

Системы информационного  
обеспечения для  
проведения распределенных  
научных исследований

# План доклада

- Задачи систем информационного обеспечения распределенных научных исследований
- Тематический поиск информации
- Технологии построения специализированных информационных систем для проведения распределенных научных исследований

# Задачи информационного обеспечения распределенных научных исследований

- Сбор информации из внешних и внутренних источников о предмете исследования
- Автоматизация процесса обмена файлами данных между участниками
- Автоматизация процесса распределенной обработки информации

# План доклада

- Задачи систем информационного обеспечения распределенных научных исследований
- Тематический поиск информации
- Технологии построения специализированных информационных систем для проведения распределенных научных исследований

# Задачи автоматизированной системы тематического анализа информации

- извлечение данных из текстовых коллекций на естественных языках (русский, английский)
- автоматическое обучение системы и отчуждение персональных знаний у эксперта (знаний об объектах и структуре используемого языка, онтологии)
- автоматический тематический анализ информации и выделение фактов
- перманентный автоматический подбор материалов по запросам пользователей
- сохранение версий интересующей пользователя информации
- Минимальная «ручная» настройка

# Возможности системы

- Разовый поиск информации
- Непрерывный сбор информации
- Группировка и предварительная обработка информации для последующего анализа
- Классификация информации по древовидным классификаторам
- Контекстный поиск
- Персонализированный поиск

# Особенности системы

- Автоматическая рубрикация текстов с использованием неограниченного количества рубрикаторов
- Наличие персональных рубрикаторов и архивов
- Автоматическое ранжирование с учетом предпочтений пользователя
- Сервисно-ориентированная архитектура (SOAP)
- Открытый код

# Перспективные направления

- Выявление новых тематических направлений
- Выделение дат, имен, географических объектов, источников и других параметров
- Выделение границ ресурсов
- Выделение фактов и определение их достоверности

