

Архитектура автоматизированной системы сбора статистики событий в распределенной информационной системе ZooSPACE.

Жижимов О.Л. *, Турчановский И.Ю. *, Панышин А.А. **, Чудинов С.А. **

* Институт вычислительных технологий СО РАН, Новосибирск, tur@hcei.tsc.ru
* Институт сильноточной электроники СО РАН, Томск

В работе описана система сбора информации о событиях, происходящих в распределенной информационной системе ZooSPACE. Представлена архитектура автоматизированной системы ZooSTAT. В работе представлены алгоритмы анализа информационного потока, а так же описана структура и способы хранения данных. Описана модульная структура ZooSTAT, с возможностью визуализации результатов статистического анализа.

Введение

Целью работы является создание автоматизированной системы обработки больших объемов данных ZooSTAT генерируемых в узлах платформы массовой интеграции ZooSPACE [1], которая отражает работоспособность распределенной информационной системы (РИС). Разрабатываемая система ZooSTAT применяет сквозную технологию сбора, обработки, хранения и предоставления статистической информации, и позволяет на запрос администратора РИС визуализировать статистическую информацию о коннективности узлов РИС, отдельных хранилищ данных, а также характер запросов клиентов к распределенной информационной системе ZooSPACE.

В соответствии с архитектурой РИС клиенты формирует запросы к узлам ZooSPACE. В результате обработки этих запросов на извлечение информации формируется отдельные потоки событий для каждого узла системы, которые накапливаются в специализированных журналах узлов ZooSPACE. Для систематизации и обработки этой информации разрабатывается система ZooSTAT.

Работа выполняется в рамках технического задания по Государственному контракту № 07.514.11.4130 от 6.06.2012 по теме: «Разработка принципов и программных средств виртуальной интеграции распределённых источников данных на основе международных стандартов для создания масштабных информационных инфраструктур» (шифр заявки «2012-1.4-07-514-0022-004»).

Функциональность системы.

Разрабатываемая система ZooSTAT предполагает сбор, обработку, хранение событий РИС, а также формирование отчетов системные функции по администрированию системы и резервное копирование базы данных ZooSTAT.

Архитектура системы ZooSTAT.

Система ZooSTAT разработана на основе модульной структуры и содержит следующие модули (Рис. 1.):

- модуль сбора информации о событиях с серверов ZooSPACE (лог-данных);
- модуль генератора статистик;
- модуль подготовки шаблонов для генератора статистик;
- модуль вывода обработанной статистической информации через web-интерфейс;
- модуль резервного копирования баз данных.

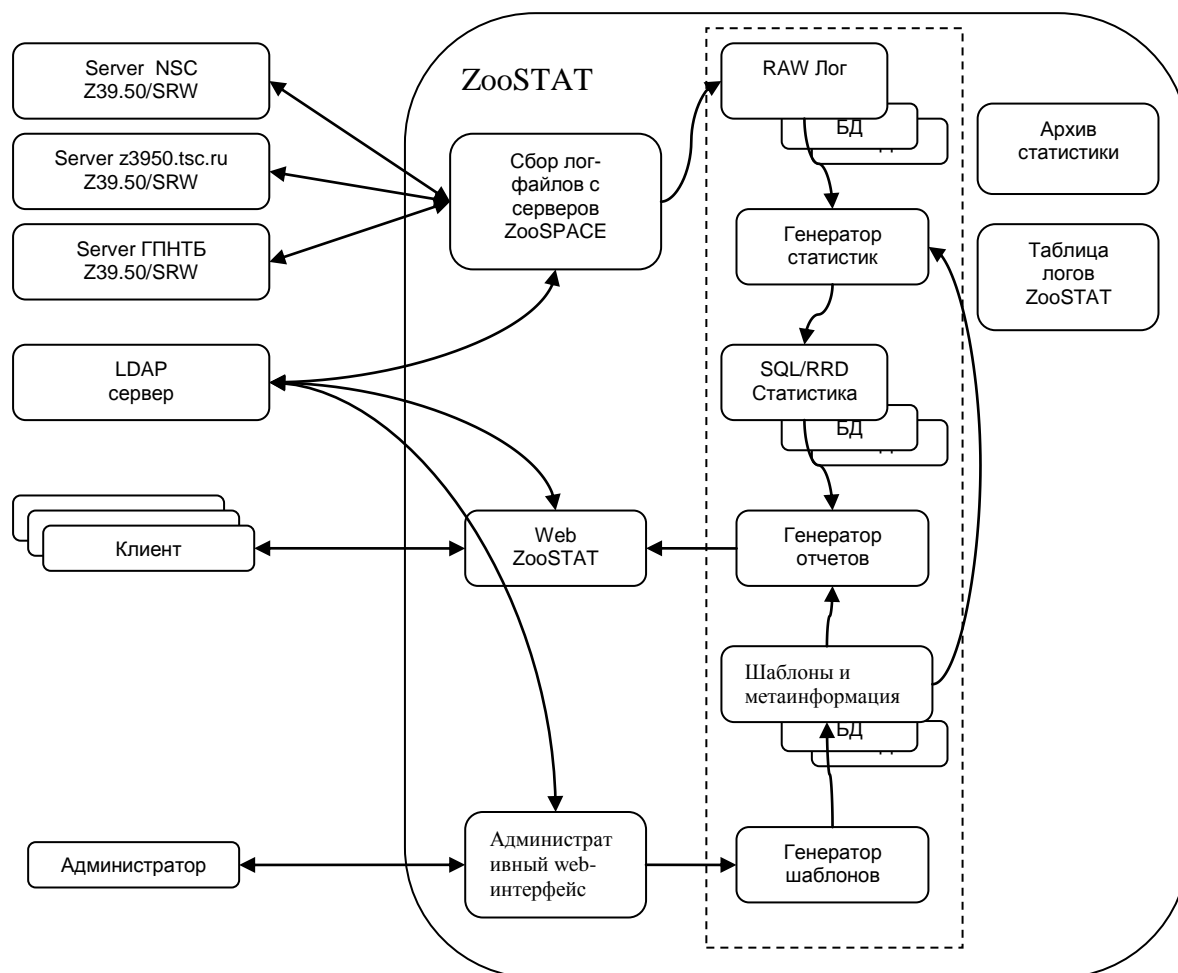


Рис. 1. Архитектура системы сбора статистики ZooSTAT.

