А.А., к.ф.н., заведующая кафедрой «МииК», АО «МУИТ», Молдагулова А.Н., к.т.н., ассоциированный профессор, АО «МУИТ», Джоламанова Б.Д., ассоциированный профессор, АО «МУИТ», Prof. Young Im Cho, PhD, Gachon University, South Korea, Prof. Michele Pagano, PhD, University of Pisa, Italy, Tadeusz Wallas, Ph.D., D.Litt., Adam Mickiewicz University in Poznań, Тихвинский В.О., д.э.н., профессор, МГУСИ, Россия, Масалович А., к.ф.-м.н., Президент Консорциума Инфорус, Россия, Lucio Tommaso De Paolis is the Research Director of the Augmented and Virtual Laboratory (AVR Lab) of the Department of Engineering for Innovation, University of Salento and the Responsible of the research group on “Advanced Virtual Reality Application in Medicine” of the DREAM, a multidisciplinary research laboratory of the Hospital of Lecce (Italy), Liz Bacon, Professor, Deputy Principal and Deputy Vice-Chancellor, Abertay University (Great Britain).

Издание зарегистрировано Министерством информации и общественного развития Республики Казахстан. Свидетельство о постановке на учет № KZ82VPY00020475 от 20.02.2020 г.

Журнал зарегистрирован в Международном центре по регистрации сериальных изданий ISSN (ЮНЕСКО, г. Париж, Франция)

Выходит 4 раза в год.

УЧРЕДИТЕЛЬ:

АО «Международный университет информационных технологий»

ISSN 2708-2032 (print)
ISSN 2708-2040 (online)

СОДЕРЖАНИЕ

РАЗРАБОТКА ПРОГРАММНОГО ОБЕСПЕЧЕНИЯ И ИНЖЕНЕРИЯ ЗНАНИЙ

Бактаев А.Б., Мукажанов Н.К.

Алгоритм решения задачи по исправлению опечаток в тексте, применяемый в поисковых системах с поддержкой казахского языка 9

Еркетаев Н.М., Мукажанов Н.К.

Эффективное хранение неструктурированных данных 19

Сагадиев Р.Т., Шайкемелев Г.Т.

Представление логической витрины данных в экосистеме Hadoop 28

Бейсенбек Е.Б., Дузбаев Н.Т.

Современные способы взлома и защиты ПО 33

Найзабаева Л.К., Алашымбаев Б.А.

Рекомендательная система для онлайн-магазинов с использованием машинного обучения 38

Мейрамбайулы Н., Дузбаев Н.Т.

Мониторинг стационарных источников выбросов загрязняющих веществ г. Алматы 47

ИНФОКОММУНИКАЦИОННЫЕ СЕТИ И КИБЕРБЕЗОПАСНОСТЬ

Айтмагамбетов А.З., Кулакаева А.Е., Койшыбай С.С., Жолшибек И.Ж.

Исследование возможностей применения низкоорбитальных спутников для радиомониторинга в республике Казахстан 54

Кемельбеков Б.Ж., Полуанов М.

Анализ метода бриллюэновской рефлектометрии в волоконно-оптических линиях связи ... 62

Турбекова К.Ж.

Анализ применения БПЛА в сетях связи при чрезвычайных ситуациях 68

ИНТЕЛЛЕКТУАЛЬНЫЕ СИСТЕМЫ

Азанов Н.П., Хабиров Р.Р., Әміров У.Е.

Конкурентная разведка и принятие решений с помощью машинного обучения для обеспечения промышленной безопасности 75

Джаныбекова С.Т., Толганбаева Г.А., Сарсембаев А.А.

Распознавание говорящего с помощью глубокого обучения 85

Салерова Д.К., Сарсембаев А.А.

Обзорная статья распознавания номерных знаков с использованием оптического распознавания символов 93

Салерова Д.К., Сарсембаев А.А.

Исследование существующих методов классификации изображений 100

Оразалин А., Мурсалиев Д.Е., Сергазина А.С.

Актуальные сверточные архитектуры нейронной сети для диагностики медицинских изображений 115

Әлімхан А.М.

Прогнозирование результатов игры в баскетбол с использованием алгоритмов глубокого обучения 112

<i>Адырбек Ж.А., Сатыбалдиева Р.Ж.</i> Анализ процессов планирования и решения проблем в логистике с помощью интеллектуальной системы	120
<i>Нурғалиев М.К., Алимжанова Л.М.</i> Геймификация в образовании	128

ЦИФРОВЫЕ ТЕХНОЛОГИИ В ЭКОНОМИКЕ И МЕНЕДЖМЕНТЕ

<i>Алимжанова Л.М., Панарина А.В.</i> Внедрение сервисной системы IT-аутсорсинга	133
<i>Жұмабай Р.Ж., Алимжанова Л.М.</i> Управление процессами работы с поставщиками на основе ERP-стандартов — подход BPM	140
<i>Бердыкулова Г.М., Төлепбергенова Д.А.</i> Менеджмент университета: практика МУИТ	146
<i>Омарова А.Ш., Алимжанова Л.М., Таштамышева А.Э.</i> Исследование и разработка методов перехода традиционного маркетинга в цифровой формат	153

