

УДК 004.3

А.А.Егошина, А.С.Вороной.  
Государственный университет информатики и искусственного интеллекта  
postmaster@suiai.edu.ua

## Повышение эффективности извлечения информации из слабо структурированных источников на основе метаданных и базы знаний

*Предлагается технология повышения эффективности поиска в слабо структурированных базах данных с Web-интерфейсом на основе метаданных, представляющих структуру баз данных и базы знаний, описывающей семантику хранимых данных. Предложена модель метаданных для рассматриваемых баз данных и структура хранилища метаданных. Хранение метаданных и базы знаний на отдельном сервере позволяет сократить время обработки запросов. Приведены результаты экспериментальных исследований в производственной среде работы с базами данных для различных профилей операций пользователей. Результаты экспериментов показали, что при реальных значениях числа запросов к базам данных предлагаемая технология обеспечивает комфортное для пользователей время обработки запросов и сокращение нагрузки на серверы системы. С ростом числа запросов преимущества предложенной технологии возрастают.*

**Ключевые слова:** метаданные, база знаний, база данных, поиск, ресурс.

### Общая постановка проблемы

Непрерывный рост количества электронных документов и их доступности, а также увеличение объемов баз данных информационных систем (ИС) наряду со слабой структурированностью информационных фондов осложняет управление информацией и работу пользователей с ней.

За последние несколько лет активное развитие получило направление в информационных технологиях, занимающееся проблемами учета семантики в рамках информационных систем. ИС, использующие семантические технологии, отличаются от традиционных ИС следующими особенностями: использованием знаний предметной области, в которой проводится обработка информации; знания предметной области выражаются явно, в виде моделей; модель выражает смысл терминов (понятий) предметной области через связи между ними.

Для решения проблемы совершенствования доступа к растущему объему информации и информационным услугам необходимо проводить обработку на семантическом уровне с использованием хорошо структурированного и постоянно применяемого стандарта метаданных. Такой стандарт позволит пользователям совершать поиск в большом количестве таблиц баз данных и уверенно определять местонахождение интересующей информации.

### Информационный поиск на основе использования метаданных

Метаданные определяют ортогональный основному уровню описания информации (который формируется такими понятиями, как классы, типы данных и др.) уровень описания свойств [1]. Метаданные также могут описывать схему информационного источника, например реляционную, объектную или слабоструктурированную.

Контекстные метаданные описывают связь объекта с другими объектами системы, а контентные метаданные описывают содержимое объекта (т. е., имеющиеся в объекте знания).

Семантические метаданные (англ. semantic metadata) — дополнительный слой информации, позволяющий автоматизировать установление связей между различными частями контента. Семантические метаданные позволяют обогатить имеющийся контент принципиально новыми связями, которые невозможно установить при стандартной кластеризации.

Использование метаданных, в особенности контентных (семантических), позволяет эффективно решать такие задачи работы со знаниями, как поиск, категоризация и рекомендация знаний.

Для поиска необходимой информации в настоящее время в основном применяется поиск по ключевым словам. Пользователю приходится проходить по нескольким десяткам ссылок и изучать огромное количество информации, основная часть которой не является релевантной. Гораздо удобнее иметь поиск, выдающий однозначный результат, который содержит в

структурованому в усіх випадках інформацію про об'єкт, сортировану за релевантністю.

Додаток запитів шаблонами, логічними операторами, спеціальними словами для пошуку в метаданих і іншими засобами уточнення картини принципово не змінює. Більш цінними виявляються якісна підтримка морфології різних мов і врахування відстані між словами. Однак при обробці великих масивів даних, і вони не дозволяють радикально знизити «шум пошуку» – це просте наявність шуканих слів далеко не завжди коректно відображає тему документа [2].

Пошук за метаданими в порівнянні з детальним пошуком за інформаційними ресурсами є більш ефективним, оскільки метадані надають інформацію, що входить до ресурсу.

### Технологія пошуку в слабо структурованих інформаційних джерелах з використанням бази знань

При наявності великих обсягів даних одне навіть просте змінення може несприятливо впливати на системи всього підприємства і зробити непридатними моделі даних, звіти або призвести до невідповідності даних. Візуальне представлення залежностей даних спрощує завдання інтеграції, аналізу, управління даними, отримуючи суттєве скорочення витрат часу при виконанні аудиту.

Ініціативи з інтеграції даних зазвичай натрапляють на відсутність знань про дані як такі. Пошук і ідентифікація потрібних даних, так само як і визначення їх місцезнаходження, можуть займати до 70% часу всього інтеграційного проекту [3].

Використання бази знань в гетерогенних інформаційних джерелах надає гнучкий підхід до інтеграції інформації.