Модуль сбора информации.

Модуль сбора информации о событиях с серверов ZooSPACE последовательно опрашивает сервера, указанные в конфигурационном списке, по протоколу Z39.50/SRW, выгружает с них лог-данные и выполняет их парсинг. Результат разбора сохраняется в реляционной базе данных (MySQL или PostgreSQL) как промежуточное представление информации, так называемый “сырой лог” – набор связанных и проиндексированных таблиц реляционной БД, структура которых позволяет выполнять быстрое построение отчета (Приложение №1 – структура таблиц “сырого лога”).

Архитектура модуля сбора информации.

Модуль сбора информации о событиях на серверах ZooSPACE состоит из следующих блоков (Рис. 2.):

- блок считывания конфигурационной информации из LDAP сервера или файла конфигурации, в случае, если LDAP сервер не доступен;
- интерфейсный блок доступа к серверам ZooSPACE по протоколу Z39.50/SRW;
- блок аутентификации, запрашивающий параметры авторизации на сервере LDAP для доступа к серверам ZooSPACE;
- блок загрузки данных с серверов ZooSPACE по протоколу Z39.50/SRW;
- блок обработки данных и сохранения распарсенных данных в таблицах реляционной БД.

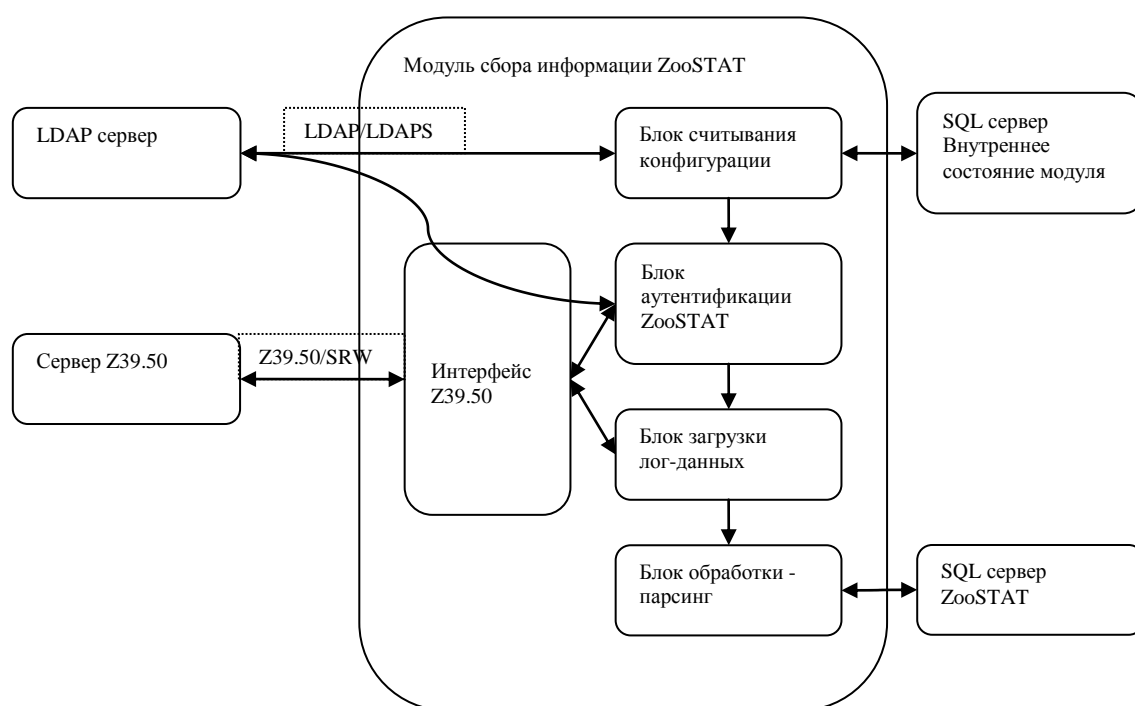


Рис. 2. Архитектура модуля сбора информации о событиях с серверов ZooSPACE.

Алгоритм функционирования модуля сбора информации.

Конфигурационные параметры модуля хранятся в LDAP каталоге, где указываются список IP-адресов опрашиваемых серверов ZooSPACE (IP/порт), аутентификационные данные для доступа к серверам, периодичность запуска модуля опроса. При отсутствии доступа к LDAP каталогу используются параметры, заданные в текстовом файле конфигурации.

Модуль хранит переменные внутреннего состояния в таблице реляционной БД. При инициализации из этой таблицы считываются для каждого из серверов: время последней обработанной записи, время доступа к серверу, количество ошибок при передаче и т.п. Список сортируется по времени доступа к серверу.

Далее, для каждого сервера ZooSPACE из списка, последовательно выполняются следующие процедуры (функциональная схема Рис. 3.):

- из LDAP каталога считываются аутентификационные параметры для доступа к серверу Z39.50;

- подключение к серверу по протоколу Z39.50/SRW;
- считывание логов сервера с определенного момента времени;
- парсинг и загрузка “сырых логов” в таблицы реляционной БД (Приложение №1);
- сохранение внутреннего состояния для текущего сервера в таблице БД.

Указанная выше последовательность процедур выполняется для каждого сервера.