CONTENTS

SOFTWARE DEVELOPMENT AND KNOWLEDGE ENGINEERING

<i>Baktayev A.B., Mukazhanov N.K.</i> Algorithm for solving the problem of correcting typos with search engines supporting the Kazakh language	9
<i>Yerketayev N.M., Mukazhanov N.K.</i> Efficient storage of unstructured data	19
<i>Sagadiyev R.T., Shaikemelev G.T.</i> Representing a logical data mart in the Hadoop ecosystem	28
<i>Beisenbek Y.B., Duzbaev N.T.</i> Modern methods of hacking and protection software	33
<i>Naizabayeva L., Alashybayev B.A.</i> A recommendation system for online stores using machine learning	38
<i>Meirambaiuly N., Duzbaev N.T.</i> Monitoring of stationary sources of pollutant emissions in Almaty	47

INFORMATION AND COMMUNICATION NETWORKS AND CYBERSECURITY

<i>Aitmagambetov A.Z., Kulakayeva A.E., Koishybai S.S., Zholshibek I.Z.</i> Study of the possibilities of using low-orbit satellites for radio monitoring in the Republic of Kazakhstan	54
<i>Kemelbekov B.J., Poluanov M.</i> Analysis of the brillouin reflectometry method in fiber-optic communication lines	62
<i>Turbekova K.Zh.</i> Analysis of the use of UAVs in emergency communication networks	68

SMART SYSTEMS

<i>Azanov N.P., Khabirov R.R., Amirov U.E.</i> Competitive intelligence and decision-making algorithm using machine learning for industrial security	75
<i>Janybekova S.T., Tolganbayeva G.A., Sarsembayev A.A.</i> Speaker recognition using deep learning	85
<i>Salerova D.K., Sarsembayev A.A.</i> Review of license plate recognition using optical character recognition	93
<i>Salerova D.K., Sarsembayev A.A.</i> Research on the existing image classification methods	100
<i>Orazalin A., Mursaliyev D.E., Sergazina A.S.</i> Current convolutional neural network architectures for diagnosing medical images.....	105
<i>Alimkhan A.M.</i> Predicting basketball results using deep learning algorithms	112
<i>Adyrbek Zh.A., Satybaldiyeva R.Zh.</i> Analysis of the planning and problem-solving processes in logistics using an intelligent system	120
<i>Nurgaliyev M.K., Alimzhanova L.M.</i> Gamification in education	128

DIGITAL TECHNOLOGIES IN ECONOMICS AND MANAGEMENT

Alimzhanova L.M., Panarina A.V.

Implementation of an IT outsourcing service system 133

Zhumabay R.Zh., Alimzhanova L.M.

Supplier process management based on ERP standards: the BPM approach 140

Berdykulova G.M., Tolepbergenova D.A.

University management: case study of IITU 146

Omarova A.Sh., Alimzhanova L.M., Tashtamysheva A.E.

Research and development of methods for the transition of traditional marketing to digital
format 153

МАЗМҰНЫ

БАҒДАРЛАМАЛЫҚ ҚАМТАМАНЫ ӨЗІРЛЕУ ЖӘНЕ БІЛІМ ИНЖЕНЕРИЯСЫ

Бактаев А.Б., Мукажанов Н.К.

Қазақ тілін қолдайтын іздеу жүйелерінде қолданылатын мәтіндегі жаңылыстарды түзету бойынша есептерді шешу алгоритмі..... 9

Еркетаев Н.М., Мукажанов Н.К.

Құрылымсыз деректерді тиімді сақтау 19

Сагадиев Р.Т., Шайкемелев Г.Т.

Надоор экожүйесінде логикалық деректер кесіндісін ұсыну 28

Бейсенбек Е.Б., Дузбаев Н.Т.

Бағдарламалық жасақтаманы бұзудың және қорғаудың заманауи әдістері 33

Найзабаева Л., Алашыбаев Б.А.

Машиналық оқытуды қолдану арқылы интернет-дүкендерге арналған ұсыныс жүйесі 38

Мейрамбайұлы Н., Дузбаев Н.Т.

Алматы қаласы бойынша ластаушы заттар шығарындыларының стационарлық дереккөздеріне мониторинг жүргізу 47

АҚПАРАТТЫҚ ЖӘНЕ КОММУНИКАЦИЯЛЫҚ ЖЕЛІЛЕР ЖӘНЕ КИБЕРҚАУПСІЗДІК

Айтмагамбетов А.З., Қулакаева А.Е., Койшыбай С.С., Жолшибек И.Ж.

Қазақстан Республикасында радиомониторинг үшін төмен орбиталық спутниктерді қолдану мүмкіндіктерін зерттеу 54

Кемельбеков Б.Ж., Полуанов М.

Талшықты-оптикалық байланыс желілеріндегі бриллюэн рефлектометрия әдісін талдау ... 62

Турбекова К.Ж.

Төтенше жағдайлар кезінде байланыс желілерінде ПҰА-ның қолданылуын талдау 68

ИНТЕЛЛЕКТУАЛДЫ ЖҮЙЕЛЕР

Азанов Н.П., Хабиров Р.Р., Әміров У.Е.

Өнеркәсіптік қауіпсіздікті қамтамасыз ету үшін машиналық оқытуды қолдана отырып, бәсекеге қабілеттілікті барлау және шешім қабылдау 75

Джаныбекова С.Т., Толғанбаева Г.А., Сарсембаев А.А.

Терең оқыту арқылы сөйлеушіні тану 85

Салерова Д.К., Сарсембаев А.А.