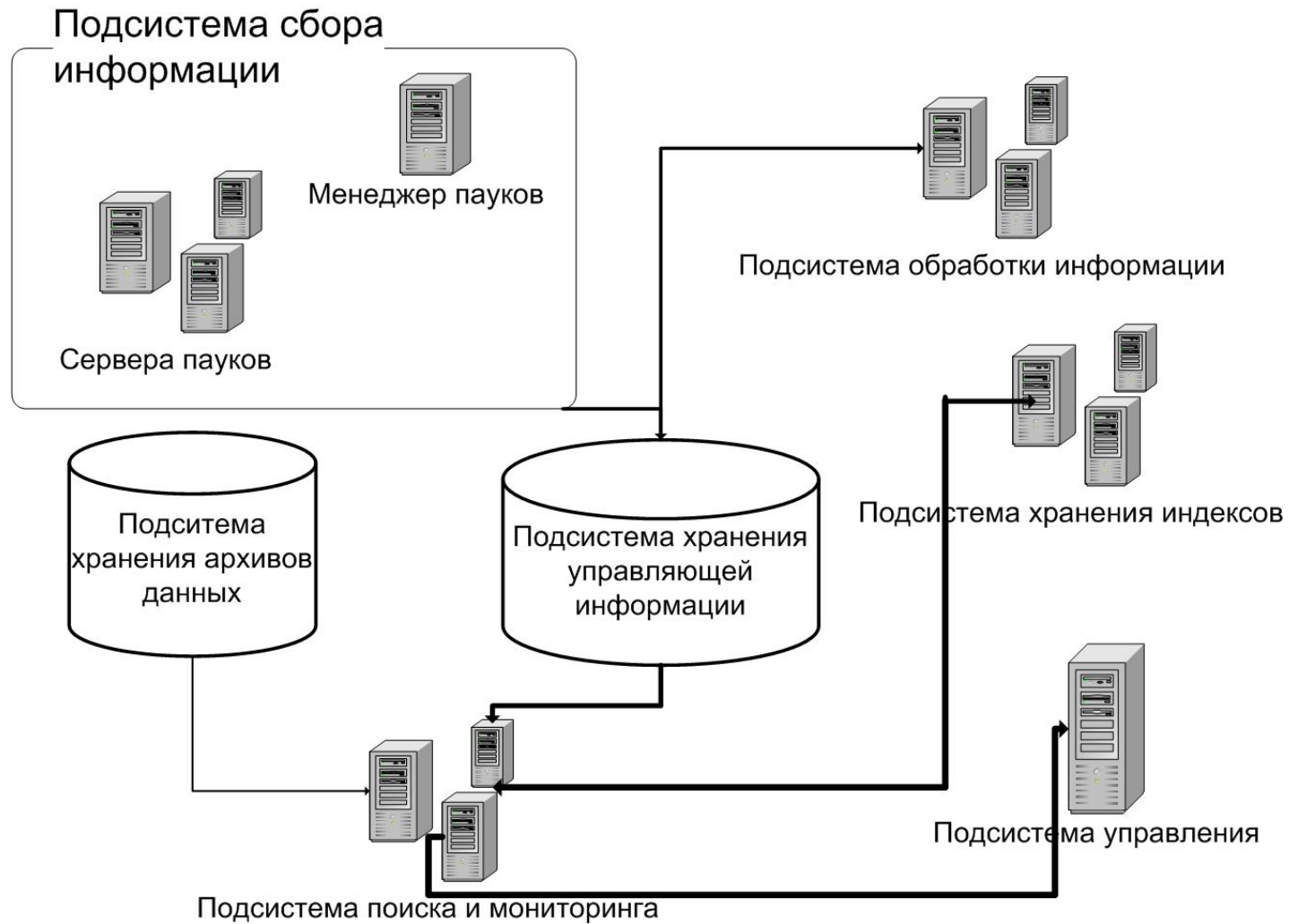


# Возможное применение

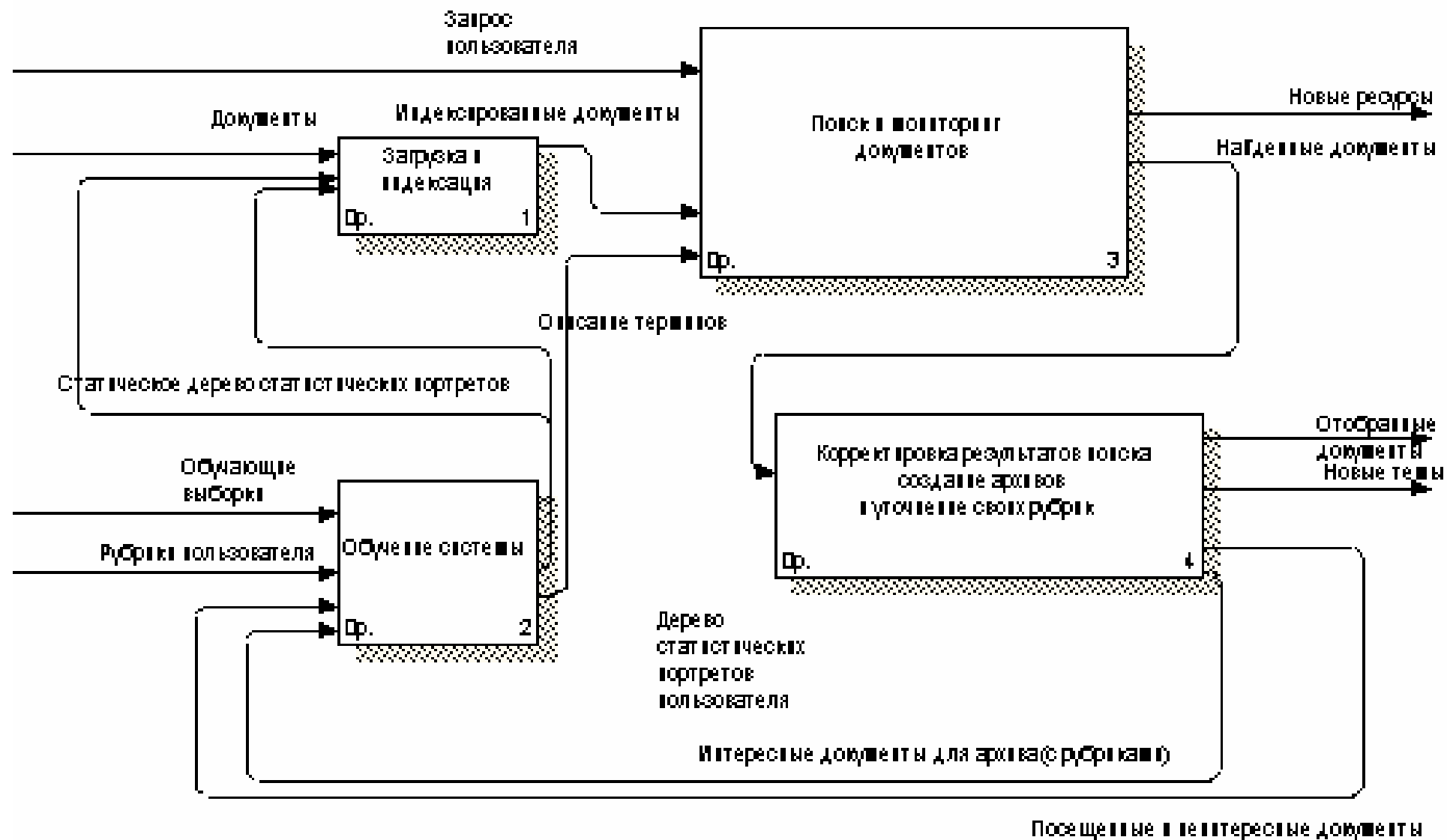
- Поиск в интернет
- Поиск по новостям, форумам, периодическим изданиям
- Поиск по специализированным хранилищам
- Внедрение в организациях, требующих сертификации ПО в области информационной безопасности



# Структура системы



# Общая схема



# Основные элементы

- Распределенное хранилище
- Модуль вноса и индексации(морфологический анализ, рубрикация, аннотирование)
- Модуль обучения
- Рубрики пользователя
- Система обратной связи
- Система поиска
- Мониторинг
- Модуль тематического анализа
- Web-интерфейс
- Модуль персонификации

# Загрузка

- Морфологический разбор текстов и загрузка пословного индекса
- Составление и загрузка аннотаций
- Выделение и загрузка ссылок документа
- Выделение и загрузка терминов документа