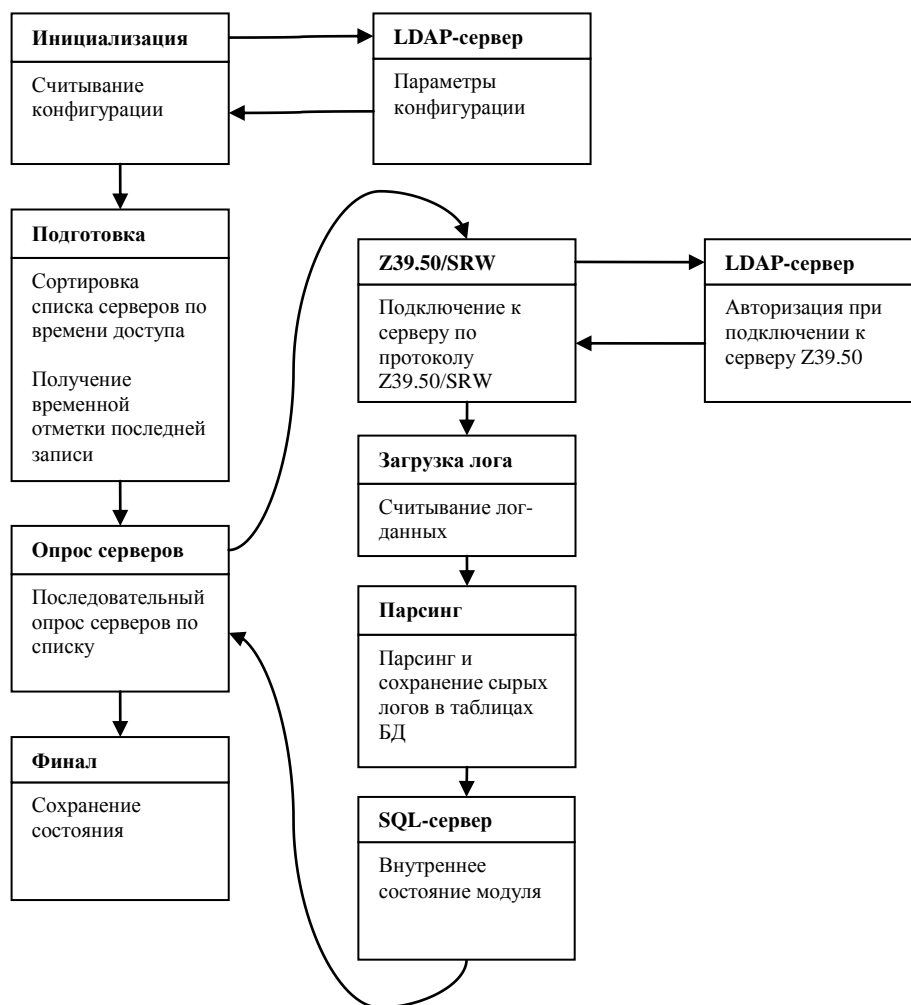


Рис. 3. Функционирование модуля сбора информации

Модуль генерации статистик событий на серверах ZooSPACE.

Модуль генератора статистики предназначен для преобразования данных, сохраненных в таблицах “сырых логов”, в соответствии с правилами, заданными администратором системы, позволяющими выделить существенную информацию о событиях в системе ZooSPACE из потока данных (“сырых логов”).

Архитектура модуля генерации статистик.

Сохраненные в таблицах реляционной БД “сырые логи” серверов обрабатываются модулем формирования статистик. Нами были опробованы различные структуры для хранения результатов статобработки (RRD/SQL/XML). Однако в текущей версии мы

склоняемся к хранению результатов в таблицах реляционной БД. Это связано в первую очередь с более универсальным подходом при формировании отчетов в сочетании с высокой скоростью выборки. Следует отметить, что RRD базы данных также заслуживают внимания, благодаря чрезвычайно высокой скорости работы, но на текущий момент, при концептуальной разработке системы ZooSTAT, желательно иметь универсальный инструмент формирования отчета, поэтому мы использовали SQL сервер.

Модуль генерации статистической информации (Рис. 4.) загружает из таблиц реляционной БД подготовленные шаблоны с соответствующей им метainформацией, описывающей условия применения шаблона и генерируемые базы статистик. Выполняет обработку “сырых логов” и генерирует необходимые для визуализации информационные структуры.



Рис. 4. Архитектура модуля генерации статистической информации.

Алгоритм функционирования модуля генерации статистик.

Модуль генератора статистических данных способен формировать таблицы реляционной БД для следующих отчетов: количество событий в единицу времени; отсортированный список – частота событий для указанного параметра. Алгоритм формирования результата статистической обработки достаточно универсален и может быть применен для широкого класса статистик.

Функциональная схема генератора статистик приведена на рис. 5.

Модуль считывает из таблицы БД список шаблонов и последовательно их обрабатывает, формируя таблицы отчетов, подготовленные для вывода.

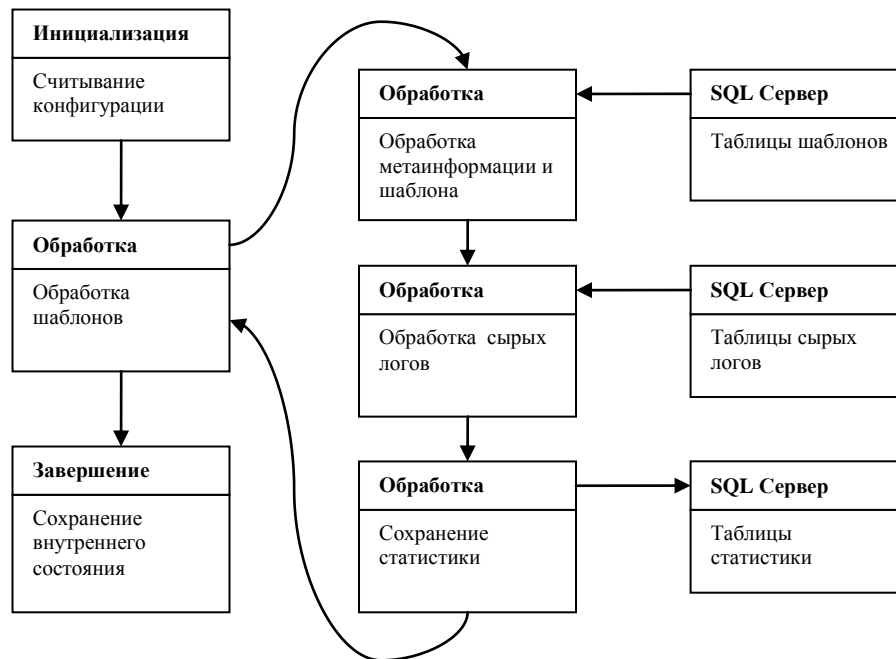


Рис. 5. Функционирование модуля генерации статистической информации.