Таңбаларды оптикалық тануды пайдалану арқылы нөмірлер белгілерін тануға шолу мақаласы 93

Салерова Д.К., Сарсембаев А.А.

Қолданыстағы бейнелерді жіктеу әдістерін зерттеу 100

Оразалин А., Мурсалиев Д.Е., Сергазина А.С.

Медициналық кейіндік диагностикаға арналған конволюциялық жүйкелік желі архитектурасы 105

Әлімхан А.М.

Терең оқыту алгоритмдерін қолдана отырып, баскетбол нәтижелерін болжау 112

<i>Адырбек Ж.А., Сатыбалдиева Р.Ж.</i> Логистикадағы жоспарлау процестерін талдау және логистикадағы интеллектуалды жүйені қолдану арқылы мәселелерді шешу	120
<i>Нұрғалиев М.Қ., Алимжанова Л.М.</i> Білім беру саласындағы геймификация	128

ЭКОНОМИКА ЖӘНЕ БАСҚАРУДАҒЫ САНДЫҚ ТЕХНОЛОГИЯЛАР

<i>Алимжанова Л.М., Панарина А.В.</i> IT-аутсорсингтің сервистік жүйесін енгізу	133
<i>Жұмабай Р.Ж., Алимжанова Л.М.</i> ERP стандарттарына негізделген жеткізушілермен жұмыс процесін басқару - BPM тәсілі	140
<i>Бердыкулова Г.М., Төлепбергенова Д.А.</i> Университетті басқару: ХАТУ практикасы	146
<i>Омарова А.Ш., Алимжанова Л.М., Таштамышева А.Э.</i> Дәстүрлі маркетингті цифрлық форматқа ауыстыру әдістерін зерттеу және әзірлеу	153

Еркетаев Н.М., Мукажанов Н.К.*

Международный университет информационных технологий, Алматы, Казахстан

ЭФФЕКТИВНОЕ ХРАНЕНИЕ НЕСТРУКТУРИРОВАННЫХ ДАННЫХ

Аннотация. В этой статье представлены результаты изучения того, как могут быть использованы различные типы информации для хранения неструктурированных данных в базе данных и классической файловой системе. Также объясняются современные методы хранения данных, предыстория проблемы и ключевые элементы, важные для этого исследования: различные методы хранения неструктурированных данных (преимущества и недостатки), используемые типы данных и экспериментальная среда. Наша основная цель — найти эффективный метод хранения неструктурированных данных.

Ключевые слова: база данных, неструктурированные данные, файловый поток, эффективность, производительность

Введение

Хорошо известно, что одним из результатов быстрого роста Интернета стало значительное увеличение объема информации, генерируемой и распространяемой организациями практически во всех отраслях и секторах. Проблема не только в генерации данных, но и в их хранении и доступе к ним. Менее известна степень потребности в больших объемах дорогостоящих ресурсов, как человеческих, так и технических, обусловленной происходящим информационным взрывом. Эти потребности, в свою очередь, создали столь же большую, но в значительной степени неудовлетворенную потребность в инструментах, которые можно использовать для управления тем, что мы называем неструктурированными данными. Следует признать, что термин неструктурированные данные может означать разные вещи в разных контекстах. Например, в контексте систем реляционных баз данных это относится к данным, которые нельзя хранить в строках и столбцах. Вместо этого такие данные должны храниться в BLOB (большом двоичном объекте), универсальном типе данных, доступном в большинстве программных систем управления реляционными базами данных (RDBMS). Здесь под неструктурированными данными понимаются файлы электронной почты, текстовые документы, презентации, файлы изображений и видеофайлы.

С другой стороны, стремительный рост объемов данных обусловлен достижениями и серьезными преобразованиями в технологии магнитной записи. Стоимость хранения резко упала с более чем 4 долларов за мегабайт в 1990 году до менее чем 0,001 доллара за мегабайт сегодня.

Согласно исследованию, проведенному IDC, ведущей фирмой по исследованию и анализу рынка информационных технологий, объем данных, которые будут собираться, храниться и воспроизводиться по всему миру, вырастет со 161 эксабайта в 2006 году до 988 эксабайт в 2010 году (1 эксабайт = 1018 байт). Два ключевых вывода этого исследования:

- Основная часть этих данных будет в виде изображений, снятых большим количеством устройств, таких как цифровые камеры, телефоны с камерами, камеры наблюдения и медицинское оборудование для визуализации. Большая часть этих данных должна храниться и управляться централизованными системами внутри организаций. Исследование показывает, что к 2010 году, хотя предприятия будут создавать, собирать и тиражировать только 30% цифровой вселенной [3], им придется хранить более 85% всех данных и управлять ими.

- Более 95% цифровой вселенной — это неструктурированные данные. Согласно этому исследованию, 80% всех хранимых организациями данных не структурированы.

Ожидается, что эта тенденция роста сохранится и в будущем, поскольку потребует эффективных способов хранения, поиска, структурирования и обеспечения безопасности неструктурированных данных. Об аналогичных тенденциях в отношении важности обработки неструктурированных данных также сообщили другие исследовательские и консультационные фирмы, такие как Gartner Group и Butler Group [4].

Очевидно, что неструктурированные данные — это наша реальность, и поиск эффективного метода хранения данных является ключевым аспектом этого исследования и будущих результатов в этой области. В настоящей статье объясняются современные методы хранения данных.