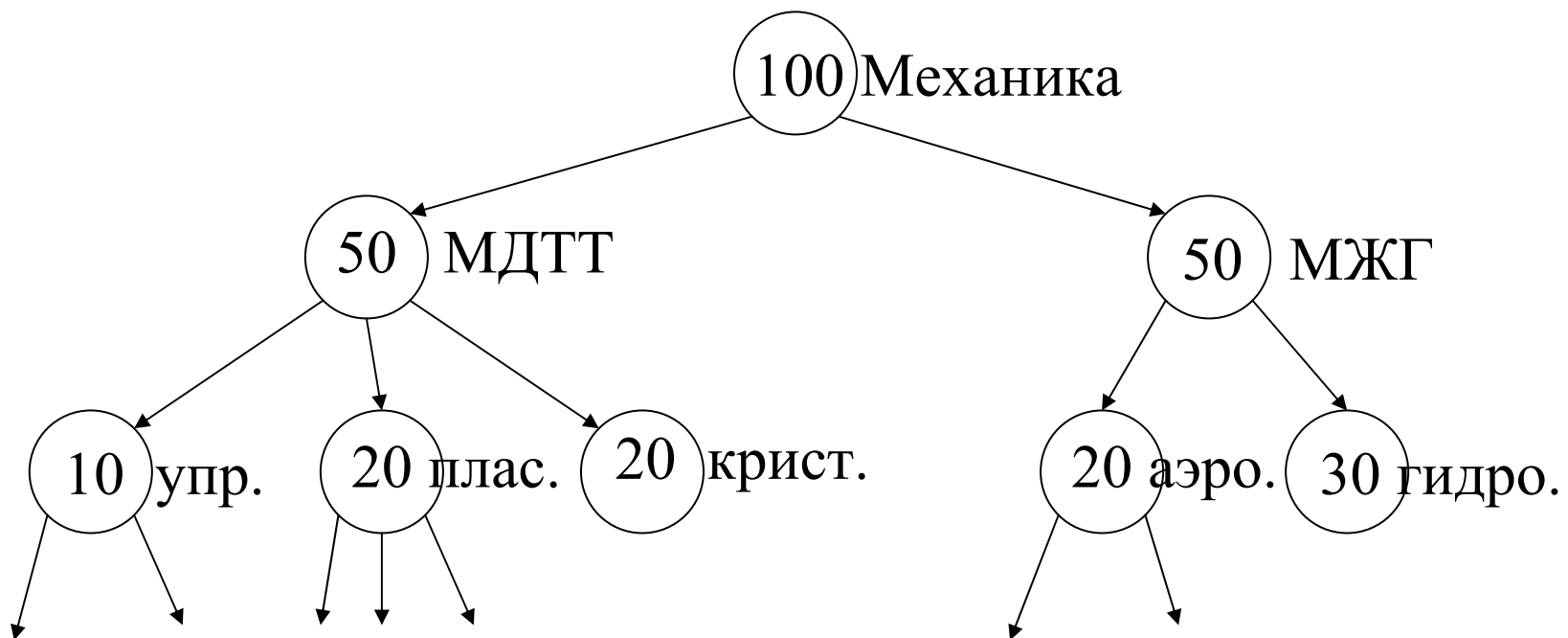
# Выделение терминов

- Вычисляются все пары, входящие в обучающие выборки
- Определяются устойчивые пары
- Определяются пары, характеризующие рубрики

# Рубрикация(1)

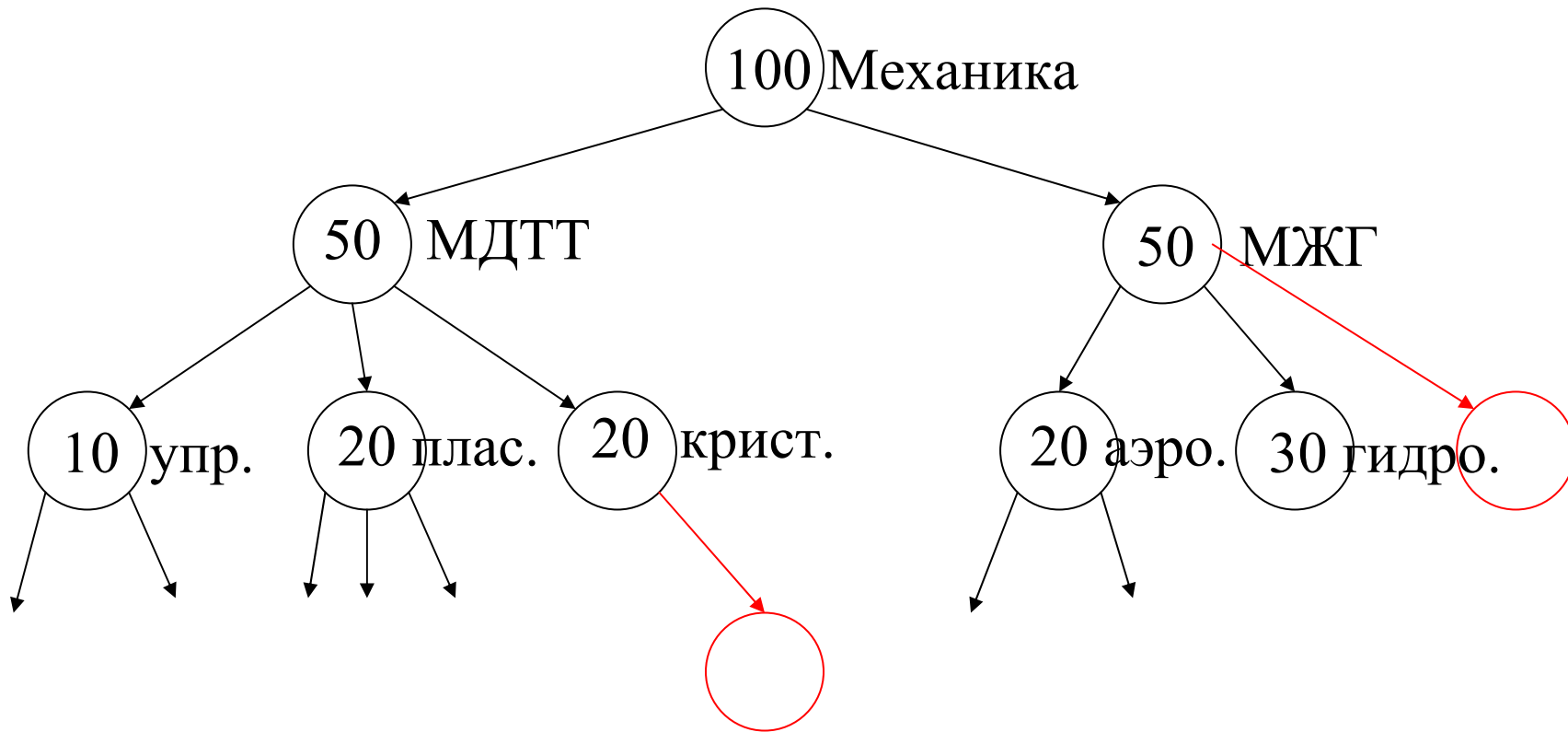
- Обучение
  - Составление классификатора
  - Подбор примеров для обучения
  - Построение статистических портретов рубрик
  - Построение дерева статистических портретов рубрик
- Автоматическая рубрикация входного потока текстов
- Дообучение системы

# Рубрикация(2)





# Рубрики пользователя



- Улучшают результаты поиска
- Позволяют использовать индекс

# Язык запросов

- Запрос может содержать выражения вида  $X\langle w1, n1 \rangle$  and  $X\langle w2, n2 \rangle$  or  $X\langle w3, n3 \rangle$   
X – тип аргумента (леммы, термины, рубрики, даты, параметры, ресурсы или ключи)  
w – символьное или числовое значение аргумента  
n – значимость при поиске

$L\langle \text{нефть} \rangle$  AND (RID $\langle 786, 0.5 \rangle$  OR RID $\langle 966 \rangle$ )

# Мониторинг

- Запоминает запрос пользователя
- Регулярно обходит интересующие пользователя ресурсы
- Создает подборки интересующих пользователя документов

# Система поиска

- Разбор запроса и построение дерева
- Считывание всех файлов индексов, упоминающихся в запросе
- Вычисление запроса для каждого документа
- Вычисление ранга документа
- Сортировка с усечением (Top100)
- Уточнение рангов по пользовательским рубрикам