База знань є проміжним шаром, розташованим між інформаційними джерелами (базами даних з метаданими) і користувачами інформації, як це показано на рисунку 1.

Використання бази знань, як певного уніфікованого інтерфейсу для вирішення завдань над багаточисельними неструктурованими інформаційними джерелами, звільняє користувача від необхідності знаходити релевантні джерела, задавати запити до кожного з них окремо і вручну порівнювати інформацію з них.



Рисунки 1 – Технологія пошуку в слабо структурованих інформаційних джерелах з використанням бази знань і метаданих

Метадані описують ресурси з допомогою малих простих пакетів інформації, які легко знайти, і до яких є доступ для великої кількості користувачів. Метадані значно покращують ступінь структуризації пошукових завдань в межах обсягу даних, наявних.

В даній роботі метадані визначаються як трійка наступного виду:

$$META = \{T, S, C\},$$

де  $T$  – сукупність відомостей про таблиці бази даних:

$$T = \{id\_t, Sch, title\_t, \{Inf\_t\}\},$$

де  $id\_t$  – код таблиці;

$Sch$  – схема таблиці (тип групування по певним загальним ознакам);

$title\_t$  – назва таблиці;

$\{Inf\_t\}$  – множина допоміжної інформації, наприклад, ім'я обробника таблиці, коди таблиць – „родителів” і т.д.

$S$  – сукупність відомостей про схеми таблиць:

$$S = \{id\_s, schema\_name\},$$

де  $id\_s$  – код схеми;

$schema\_name$  – назва схеми.

$C$  – сукупність відомостей про стовпці всіх таблиць бази даних:

$$C = \{id\_c, type\_data, title\_table, title\_column, \{Inf\_c\}\},$$

де  $id\_c$  – код стовпця;

$type\_data$  – тип даних, що містяться в стовпці;

$title\_table$  – назва таблиці, що містить даний стовпець;

$title\_column$  – назва стовпця;

$\{Inf\_c\}$  – додаткова інформація про стовпець.

При великому обсязі ресурсів (таблиць бази даних) і їх слабкій структуризованості ми будемо зберігати метадані окремо від ресурсу в сховищі метаданих на окремому сервері, як показано на рисунку 2.

Сховища метаданих використовуються для зберігання і управління метаданими. Такий механізм надає більшу гнучкість, так як метадані можна зробити відкритими в різних

установках и синтаксисах, которые со временем можно легко модифицировать. Глобальные изменения и поправки можно осуществлять после изначального создания.

Схема хранилища метаданных приведена на рисунке 2.

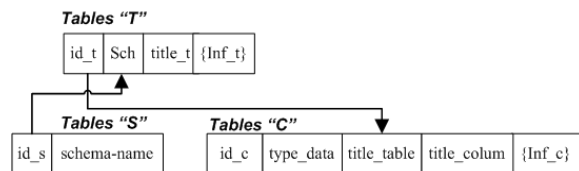


Рисунок 2 – Схема хранилища метаданных

Создание и управление предлагаемыми в данной работе метаданными позволяет решить следующие задачи.

1. Поддержка интеграции информационной системы. Схемы и интеграция данных зависят от метаданных, описывающих структуру и смысл отдельных источников данных и целевых систем.

2. Поддержка анализа и проектирования новых приложений. Метаданные повышают контролируемость и надежность процесса разработки приложений, обеспечивая информацию о смысле данных, их структуре и источниках. Кроме этого, метаданные, касающиеся решений по проектированию приложений, можно использовать повторно.

3. Повышение гибкости информационной системы и возможности повторного использования существующих программных модулей. Это возможно только для активного и полуактивного использования метаданных. Быстро изменяющиеся семантические аспекты явным образом хранятся в виде метаданных вне прикладных программ. Систему можно расширить и адаптировать без всяких трудностей. Данный подход также дает возможность повторного использования «фрагментов кода»;

4. Автоматизация административных процессов. Метаданные управляют запуском различных процессов (например, загрузки и обновления). Информация об их исполнении (журналы доступа, количество добавленных записей и т.п.) также содержится в репозитории, легко доступном администратору;

5. Эффективное извлечение информации и повышение качества данных, определяемого:

- согласованностью (является ли представление данных однородным, нет ли дубликатов, данных с пересекающимися или конфликтующими определениями);
- полнотой (все ли данные присутствуют);
- точностью (совпадением хранимых и фактических значений);
- своевременностью (актуально ли хранимое значение).

## Экспериментальные исследования

С целью определения эффективности использования предложенного в работе подхода были проведены экспериментальные исследования на серверах компании Prometheus Research [4]. Основополагающим программным обеспечением компании является инновационная технология HTSQL, предоставляющая предприятиям и организациям, проводящим биомедицинские исследования, средства управления web-ориентированными базами данных в режиме реального времени.