Связанная работа

Существует множество исследований, основанных на тестировании систем баз данных. Одна статья, представляющая собой хорошую отправную точку для нашей работы, особенно интересна: «Набор тестов для обработки неструктурированных данных» [2].

Основной целью авторов было собрать набор рабочих нагрузок, выявить и изучить широкий спектр приложений для обработки неструктурированных данных, выявить их ключевые характеристики обработки и доступа к вводу-выводу, а также составить рабочие нагрузки, воплощающие эти характеристики. Поэтому, чтобы спроектировать системы хранения для этого важного развивающегося класса приложений, они создают набор тестов, который может фиксировать их характеристики обработки и ввода-вывода.

Пакет Benchmark состоит из четырех рабочих нагрузок:

- Обнаружение края;
- Поиск близости;
- Сканирование данных;
- Слияние данных.

Был сделан вывод, что приложения, использующие неструктурированные данные для бизнес-процессов, очень интенсивны по вводу-выводу и предъявляют высокие требования к системе хранения [2].

Но предварительные исследования оставляют возможность продолжить подобные эксперименты с неструктурированными данными.

История исследований

В следующем разделе будут объяснены предыстория и ключевые элементы, важные для этого исследования: различные методы хранения неструктурированных данных (преимущества и недостатки), используемые типы данных и экспериментальная среда. Наша основная цель — найти эффективный метод хранения неструктурированных данных. В ходе этого процесса мы проведем обширный сравнительный анализ на основе четко определенных параметров.

Исследование основано на следующих методах хранения неструктурированных данных:

- Неструктурированные данные в среде реляционных баз данных;
- Неструктурированные данные вне реляционной среды.

Мы начнем с объяснения этих методов, каждого преимущества и побочных эффектов.

А. Неструктурированные данные в реляционной среде

Спор обычно возникает относительно темы реляционных баз данных и данных больших двоичных объектов (BLOB): интегрировать большие двоичные объекты в базу данных или хранить их в файловой системе? Каждый метод имеет свои преимущества и недостатки.

Хранение неструктурированных данных, таких как изображения, аудиофайлы и исполняемые файлы в базе данных с типичными текстовыми и числовыми данными, позволяет вам хранить вместе всю связанную информацию для данного объекта базы данных

(рис. 1). И этот подход позволяет легко искать и извлекать данные BLOB; вы просто запрашиваете соответствующую текстовую информацию. Однако хранение неструктурированных данных может значительно увеличить размер ваших баз данных.

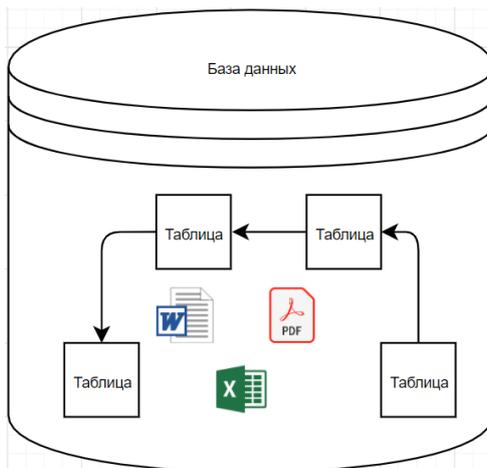


Рисунок 1 - Среда базы данных с неструктурированными данными внутри

Таблица-1. Характеристики базы данных

Характеристики базы данных					
	Назначение базы данных	Размер с неструктурированными данными	Количество пользователей	Количество BLOB-записей	Время резервного копирования
1.	LMS/CM S	≈ 12 ГБ	≈ 5000 студентов	≈ 7000 файлов размером от 100 Кб до 15 Мб	25 мин.

Основные преимущества использования этого метода:

- Когда данные хранятся в базе данных, резервная копия является согласованной. Нет необходимости в отдельной политике резервного копирования;

- Когда данные хранятся в базе данных, это часть транзакции. Так, например, откат включает все традиционные операции с базой данных вместе с операциями с двоичными данными. Обычно это делает клиентское решение более надежным и с меньшим количеством кода.

Основные недостатки использования этого метода:

- Хранение неструктурированных данных позволяет значительно увеличить размер баз данных;

- Резервное копирование и восстановление может занять много времени;

- Проблемы с производительностью подсистем ввода-вывода;

Б. Неструктурированные данные вне реляционной среды

Распространенной альтернативой вышеописанному методу является хранение двоичных файлов вне базы данных с последующим включением в качестве данных в базу данных пути к файлу или URL-адреса объекта.

Этот отдельный метод хранения имеет несколько преимуществ по сравнению с интеграцией данных BLOB в базе данных. Это несколько быстрее, потому что чтение данных из файловой системы требует немного меньше накладных расходов, чем чтение

данных из базы данных. А без больших двоичных объектов базы данных, как правило, меньше.

Однако нам необходимо вручную создать и поддерживать связь между базой данных и файлами внешней файловой системы, которые могут рассинхронизироваться. Кроме того, обычно нам требуется уникальная схема именования или хранения файлов ОС, чтобы четко идентифицировать потенциально сотни или даже тысячи файлов BLOB [1].

Хранение данных BLOB в базе данных устраняет эти проблемы, позволяя хранить данные BLOB вместе с соответствующими реляционными данными (рис. 2).

