

1. Введение

В настоящий момент существуют некоторые проблемы организации структурированного хранения экспериментальных данных в небольших научных группах, ведущих исследования в области естественных науки (физика, химия, биология, астрономия и т. д.). Проблемы можно условно разбить на три категории:

1. **Локальность.** Данные хранятся локально, без возможности какого-либо внешнего доступа к ним. Пример: результаты эксперимента хранятся на носителе одного из членов научной группы, в то время как другие члены группы работают рассредоточено, в результате чего передача данных может быть реализована только через личную встречу или через пересылку данных, непрозрачную для остальных членов группы.
2. **Неструктурированность.** Зачастую даже в пределах одной научной группы не существует жестко заданного формата хранения метаданных об эксперименте (протокола эксперимента), в результате чего возможна неправильная трактовка результатов эксперимента, а также невозможно его воспроизведение. Кроме того, такой подход к хранению метаданных затрудняет подключение новых участников к научной работе.
3. **Потеря актуальности.** В результате первых двух пунктов появляется проблема потери актуальности данных, так как при рассредоточенной работе нет возможности получать последние данные об эксперименте с малой задержкой.

2. Постановка задачи

Требуется разработать систему, предназначенную для хранения и первичной обработки данных экспериментальных исследований в области естественных наук.

Система должна предоставлять следующие функции:

- Возможность накопления и сортировки экспериментальных данных. Под экспериментальными данными подразумеваются данные, полученные с экспериментальной установки (далее сырые данные), а также исчерпывающие метаданные, необходимые для их интерпретации (далее метаданные).
- Возможность управления форматами экспериментальных данных. Сырые данные могут представлять собой как файлы определенного расширения, так и бинарные файлы. Метаданные могут состоять из произвольной информации, то есть необходимо исключить строгие (например, табличное представление) ограничения на формат метаданных.
- Возможность написания и "горячей" установки отдельных модулей, позволяющих проводить парсинг сырых данных, расчет необходимых метрик и их графическое представление (далее парсеров). Модули должны быть полностью обособлены от самой системы.
- Возможность генерации простейших отчетов.
- Возможность экспорта данных для их автономного перемещения и хранения.

Требования к системе:

- Должна быть реализована клиент-серверная архитектура, причем клиент должен предоставлять GUI.
- Должна быть обеспечена неизменяемость сохраненных данных.
- Должна быть реализована система пользователей и прав доступа.
- Механизм разработки парсеров может быть реализован как с использованием разработанного DSL, так и с использованием уже существующего языка программирования.
- Должен быть реализован API, позволяющий производить накопление сырых данных непосредственно в хранилище в процессе эксперимента (путем вызова методов API из некоторого связующего ПО для экспериментальной установки).