Тестирование этого сценария должно помочь в оценке того, как изменение различных факторов и наличие метаданных влияет на производительность системы Prometheus Research Informatics Services And Solutions.

К основным характеристикам эксперимента относятся: количество пользователей, одновременно работающих с системой; типы выполняемых пользователями операций; количество опрашиваемых документов в индексе.

Для тестирования использовалась конфигурация, в состав которой входил один сервер индексирования (Xeon MP 3.0, 16 Gb RAM, 1Tb HDD) и один сервер баз данных (Xeon MP 2.2 8Gb RAM, 1Tb HDD).

В ходе тестирования был произведен обход около 500 элементов (таблиц базы данных).

В таблице 1 представлены примеры элементов, обход которых был осуществлен.

Размер элементов составляет от 10 КБ до 100 МБ.

Таблица 1. Характеристики элементов

Название элемента	Число столбцов	Число строк
ssc_med_hx_birth_defects	96	808
ssc_med_hx_diet_medication_sleep_puberty	221	800
ssc_med_hx_language_disorders	139	806
ssc_med_hx_v2_labor_delivery_birth_feeding	130	802
ssc_med_hx_v2_medication_drugs_mother	225	810
ssc_med_hx_v2_pregnancy_history	247	790
ssc_med_hx_v2_psychiatric_disorders	182	806

В таблице 2 показано использование дискового пространства.

Таблица 2. Использование дискового пространства

Размер индекса на сервере индексирования	100 МБ
Размер базы данных поиска	2 ГБ

Следует отметить, что экспериментальные исследования проводились в производственной среде, где задержка сети и реакция хранилищ, обход которых осуществлялся, влияли на скорость обхода. Скорость обхода, измеряемая количеством таблиц в секунду, могла быть значительно выше в исключительно тестовой среде или в средах с более высокой пропускной способностью и лучшей реакцией хранилищ, обход которых осуществлялся.

В таблице 3 представлены результаты тестирования для различных профилей операций пользователей на основе оборудования и профиля использования, указанного выше. Таблица содержит следующие столбцы: № - номер эксперимента; з/с - запросов в секунду; %ЦПСИ - процент использования ресурсов центрального процессора сервера индексирования; %ЦПСБД - процент использования ресурсов центрального процессора сервера баз данных; о/с - средняя скорость операции (поиска, извлечения, записи) в секунду на диск сервера баз данных.

Таблица 3. Результаты экспериментов

№	з/с	%ЦПСИ	%ЦПСБД	о/с
1	1	0,038	0,016	0,78
2	10	0,56	0,17	1,019
3	50	0,87	0,48	1,098
4	100	5,2	6,7	5,2
5	500	11	10,27	7,8
6	1000	36,4	28,3	9,2
7	10000	54,2	46,3	11,6

В таблице 4 приведены результаты аналогичных экспериментов, но с использованием метаданных, предложенных в данной работе.

Результаты экспериментов показали, что при реальных значениях числа пользовательских запросов к базам данных (>1000 запросов в сек.) предлагаемая авторами технология обеспечивает комфортное для пользователей время обработки запросов.

Таблица 4. Результаты экспериментов с использованием метаданных

№	з/с	%ЦПСИ	%ЦПСБД	о/с
1	1	0,031	0,01	0,72
2	10	0,5	0,11	1,01
3	50	0,72	0,3	1,09
4	100	1	0,5	1,5
5	500	1,2	0,8	1,7
6	1000	2	1,5	3
7	10000	5	4,7	7

Кроме того, сокращается нагрузка на серверы системы, что позволяет снизить требования к аппаратному обеспечению. С ростом числа запросов преимущества предложенной технологии возрастают.

### Заключение

В настоящей работе предложен новый подход к извлечению информации из слабо структурированных источников на основе метаданных и базы знаний.

Использование метаданных и базы знаний в информационной системе позволяет эффективно проводить мониторинг и оценивать качество данных, их содержание, применение в проектах и производительность. Компании и предприятия могут быстро и эффективно создавать собственные отчеты для своих интеграционных нужд.

### Литература

1. <http://www.w3.org/Metadata/>
2. Даниэла Флореску Технологии баз данных для World-Wide Web: Обзор. Системы управления базами данных / Даниэла Флореску, Алон Леви, Альберто Мендельсон. – 1998. – 286 с.
3. Allison Powell James. The impact of database selection on distributed searching. In Proc. of the SIGIR'00, 2000.
4. <http://www.prometheusresearch.com>

Надійшла до редакції 01.03.2011