# Web-интерфейс

- Поиск
- Управление мониторингом
- Управление поиском новых ресурсов
- Конфигурация рубрикатора пользователя
- Управление архивом пользователя
- Обратная связь (посещенные и неинтересные документы)

# Паук

- Менеджер пауков и группы пауков с разными приоритетами
- Правило включает маску URL, разрешение или запрещение сканирования, принадлежность рубрикам, начальный интервал переиндексации, приоритет, глубина сканирования, общее количество
- Время следующего сканирования определяется по приоритету и времени последнего изменения
- Возможно подключать пауки других производителей

# План доклада

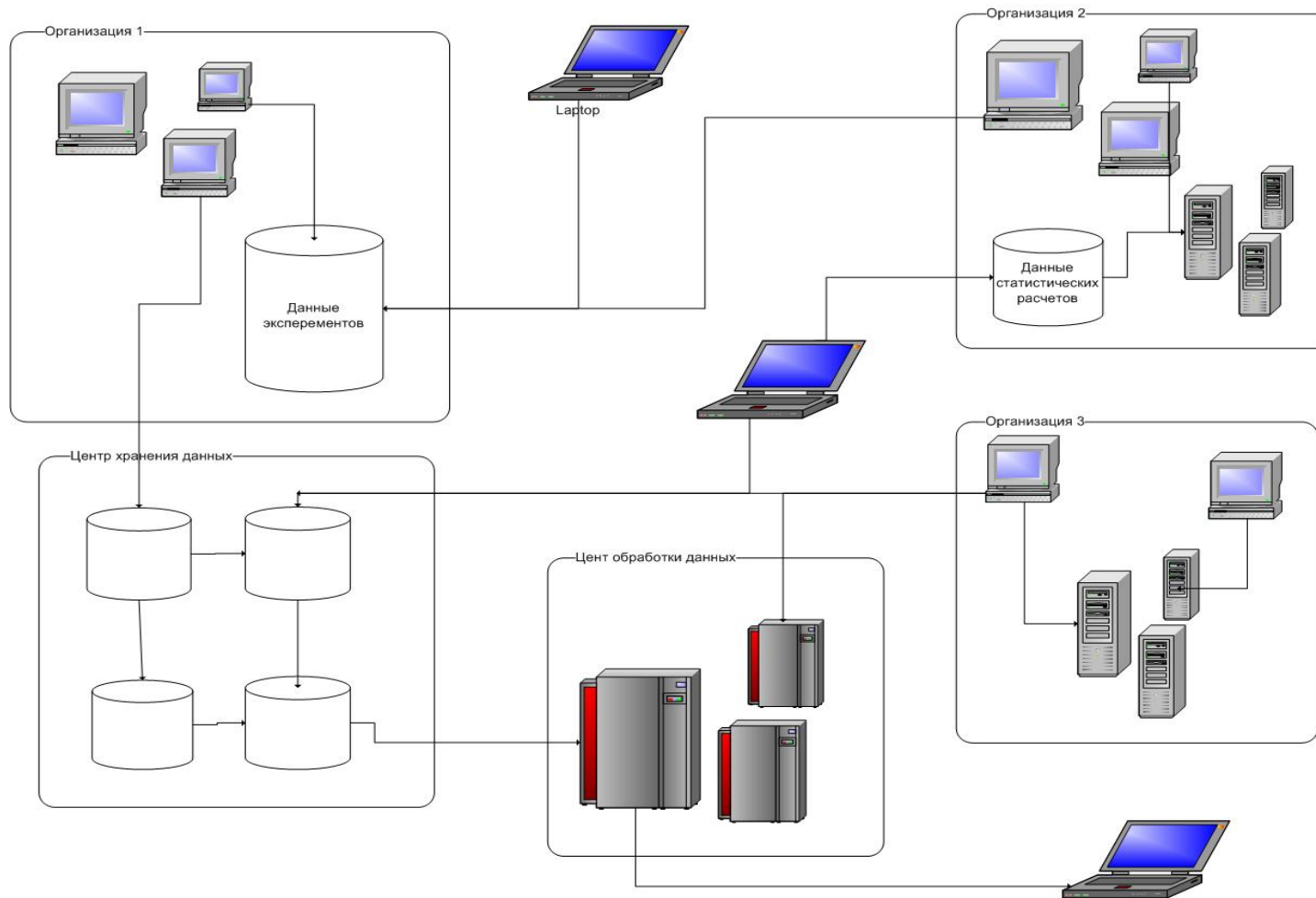
- Задачи систем информационного обеспечения распределенных научных исследований
- Тематический поиск информации
- Технологии построения специализированных информационных систем для проведения распределенных научных исследований

# Особенности строго структурированных данных

- Строгость описания
- Автоматический контроль целостности и корректности
- Простота и высокая скорость доступа при автоматической обработке
- -----
- Невозможность создания «универсальных» структур данных
- Необходимость создания специальных информационных систем под каждую конкретную задачу