Модуль подготовки шаблонов генератора статистики событий на серверах ZooSPACE.

Модуль предназначен для подготовки шаблонов и соответствующей им метайнформации описания, проверки их синтаксической и семантической корректности.

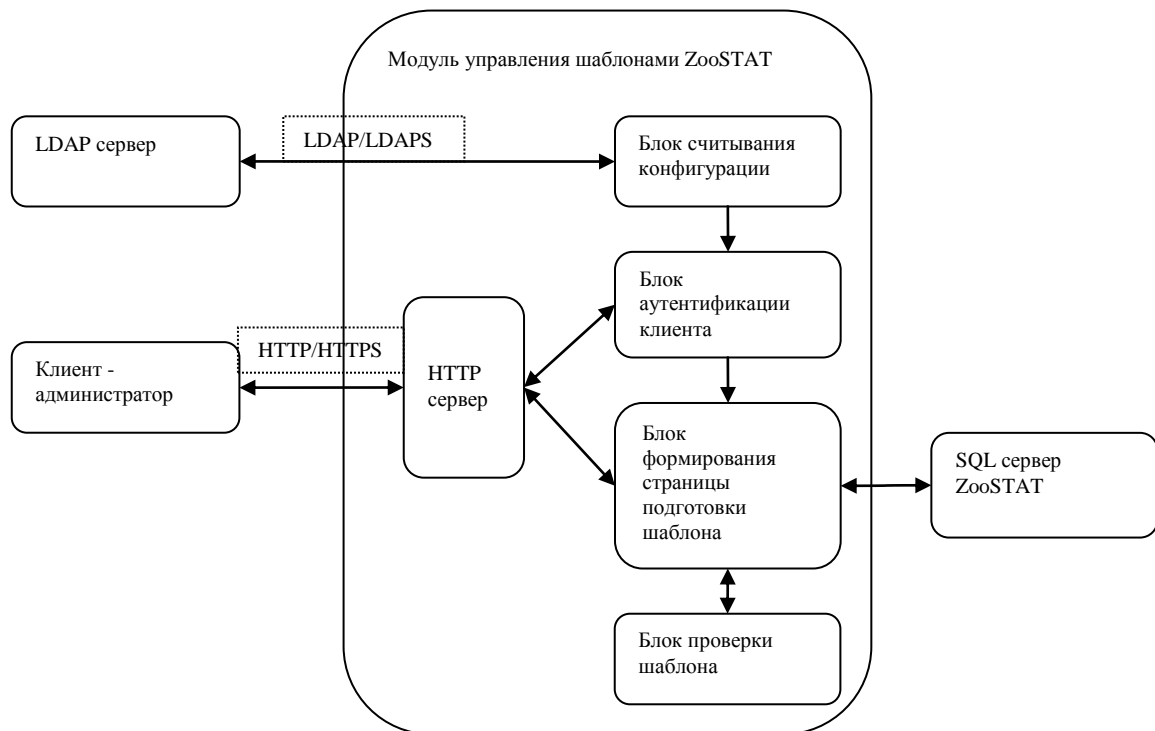


Рис. 6. Архитектура модуля подготовки шаблонов для генератора статистик событий на серверах ZooSPACE

Архитектура модуля подготовки шаблонов.

Модуль представляет собой web-интерфейс администратора, разработанный на языке программирования PHP и состоящий из следующих блоков (Рис. 6.):

- блок считывания конфигурационной информации из LDAP сервера или файла конфигурации, в случае, если LDAP сервер не доступен;
- web-сервер, для взаимодействия с клиентом-администратором;
- блок аутентификации клиента;
- блок формирования страницы подготовки шаблона;
- блок проверки корректности шаблона.

Модуль подготовки шаблонов позволяет генерировать статистики 2-х типов:

- кол-во событий в единицу времени, отображаемых в виде графиков;
- список событий, зарегистрированных для фиксированного параметра, отсортированный по кол-ву событий, отображаемых в виде таблицы или круговых диаграмм.

Статистики первого типа могут быть наглядно представлены в виде графиков, с временной координатой. Например, в качестве такой статистики можно предложить “число запросов к конкретному серверу в течении последнего месяца”, “число ошибочных запросов к серверу в течении дня” и т.п. При формировании таблицы для статистики первого типа указывается логическое выражение, например, (APDU:common: typeAPDU = 1), где число 1 – числовой идентификатор типа APDU – initRequest. Агрегирование по времени выполняется с фиксированной величиной 1 минута. Величина агрегации – 1 минута – может быть изменена в конфигурации. Вероятнее всего, в рабочем проекте для уменьшения размера таблиц эту величину можно установить 10 или 30 минут. Для указанного выше выражения генератор сформирует таблицу реляционной БД, в первом поле которой – отметка времени с шагом в 1 минуту, вот втором – число событий, удовлетворяющих указанному выражению за этот промежуток времени. При просмотре страницы статистики пользователем возможна любая агрегация по времени, кратная времени агрегации. Пользователем могут указываться дополнительные условия для генерации статистики (например, сервер, база данных).

Статистики второго типа могут представляться в виде списка. Примером может служить таблица “Количество обращений к серверу с определенного IP за последний месяц”, отсортированная по кол-ву обращений. Вывод информации для конечного пользователя может быть либо в виде списка, например “Топ 50”, либо в виде круговых диаграмм.