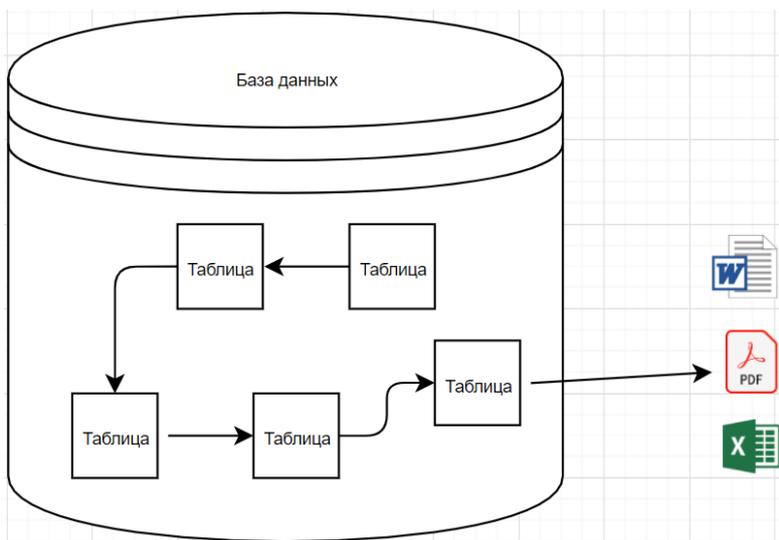


Рисунок 2 - Среда базы данных без неструктурированных данных внутри

В таблице 2 у нас такая же среда базы данных, как и в таблице 1. Единственное отличие состоит в том, что мы исключаем пару таблиц, которые содержат неструктурированные данные внутри физического файла базы данных. Разница более чем очевидна: размер базы данных при одинаковом количестве пользователей составляет половину от первоначального. Давайте подробнее рассмотрим преимущества этого метода.

Таблица-2. Характеристики базы данных

	Характеристики базы данных				
	Назначение базы данных	Размер с неструктурированными данными	Количество пользователей	Количество BLOB-записей	Время резервного копирования
1.	LMS/CMS	≈ 6 ГБ	≈ 5000 студентов	0	11 мин.

Основные недостатки использования этого метода:

- Когда данные хранятся вне базы данных, резервная копия не согласована;
- Неструктурированные данные не являются частью транзакции.

Основные преимущества использования этого метода:

- Хранение неструктурированных данных за пределами базы данных может снизить размер баз данных и снизить пропускную способность ввода-вывода;

- Резервное копирование и время восстановления могут занять меньше времени;

С. Гибридный способ хранения неструктурированных данных

Проблема с первыми двумя методами заключается в том, что мы не знаем, как фактический размер данных влияет на производительность базы данных. Результаты пока не говорят нам, есть ли разница в хранении BLOBS: 10 КБ, 1 МБ, 100 МБ и т. д.

Основные поставщики баз данных теперь поддерживают гибридный способ хранения неструктурированных данных. В этом случае данные находятся «вне» среды базы данных. Но главное преимущество заключается в том, что BLOB объекты находятся в условиях транзакционной согласованности базы данных (рис. 3).

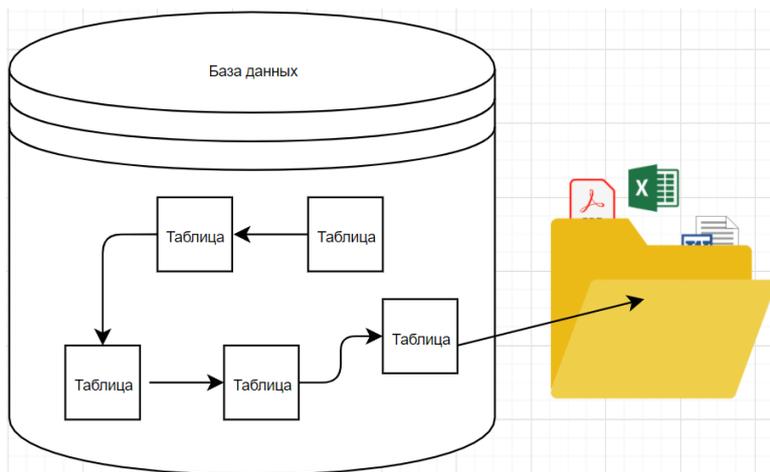


Рисунок 3 - Гибридный способ хранения неструктурированных данных

В нашем случае используются те типы данных, которые предлагаются в новой версии SQL Server 2008. Использование этой новой технологии фактически делает возможным гибридный метод хранения неструктурированных данных. Ядро базы данных поддерживает новый тип данных с именем *filestream*.

Ключевым аспектом исследования является выяснение эффективности этого метода хранения данных.