3. UML модель системы

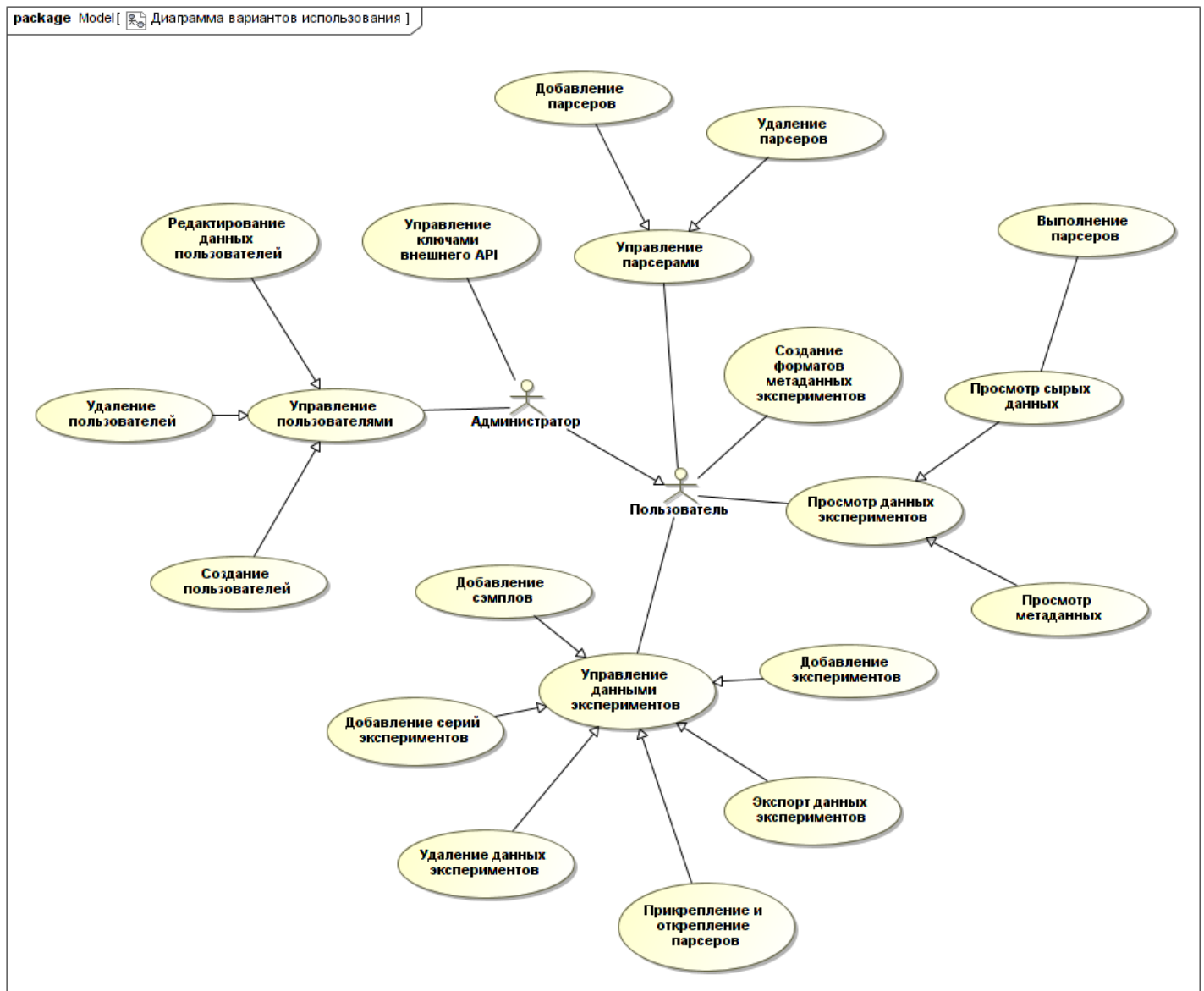


Рис. 1 Диаграмма вариантов использования

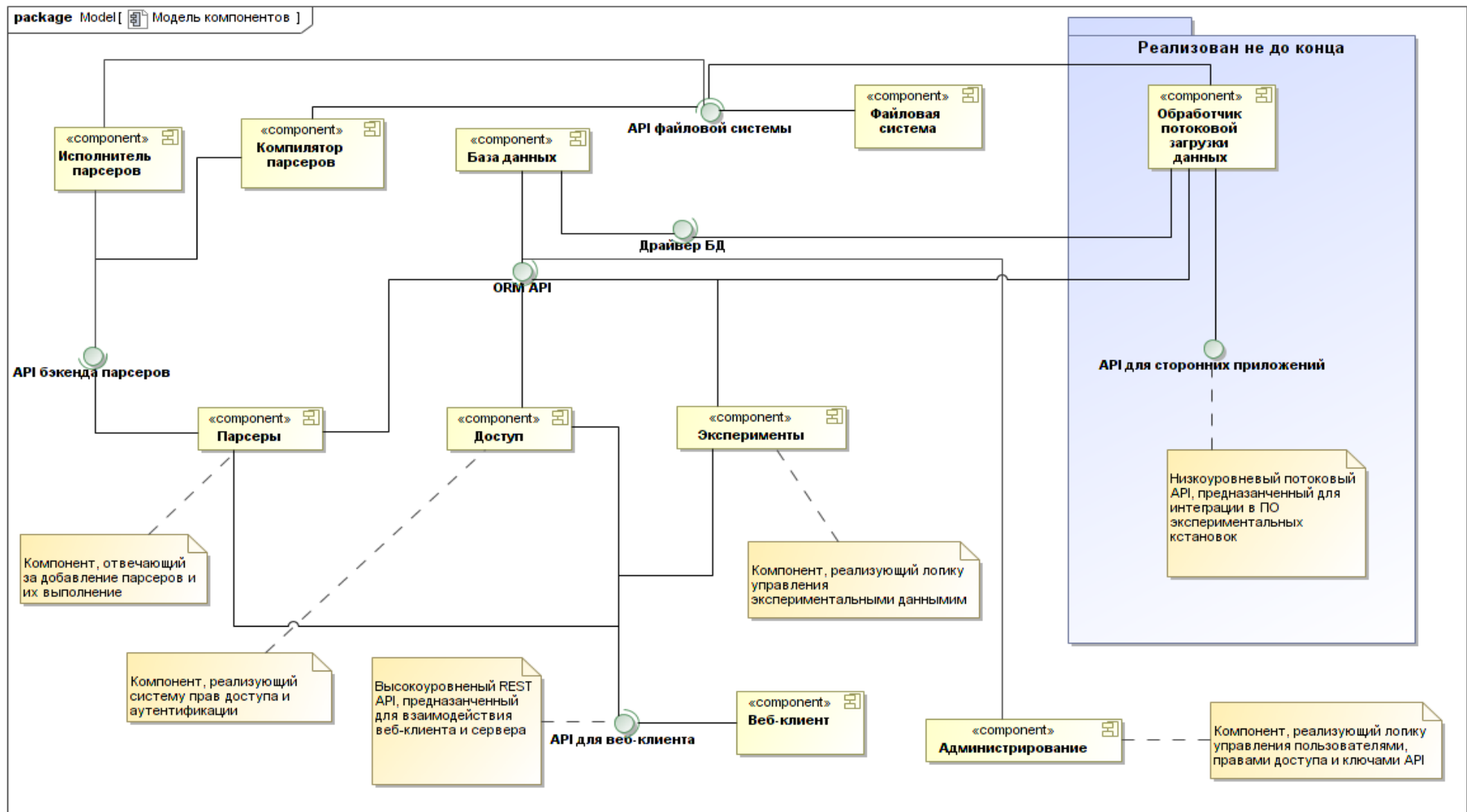


Рис. 2 Диаграмма компонентов

Пояснение к рис.2 (диаграмме компонентов):

Описание компонентов:

- **Эксперименты.** Отвечает за управление форматами метаданных; записями экспериментов и серий в БД; отображение и управление сериями, экспериментами и семплами; экспорт данных.
- **Доступ.** Отвечает за аутентификацию пользователей, проверку прав доступа пользователей.
- **Парсеры.** Отвечает за создание и отображение парсеров.
- **Исполнитель парсеров.** Отвечает за валидацию парсеров и их выполнение.
- **Компилятор парсеров.** Отвечает за валидацию парсеров и их компиляцию.
- **База данных.** Отвечает за хранение данных, предоставляет обертку в виде ORM.
- **Файловая система.** Отвечает за хранение «сырых» данных, выходных данных парсеров и их выдачу.
- **Администрирование.** Отвечает за создание и редактирование пользователей; назначение прав доступа; управление ключами потокового API.
- **Обработчик потоковой загрузки данных.** Отвечает за прием потоковых данных от экспериментальных установок через TCP сокет, обеспечение их целостности и валидности.