Сгенерированный шаблон сохраняется в БД вместе с метаинформацией, описывающей условия применения и использования статистики.

Алгоритм функционирования модуля подготовки шаблонов.

Конфигурационные параметры модуля хранятся в LDAP каталоге. При отсутствии доступа к LDAP каталогу используются параметры, заданные в текстовом файле конфигурации.

Для получения доступа к страницам сервиса требуется авторизация клиента-администратора. Проверка пароля и прав доступа выполняется через запрос к LDAP каталогу.

После успешного входа в систему администратору предоставляется возможность просмотреть уже имеющийся список шаблонов, создать новый или модифицировать уже имеющийся шаблон.

Функциональная схема модуля представлена на рис. 7.

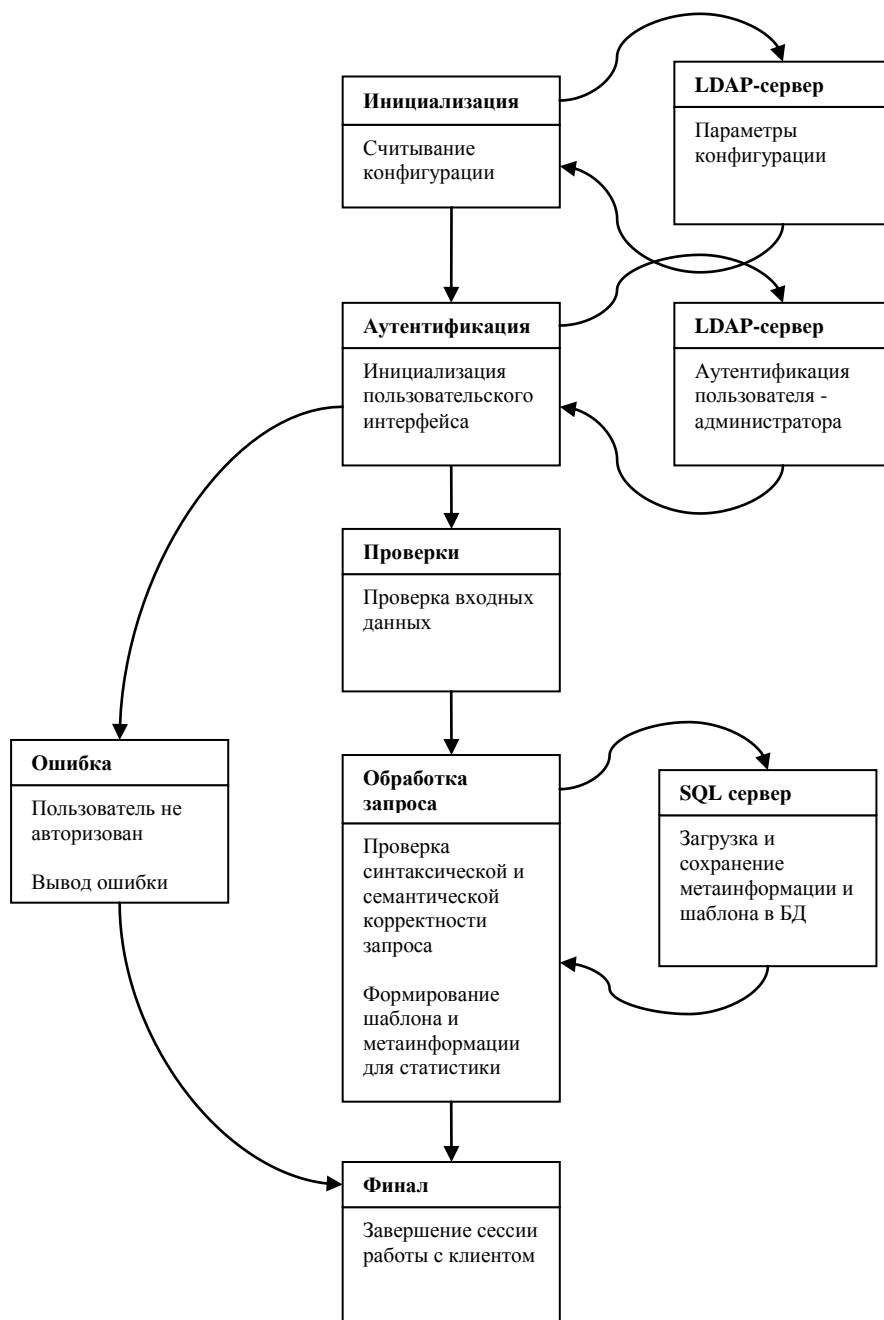


Рис. 7. Функционирование модуля подготовки шаблонов для генератора статистик событий на серверах ZooSPACE

Пример страницы генерации нового шаблона “Число запросов ко всем серверам в единицу времени” показан на рис. 8.