Экспериментальная среда

Для этого мы настраиваем тестовую среду со следующими компонентами:

- Процессор: AMD X2 3.0 Ghz 4 ГБ;
- Физическая память;
- Файлы базы данных на С;
- Компонент файлового потока на диске Е;
- Диски С: и Е: на отдельных физических дисках SATA;
- SQL Server 2008 R2 и клиентское приложение для настраиваемого тестирования

находятся на одном компьютере;

На следующих диаграммах показано среднее время загрузки для следующих условий:

- Файл размером 10 КБ, повторяется 3 раза для каждого измерения (рис. 4);
- Файл размером 1 МБ, повторяется 3 раза для каждого измерения (рис. 5);
- Файл размером 10 МБ, повторяется 3 раза для каждого измерения (рис. 6);

Результаты для файла размером 10 КБ

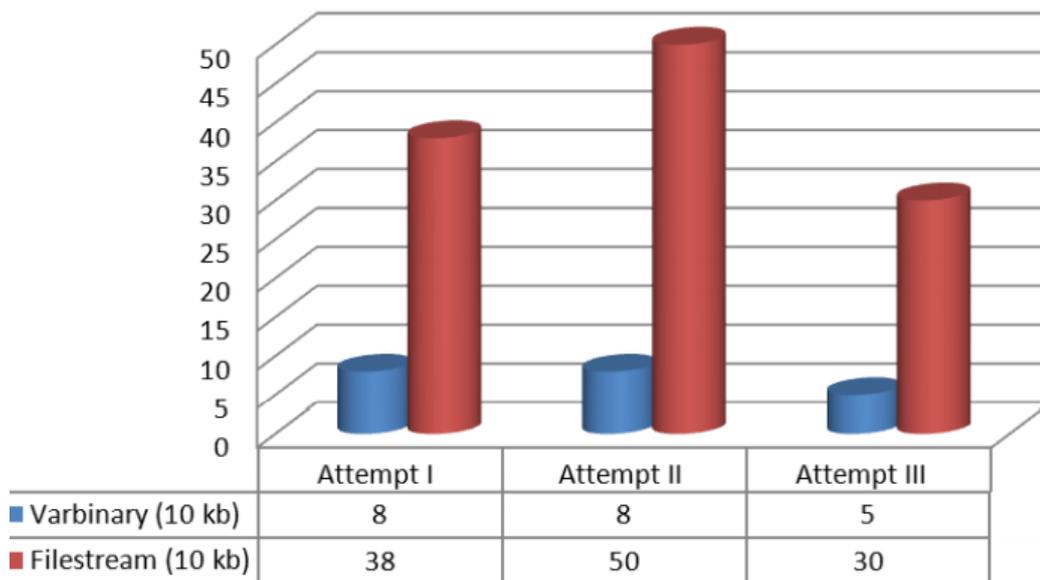


Рисунок 4 - Результаты сохранения файла размером 10 КБ внутри базы данных и файловой системы.

В этом измерении вы можете четко увидеть накладные расходы и влияние на производительность, вызванные использованием файлового потока на “маленькие” файлы. Время при хранении внутри базы данных было каждый раз меньше 10 миллисекунд. С другой стороны, время хранения с использованием файлового потока в 4-5 раз больше.

Результаты для файла размером 1 МБ

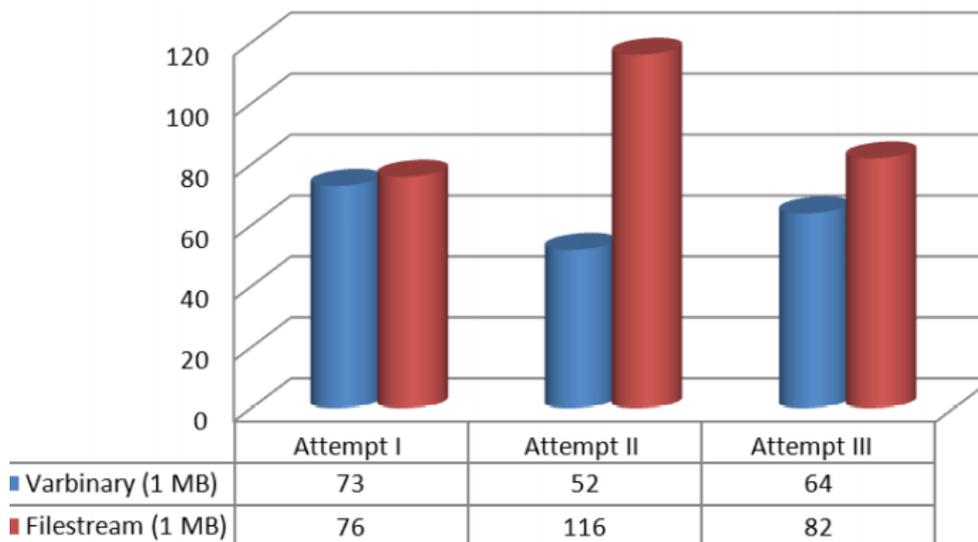


Рисунок 5 - Результаты сохранения файла размером 1 МБ внутри базы данных и файловой системы.

При объеме данных 1 МБ традиционный двоичный файл и файловый поток действуют одинаково, и разница между ними — максимум на один раз быстрее.