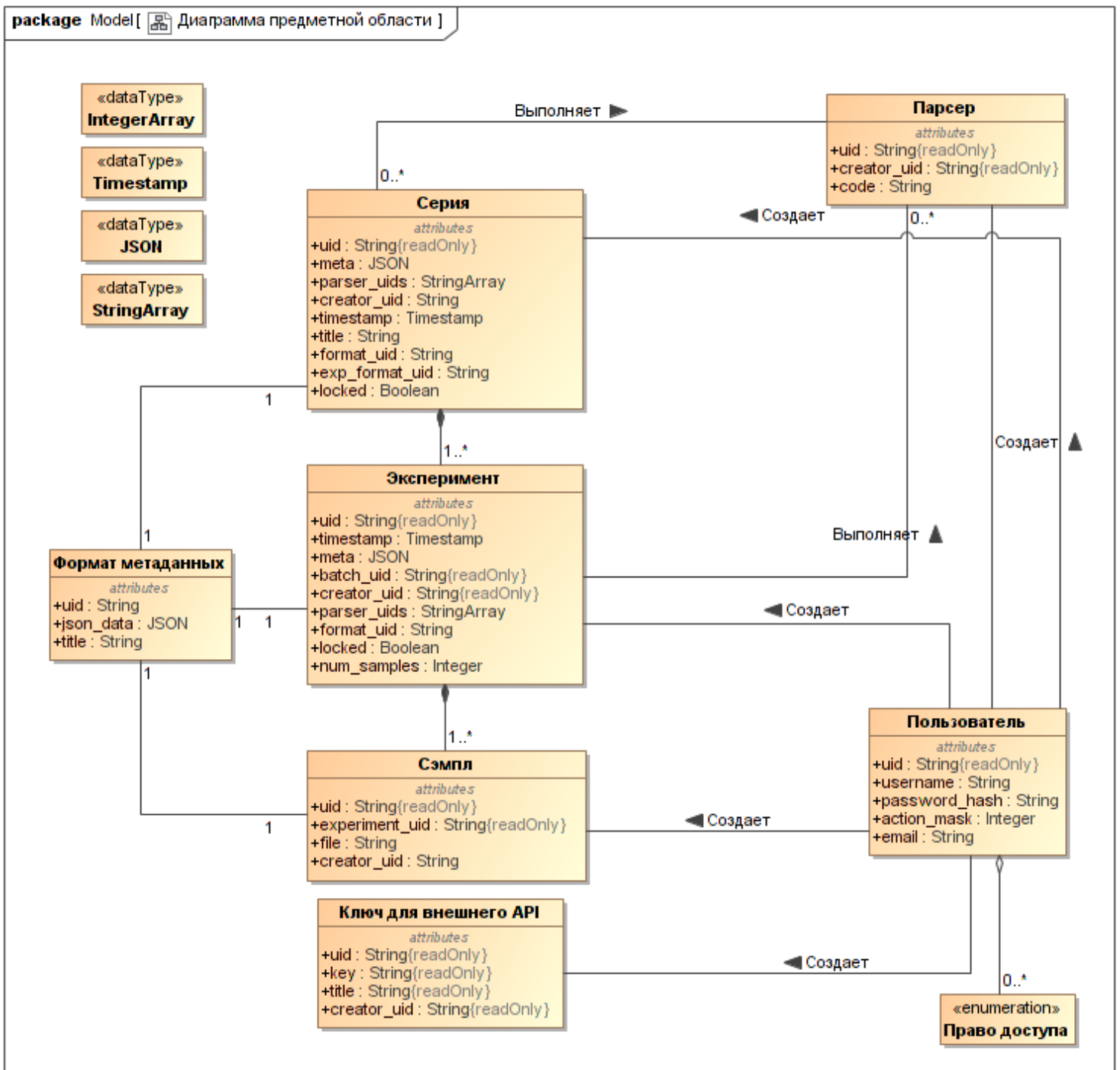


Рис. 3 Модель предметной области

Пояснение к рис. 3 (модели предметной области):

Сущности:

- **Формат метаданных.** Формат протокола эксперимента в виде JSON.
- **Серия.** Набор экспериментов, объединенных общими свойствами. Имеет свои метаданные, а также определяет формат метаданных экспериментов, в нее входящих.

- **Семпл.** Единица «сырых» данных эксперимента. В системе представляется произвольным файлом. В отрыве от эксперимента и его метаданных смысла не имеет.
- **Ключ для внешнего API.** Ключ-токен для авторизации клиентских приложений во внешнем потоковом API.
- **Парсер.** Программа на языке программирования Python, реализующая определенный интерфейс.
- **Право доступа.** Определенное действие, к которому пользователь может иметь или не иметь доступ. В системе представлено значением сдвига битовой маски всех прав доступа пользователя.
- **Пользователь.** Профиль пользователя. Используется для аутентификации и идентификации в системе.