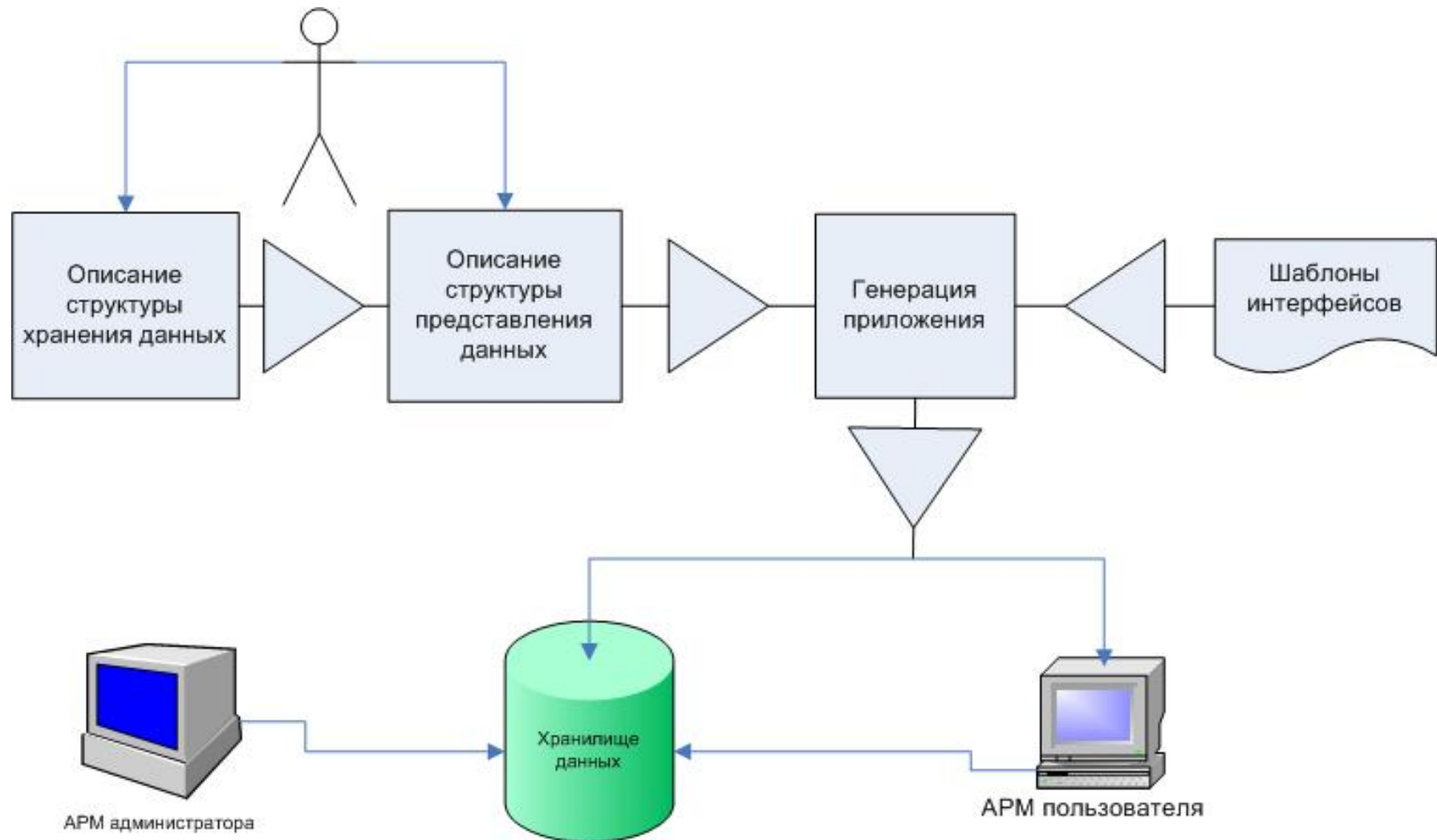
# Использование справочных информационных систем



# Требование к информационным системам

- Распределенный доступ к данным
- Разделение привилегий для доступа к данным
- Удобный интерфейс (поиск, сортировка, фильтрация, создание отчетов, экспорт данных)
- Быстрое создание и модификация структуры пользовательских данных

# Технология быстрого создания справочных информационных систем



# Паспорта материалов(1)

информационная система "Паспорт", User : (1)

Файл Правка Модули Данные Операции Окно ?

Найти  Все

Шаблон паспорта

Название шаблона
Металлы
Старинные пушки

Характеристики шаблона паспорта

Название	Номер	Количество
Назначение по видам	6	1
Химический состав	4	10
Исходные материалы	5	1
Рекомендуемые обле	7	1
Какие сплавы могут б	8	1
Модуль упругости пр	11	10

# Паспорта материалов(2)

информационная система "Паспорт", User : (1)

Файл Правка Модули Данные Операции Окно ?

Найти  Все

Материал

Название
Алюминий
Пушка
титан

Марка

Название марки
ВК-1251

Незаполненные характеристики

Номер	Название	Значение	Значение2	Значение3
4	Химический состав	5.8	7	Al
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
4	Химический состав			
5	Исходные материалы			
6	Назначение по видам			
7	Рекомендуемые объемы			
8	Какие сплавы могут использоваться			
1	Модуль упругости при растяжении	33	54	
1	Модуль упругости при сжатии	2	6	
1	Модуль упругости при изгибе	4	44	
1	Модуль упругости при кручении	4	8	
1	Модуль упругости при ударном изгибе			
1	Модуль упругости при ударном кручении			
1	Модуль упругости при ударном изгибе			
1	Модуль упругости при ударном кручении			

Паспорт

Номер паспорта	Название шаблона	Название источника
1111111	Металлы	???