Результаты для файла размером 10 МБ

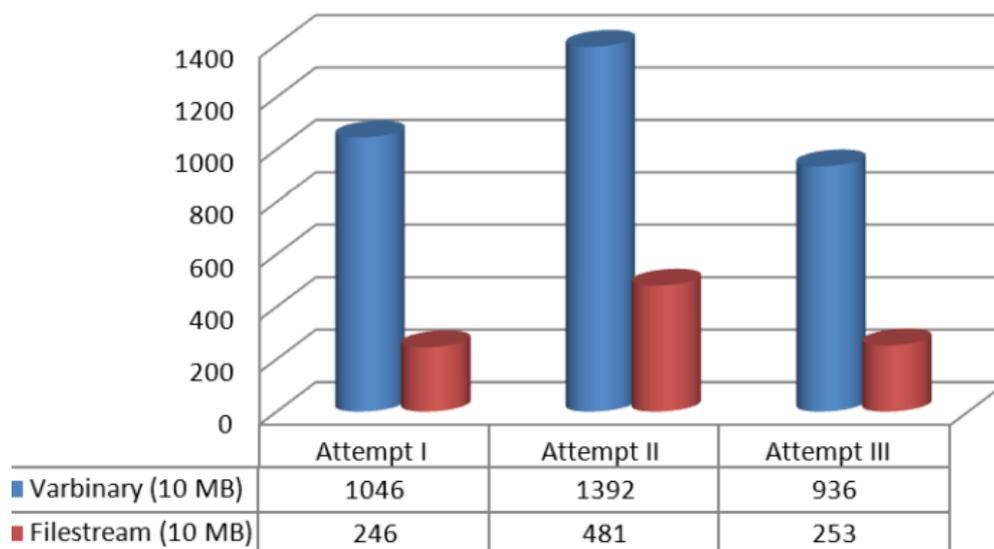


Рисунок 6 - Результаты сохранения файла размером 10 МБ внутри базы данных и файловой системы.

Когда используются файлы размером 10 МБ, хранение данных в традиционном файле базы данных происходит намного медленнее. Основываясь на этих измерениях, можно сказать, что более эффективно использовать файловый поток, когда типичный размер файла составляет около 1 МБ или более. Если файлы имеют небольшой размер (явно менее 1 МБ), традиционное хранилище работает лучше. При измерении времени мы обнаружили, что удаление строк на основе файлового потока происходит намного быстрее, чем при хранении внутри таблицы.

Заключение

В исследовании представлена гипотеза о неэффективности существующих методов хранения и извлечения неструктурированных данных.

Новый способ использования файловой системы при согласованности базы данных был хорошей возможностью опробовать новую технологию в разных областях хранения данных. Примеры тестов сохранения файлов размерами 10 КБ, 1 МБ и 10 МБ в реальных средах информационных систем могут относиться к фотографиям пользователей, файлам резюме, небольшим офисным документам, видео - и аудиофайлам. Результаты исследования могут использоваться для моделирования системы и четкого определения потребностей в хранении данных. Преимущество заключается в получении максимальной производительности на основе аппаратной и программной инфраструктуры. Также в системах, в которых проблемы с производительностью уже существуют, эта модель может помочь выявить узкие места и найти способ их улучшить.

Результаты исследования позволяют создать модель тестирования системы для хранения неструктурированных данных. Дальнейшие шаги по улучшению — реализация этой модели в реальных средах, где важны неструктурированные данные (платформы электронного обучения, социальные сети, порталы хранения видео и т. д.). После этого настоящее исследование может стать официальной методологией анализа системы с потребностями в неструктурированных данных. Современные тенденции [3] показывают нам, что www становится глобальным хранилищем для всех видов данных, где преобладают неструктурированные данные.

СПИСОК ЛИТЕРАТУРЫ

1. Hitachi Global Storage Technologies - Обзорные схемы технологии HDDT. [Электронный ресурс] URL:<http://www.hitachigst.com/hdd/technolo/overview/storagetechchart.html>. (дата обращения: 10.02.2021)
2. Набор тестов для обработки неструктурированных данных Smullen, C.W.; Tarapore, S.R.; Gurumurthi, S.; Архитектура сети хранения данных и Параллельный ввод-вывод, 2007г. SNAPI. Международный семинар 24 сентября. 2007 Страниц: 79 – 83.
3. J. Gantzandetal. Расширение цифровой вселенной - прогноз роста мировой информации до 2010 г., март 2007 г. Белая книга IDC.
4. C.White. Консолидация, доступ и анализ неструктурированных данных, O'Reilly Media, 2005 г. Страницы 118–124.
5. Р. Чемберлен, М. Франклин. Использование реконфигурируемости для текстового поиска. В материалах семинара по высокопроизводительным встроенным вычислениям (HPEC), O'Reilly Media, 2006 г. Страницы 1178–1181.
6. Повышение доступности данных с помощью файловых сетей Geer, В; Компьютер. O'Reilly Media, 2017 г. Страницы 1356–1359.
7. Основы классификации неструктурированных данных Островски, Д.А.; Семантические вычисления, 2009. ICSC '09. Международная конференция IEEE, 14–16 сентября 2009 г. Страницы: 373–377.
8. Анализ неструктурированных данных и оптимизация их хранения. [Электронный ресурс] URL: <https://habr.com/ru/company/hpe/blog/265499>. (дата обращения: 19.02.2021)
9. Преобразование неструктурированных данных из разрозненных источников в знания Plejic, B.; Vujnovic, B.; Penco, R.; Семинар по приобретению знаний и моделированию, 2008 г. Семинар по КАМ, 2008 г. Международный симпозиум IEEE 21–22 декабря 2008 г. Страницы: 924 – 927.
10. G.Fountainand и S.Drager. Высокопроизводительная архитектура слияния в реальном времени. O'Reilly Media, 2002 г. Страницы 1478–1485.