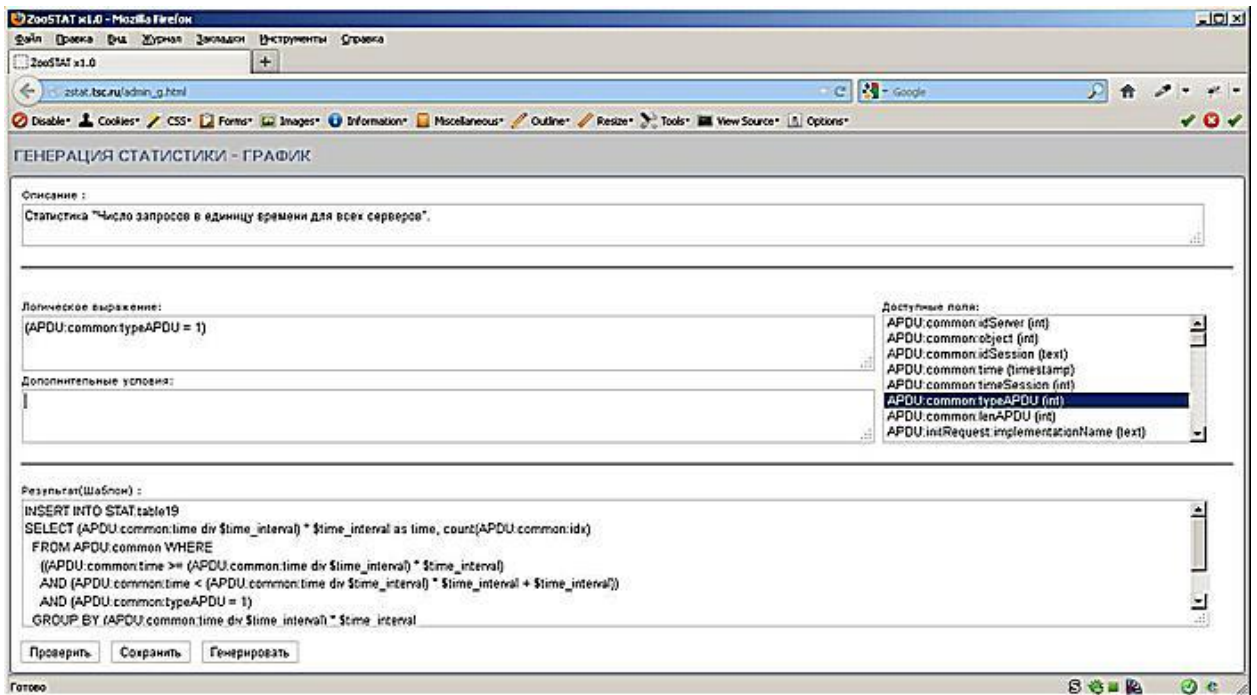


Рис. 8. Пример формирования шаблона для статистики “Число запросов ко всем серверам в единицу времени”.

Модуль вывода статистической информации.

Модуль вывода предназначен для предоставления web-доступа к статистической информации конечных пользователей. И позволяет в удобной, доступной форме получить информацию о событиях в системе ZooSPACE.

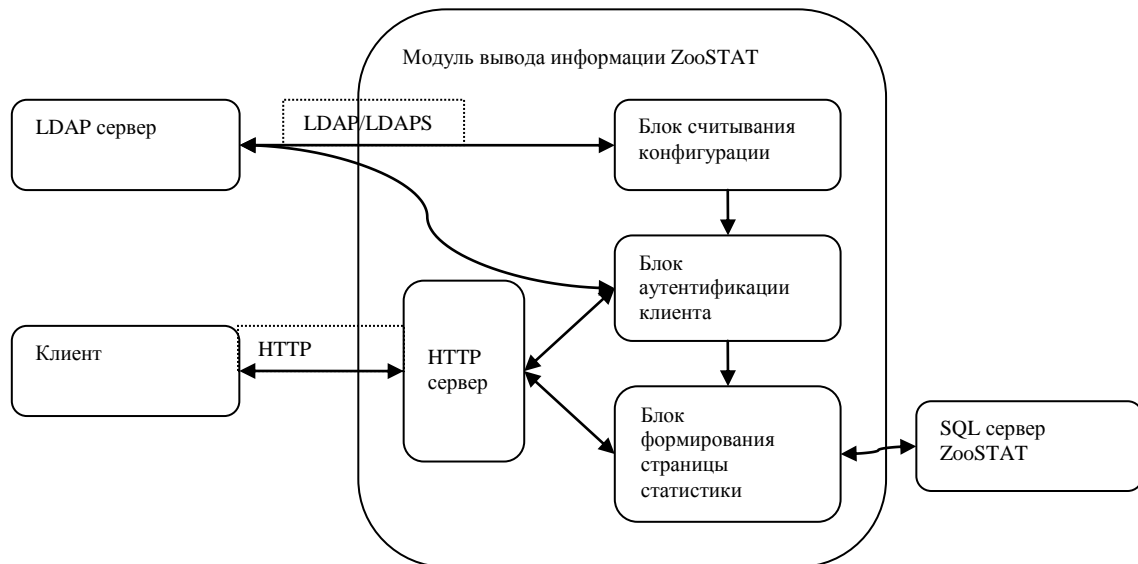


Рис. 9. Архитектура модуля вывода статистической информации о событиях на серверах ZooSPACE.

Архитектура модуля вывода статистической информации.

Модуль вывода статистической информации о событиях на серверах ZooSPACE состоит из следующих блоков (Рис. 9.):

- блок считывания конфигурационной информации из LDAP сервера или файла конфигурации, в случае, если LDAP сервер не доступен;
 - web-сервер, для взаимодействия с клиентом;
 - блок аутентификации клиента;
 - блок формирования страницы статистики.
- Модуль реализован на PHP.