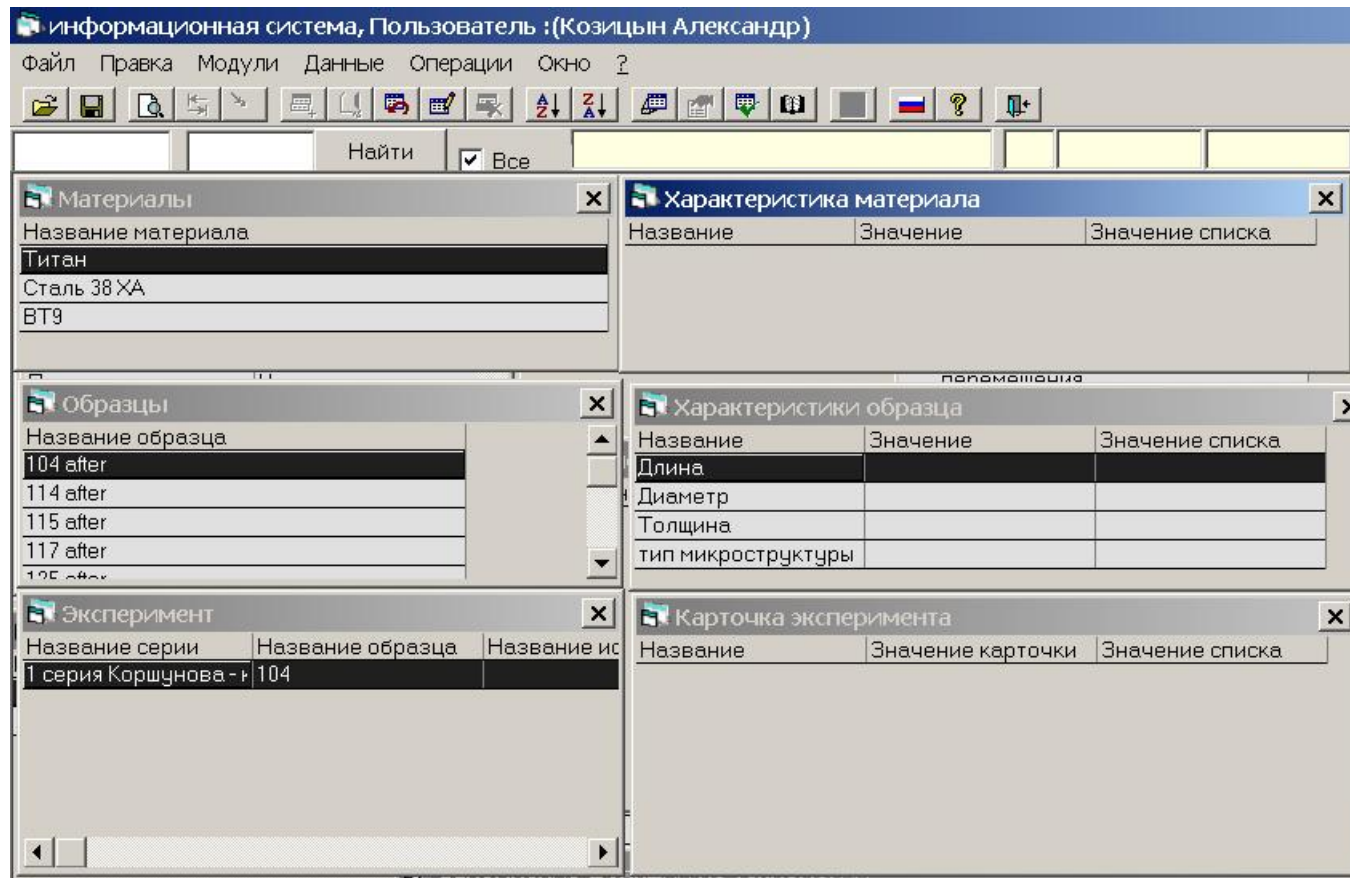
5 — 6.9

5.95

Ni

Сохранить Отменить

# Система сопровождения экспериментов(1)



# Система сопровождения экспериментов(2)



Файл Правка Модули Данные Операции Окно ?

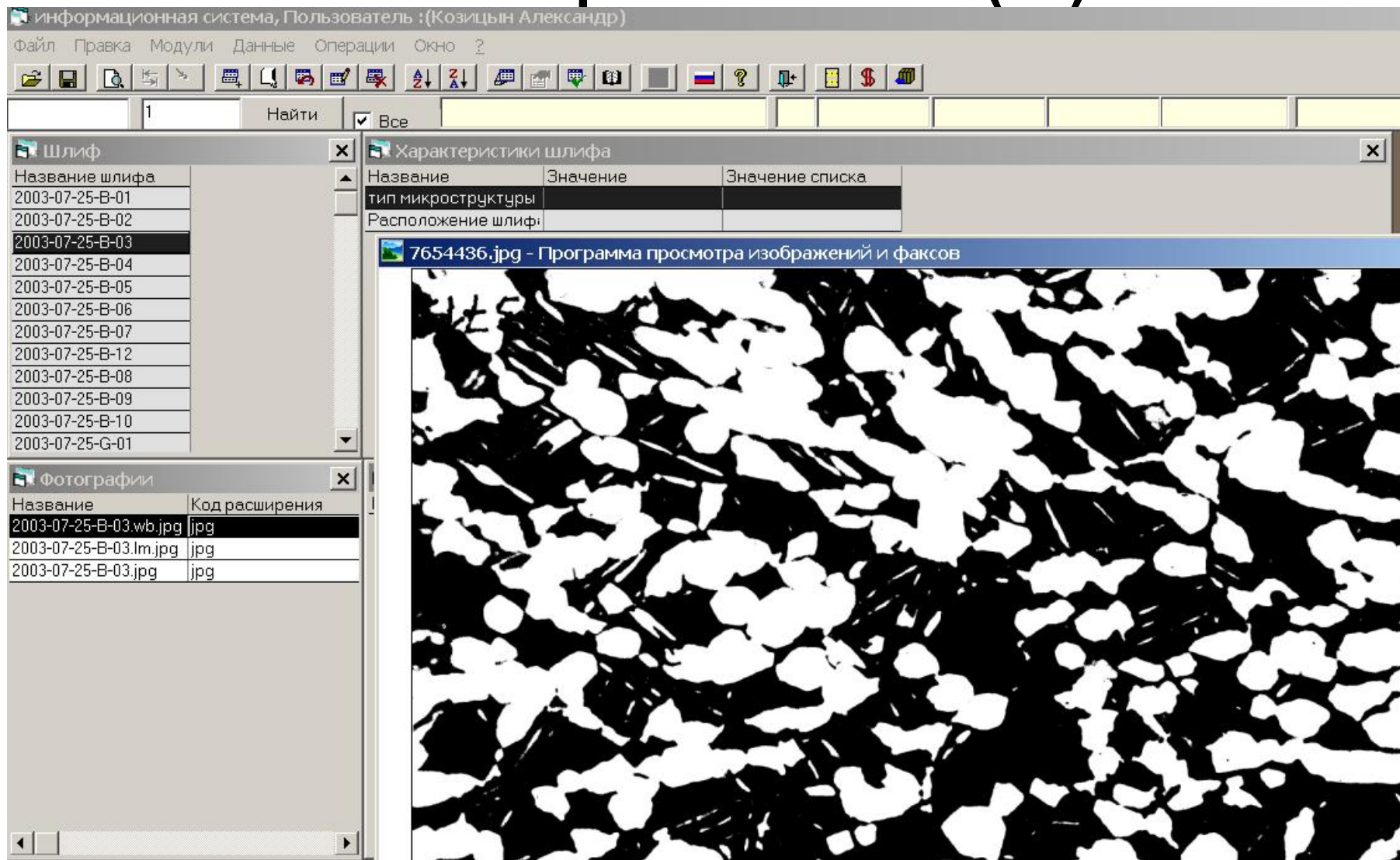
Найти  Все

II-A3 Однокомпонентный M

t	dL	dR	M
0	0.		0.
1	1.667		0.1154901960
2	3.333		0.2421568627
3	5.		0.4656862745
4	6.667		0.6109803921
5	8.333		0.7450980392
6	10.		0.8382352941
7	11.667		0.9052941176
8	13.333		0.9462745098
9	15.		0.9611764705
10	16.667		0.9909803921
11	18.333		1.0021568627
12	20.		1.0170588235
13	21.667		1.0282352941
14	23.333		1.0319607843
15	25.		1.0468627450
16	26.667		1.0580392156
17	28.333		1.0617647058
18	30.		1.0692156862
19	31.667		1.0766666666
20	33.333		1.0841176470
21	35.		1.0915686274
22	36.667		1.0990196078
23	38.333		1.1064705882
24	40.		1.1101960784
25	41.667		1.1176470588