REFERENCES

1. Hitachi Global Storage Technologies – Overview diagrams of HDDT technology. [Electronic resource] URL:<http://www.hitachigst.com/hdd/technolo/overview/storagetechchart.html>. (date of the application: 10.02.2021)
2. A Benchmark Suite for Unstructured Data Processing Smullen, C.W.; Tarapore, S.R.; Gurumurthi, S.; Storage Network Architecture and Parallel I/Os, 2007. SNAPI.International Workshop on 24 Sept. 2007 Page(s): 79 – 83
3. J.Gantzandetal. The Expanding Digital Universe – A Forecast of Worldwide Information Growth Through 2010, March 2007. IDC Whitepaper
4. C.White. Consolidating, Accessing, and Analyzing Unstructured Data, O'Reilly Media, 2005 Page(s): 118 – 124
5. R. Chamberlain, M. Franklin and R. Index. Using reconfigurability for text search. In High Performance Embedded Computing (HPEC) Workshop Proceedings. O'Reilly Media, 2006 Page(s): 1178 – 1181
6. Increasing data availability with Geer, D. File networks; Computer. O'Reilly Media, 2017 Page(s): 1356 – 1359
7. Fundamentals of unstructured data classification Ostrowski, D.A.; Semantic Computing, 2009. ICSC '09. IEEE International Conference, September 14-16, 2009, Pages: 373-377
8. Analysis of unstructured data and optimization of their storage. [Electronic resource] URL: <https://habr.com/ru/company/hpe/blog/265499>. (Date accessed: 19.02.2021)
9. Converting unstructured data from disparate sources to knowledge Plejic, B.; Vujnovich, B.; Penco, R. ; Knowledge Acquisition and Modeling Workshop 2008 KAM Workshop 2008 IEEE International Symposium December 21-22, 2008 Pages: 924 - 927

10. J.Fountain and S.Drager. High-performance real-time merge architecture. O'Reilly Media, 2002 Page(s): 1478 – 1485

Н.М. Еркетаев, Н.К. Мұқажанов
Құрылымсыз деректерді тиімді сақтау

Аңдатпа. Бұл мақалада құрылымсыз деректерді мәлімет базасында және классикалық файлдық жүйеде сақтау үшін сан түрлі типтегі деректерді пайдалану бойынша зерттеу нәтижелерін ұсынамыз. Ол сонымен қатар деректерді сақтаудың заманауи әдістерін, фондық тарихты және негізгі элементтерді түсіндіреді. Осы зерттеуге қатысты: құрылымдалмаған мәліметті сақтаудың әртүрлі әдістері (артықшылықтары мен кемшіліктері), қолданылатын мәлімет типтері және эксперименттік орта. Біздің басты мақсат – құрылымдалмаған деректерді сақтаудың тиімді әдісін табу.

Түйінді сөздер: мәлімет қоры, құрылымдалмаған мәлімет, файлдар ағыны, тиімділік, өнімділік.

N.M. Yerketayev, N.K. Mukazhanov
Efficient storage of unstructured data

Abstract. In this article we present the research findings on the use of different types of unstructured data stored in a database and the classic file system. It also explains modern data storage methods, background history and key elements. The issues relevant to this research: different methods of storing unstructured data (advantages and disadvantages), the used data types and the experimental environment. Our main goal is to find an efficient method for storing unstructured data.

Keywords: database, unstructured data, file stream, efficiency, performance.

Авторлар туралы мәлімет:

Еркетаев Нұрзат Мейірханұлы, «Компьютерлік инженерия және ақпараттық қауіпсіздік» кафедрасының магистранті, Халықаралық ақпараттық технологиялар университеті.

Мұқажанов Нуржан Какенұлы, PhD, «Компьютерлік инженерия және ақпараттық қауіпсіздік» кафедрасының ассистент-профессоры, Халықаралық ақпараттық технологиялар университеті.

Сведения об авторах:

Еркетаев Нұрзат Мейірханұлы, магистрант кафедрасы «Компьютерная инженерия и информационная безопасность», Международный университет информационных технологий.

Мұқажанов Нуржан Какенович, PhD, ассистент-профессор кафедрасы «Компьютерная инженерия и информационная безопасность», Международный университет информационных технологий.

About the authors:

Nurzat M. Yerketayev, master student, Department of Computer Engineering and Information Security, International Information Technology University.

Nurzhan K. Mukazhanov, PhD, Assistant-Professor, Department of Computer Engineering and Information Security, International Information Technology University.

INTERNATIONAL JOURNAL OF INFORMATION AND
COMMUNICATION TECHNOLOGIES

МЕЖДУНАРОДНЫЙ ЖУРНАЛ ИНФОРМАЦИОННЫХ И
КОММУНИКАЦИОННЫХ ТЕХНОЛОГИЙ

ХАЛЫҚАРАЛЫҚ АҚПАРАТТЫҚ ЖӘНЕ
КОММУНИКАЦИЯЛЫҚ ТЕХНОЛОГИЯЛАР ЖУРНАЛЫ

Ответственный за выпуск	Есбергенов Досым Бектенович
Редакторы	Далабаева Айсара Касымбековна Джоламанова Балия Джалгасбаевна Медведев Евгений Юрьевич
Компьютерная верстка	Туратауова Айжаркын Ахметовна
Компьютерный дизайн	Туратауова Айжаркын Ахметовна

Редакция журнала не несет ответственности за
недостоверные сведения в статье и
неточную информацию по цитируемой литературе

Подписано в печать 26.06.2021 г.
Тираж 500 экз. Формат 60x84 1/16. Бумага тип.
Уч.-изд.л. 10.1. Заказ №165

Издание Международный университет информационных технологий
Издательский центр КБТУ, Алматы, ул. Толе би, 59