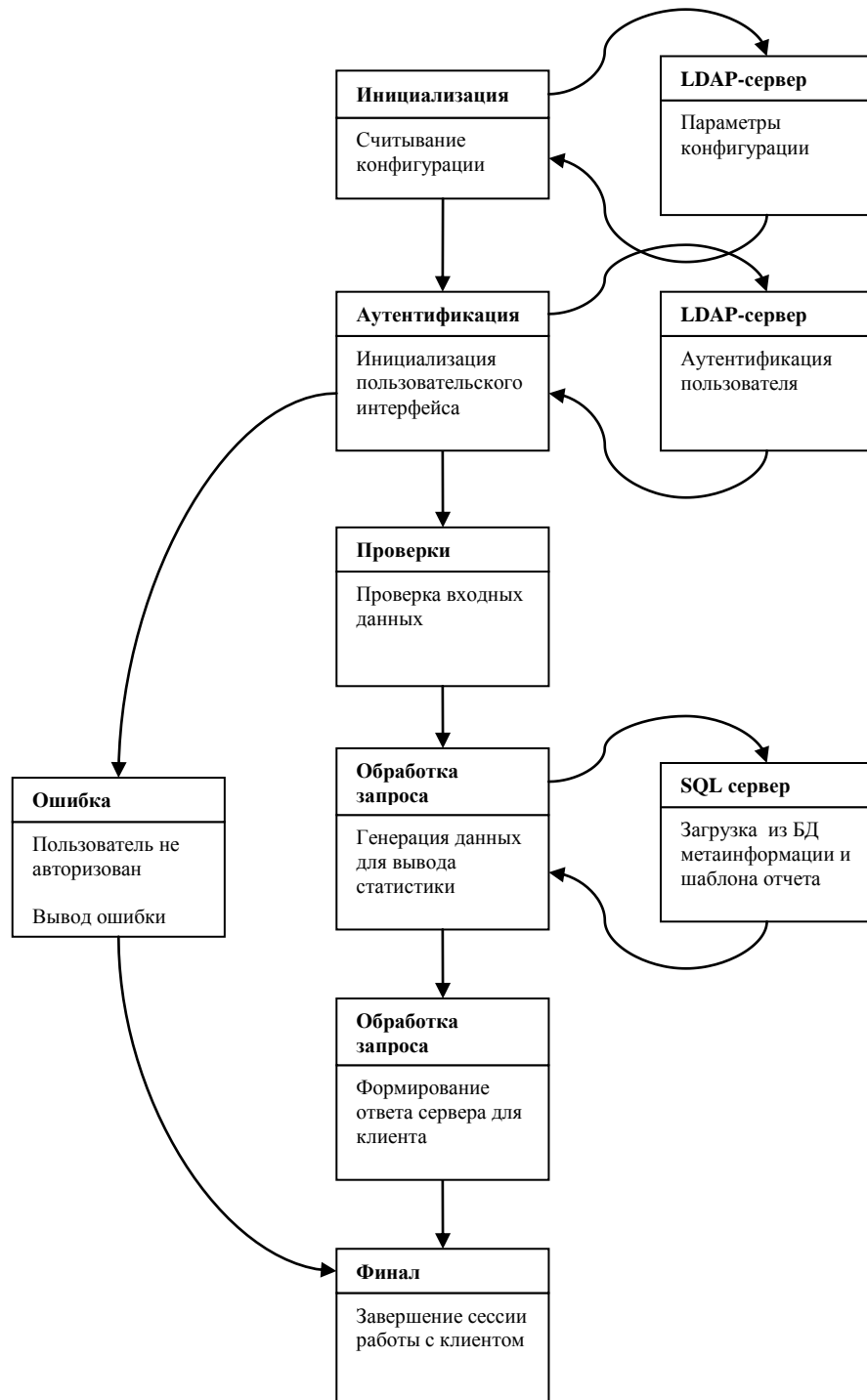


Рис. 10. Функционирование модуля вывода статистической информации.

Алгоритм функционирования модуля вывода статистической информации.

Конфигурационные параметры модуля хранятся в LDAP каталоге. При отсутствии доступа к LDAP каталогу используются параметры, заданные в текстовом файле конфигурации.

Для получения доступа к страницам сервиса требуется авторизация клиента. Проверка пароля и прав доступа выполняется через запрос к LDAP каталогу.

После успешного входа в систему пользователю предоставляется возможность посмотреть накопленные статистические данные по любому из серверов ZooSPACE, например, количество запросов к серверу в единицу времени, количество удачных/неудачных запросов в единицу времени, список наиболее активных клиентов, минимальное, среднее, максимальное время сессии и т.д. Однако список возможных статистик определяется администратором системы, который сформировал соответствующие шаблоны и условия для модуля генератора отчетов.

На рис. 10. показана функциональная схема модуля вывода. Генерация страницы, с запрошенной клиентом статистикой, выполняется на основе метаинформации, описывающей шаблон, условия применимости и имени таблицы, содержащей подготовленные для вывода данные.

При формировании web-страницы, из БД считывается метаинформация и шаблон запроса к реляционной БД, выполняются необходимые подстановки и проверки, после чего формируется запрос к сгенерированной заранее таблице, содержащей статистические данные. Результат выводится пользователю в виде графика, круговой диаграммы или списка. На рис. 11. приведен пример генерируемой модулем вывода статистики “Число запросов к серверу в минуту”.