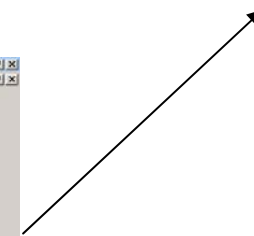
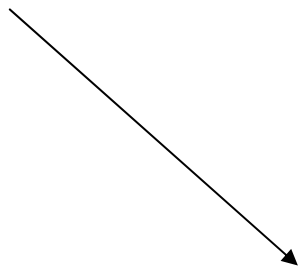
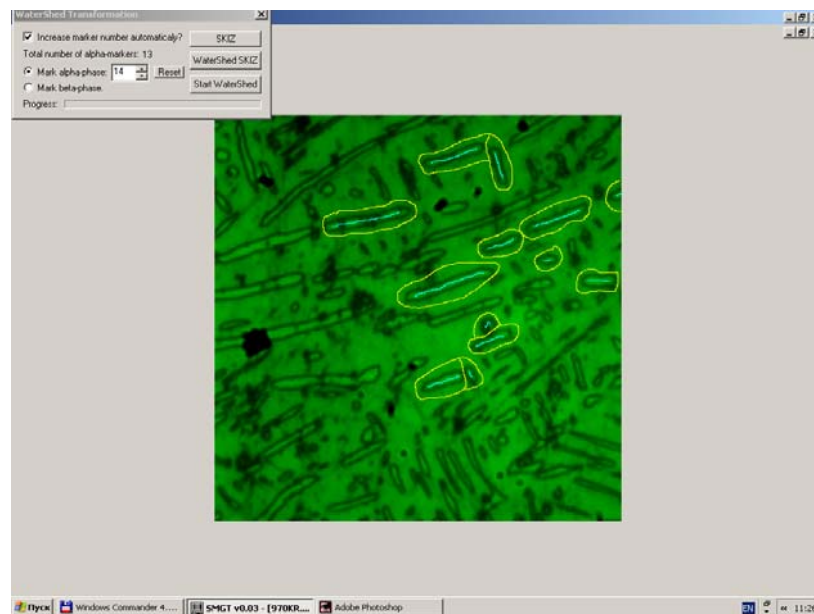
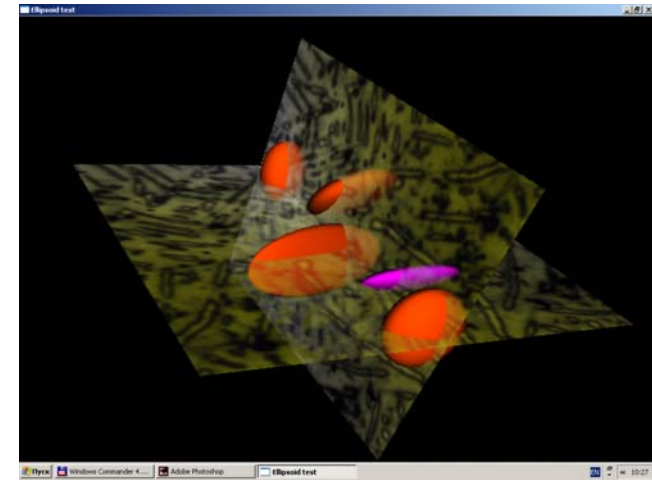
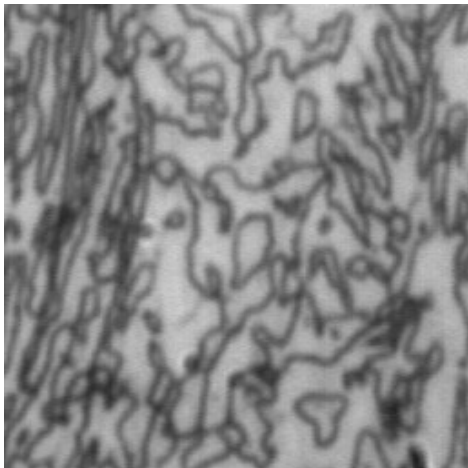