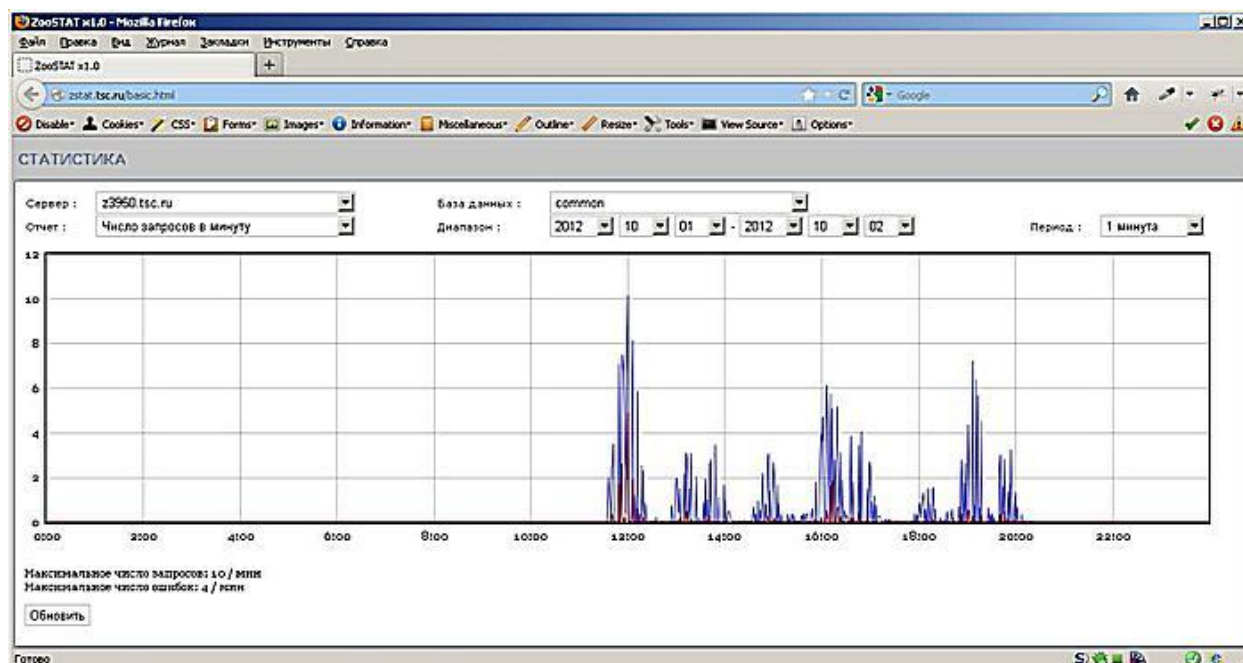


Рис. 11. Пример вывода статистической информации “Число запросов в минуту” работы тестового сервера z3950.tsc.ru в течении суток. Синим цветом – число запросов к серверу

в единицу времени, красным – число запросов с ошибками в единицу времени. Период агрегации – 1 минута.

Модуль резервного копирования.

Модуль резервного копирования предназначен для сохранения дампов баз данных конфигураций модулей ZooSTAT, шаблонов статистики и сырых логов. Информация, сохраняемая в бэкапе, позволяет полностью восстановить работоспособность системы при возможных сбоях оборудования.

Архитектура

Модуль резервного копирования (Рис. 12.) запускается по крону и выполняет сохранение дампов БД конфигураций модулей, шаблонов статистики и сырых логов. Базы данных содержащие результаты статобработки не сохраняются в бэкапе. После восстановления конфигураций модулей и шаблонов из бэкапа, модуль генератора статистики самостоятельно сформирует все необходимые таблицы реляционной БД. Периодичность сохранения и список сохраняемых таблиц указываются в LDAP каталоге и конфигурационном файле модуля. Для уменьшения объема занимаемого бэкапами дискового пространства используется системный вызов утилиты bzip2.

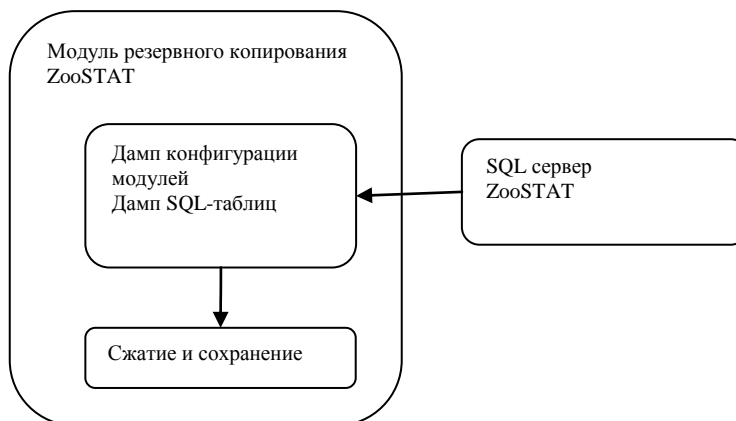


Рис. 12. Архитектура модуля резервного копирования.

Алгоритм функционирования

На рис 13. приводится алгоритм функционирования модуля резервного копирования.

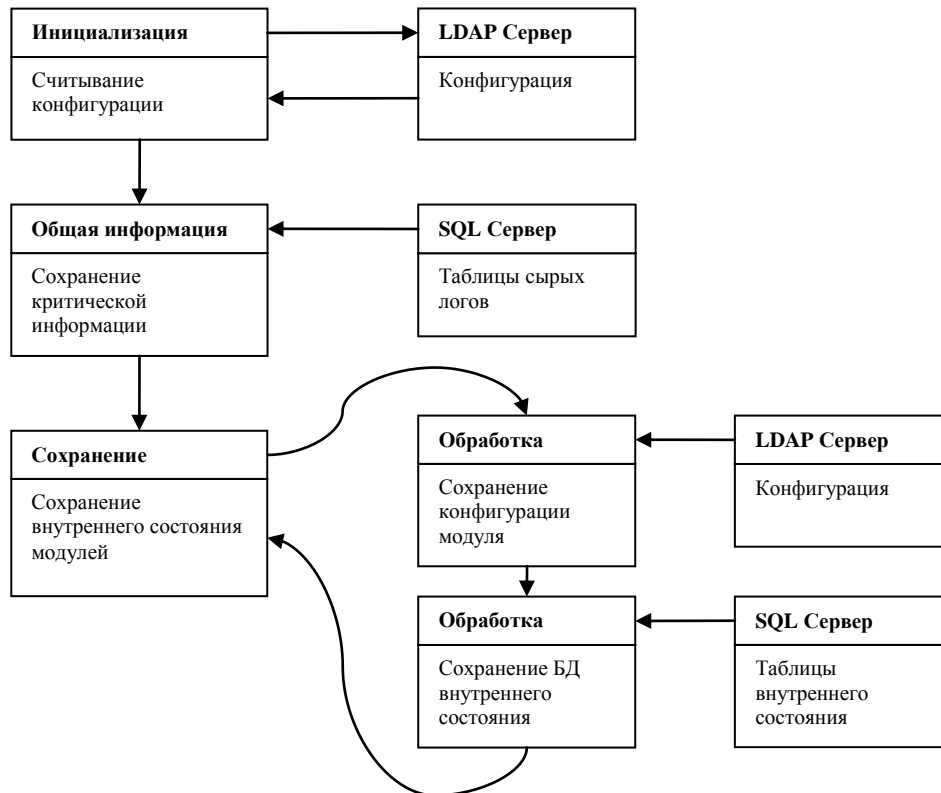


Рис 13. алгоритм функционирования модуля резервного копирования.

Литература.