# Система сопровождения экспериментов(3)



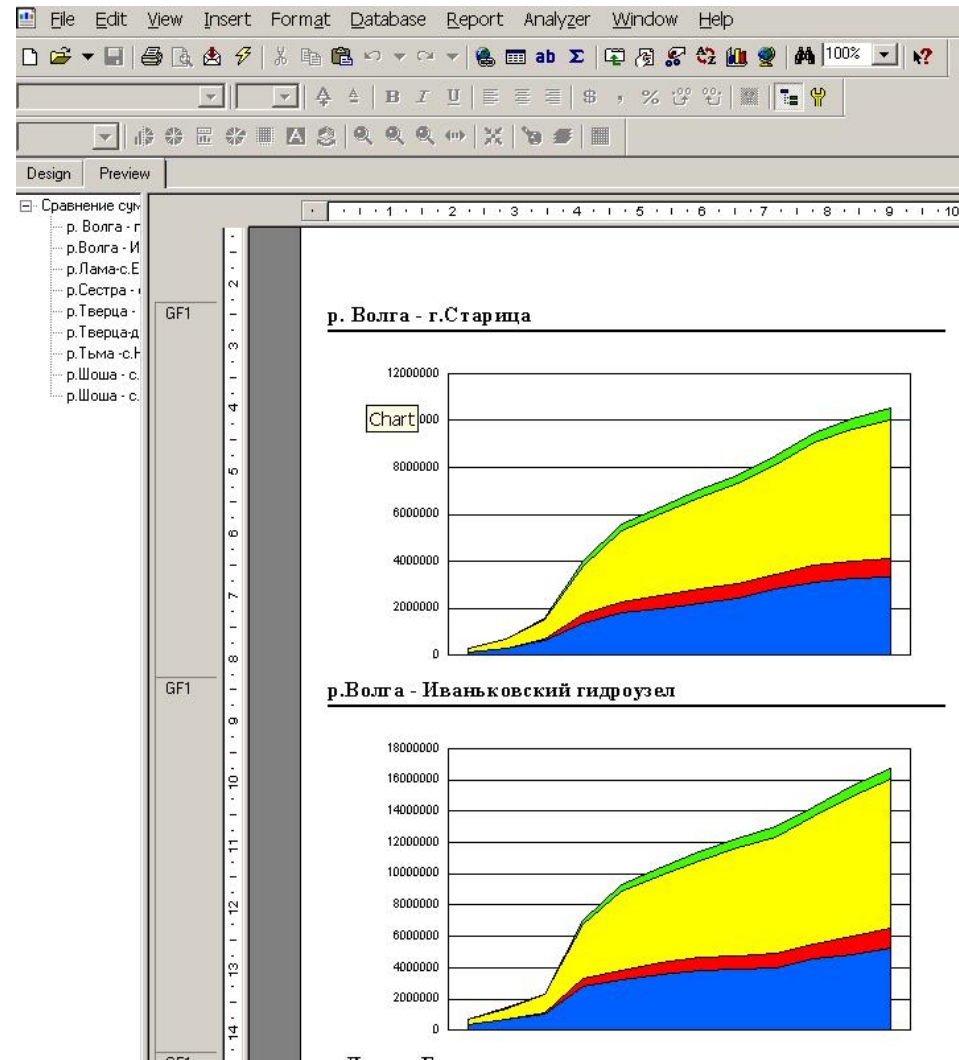


# Автоматическая обработка данных

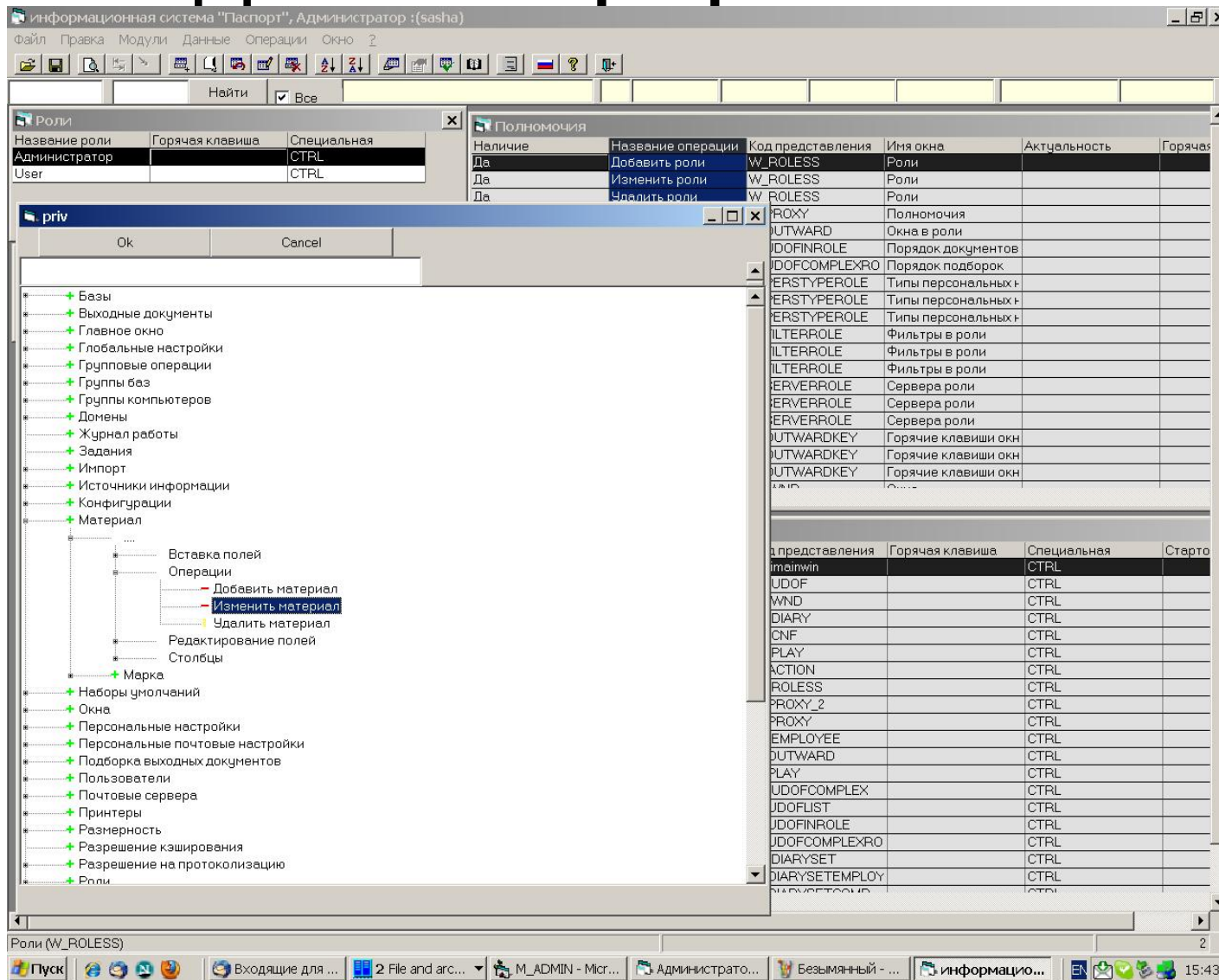


# Аналитические отчеты

		р.Волга - Иваньковск ий гидроузел	Total
937	1	0.00	0.00
	2	0.00	0.00
	3	0.00	0.00
	4	670 000.00	670 000.00
	5	231 000.00	231 000.00
	6	74 000.00	74 000.00
	7	55 000.00	55 000.00
	8	141 000.00	141 000.00
	9	132 000.00	132 000.00
	10	95 000.00	95 000.00
	11	67 000.00	67 000.00
	12	50 000.00	50 000.00



# Подсистема администрирования



# Особенности технологии

- Простота разработки
- Визуальное представление структур данных
- Простота модификации
- Наличие шаблонов интерфейсов
- Наличие готового модуля администрирования
- -----
- Создание только справочных систем

**Спасибо за внимание**