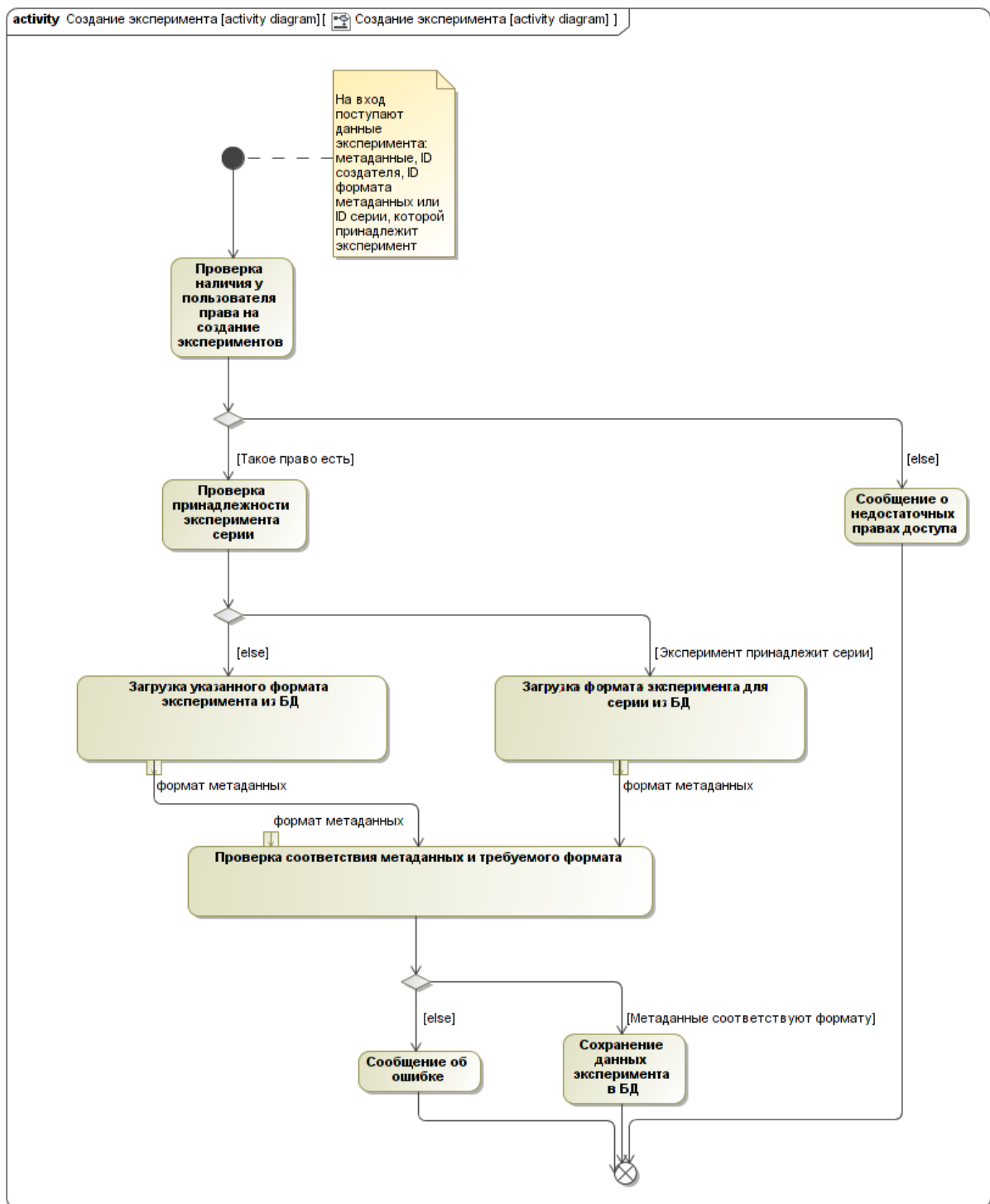


Рис. 4 Диаграмма деятельности процесса создания эксперимента

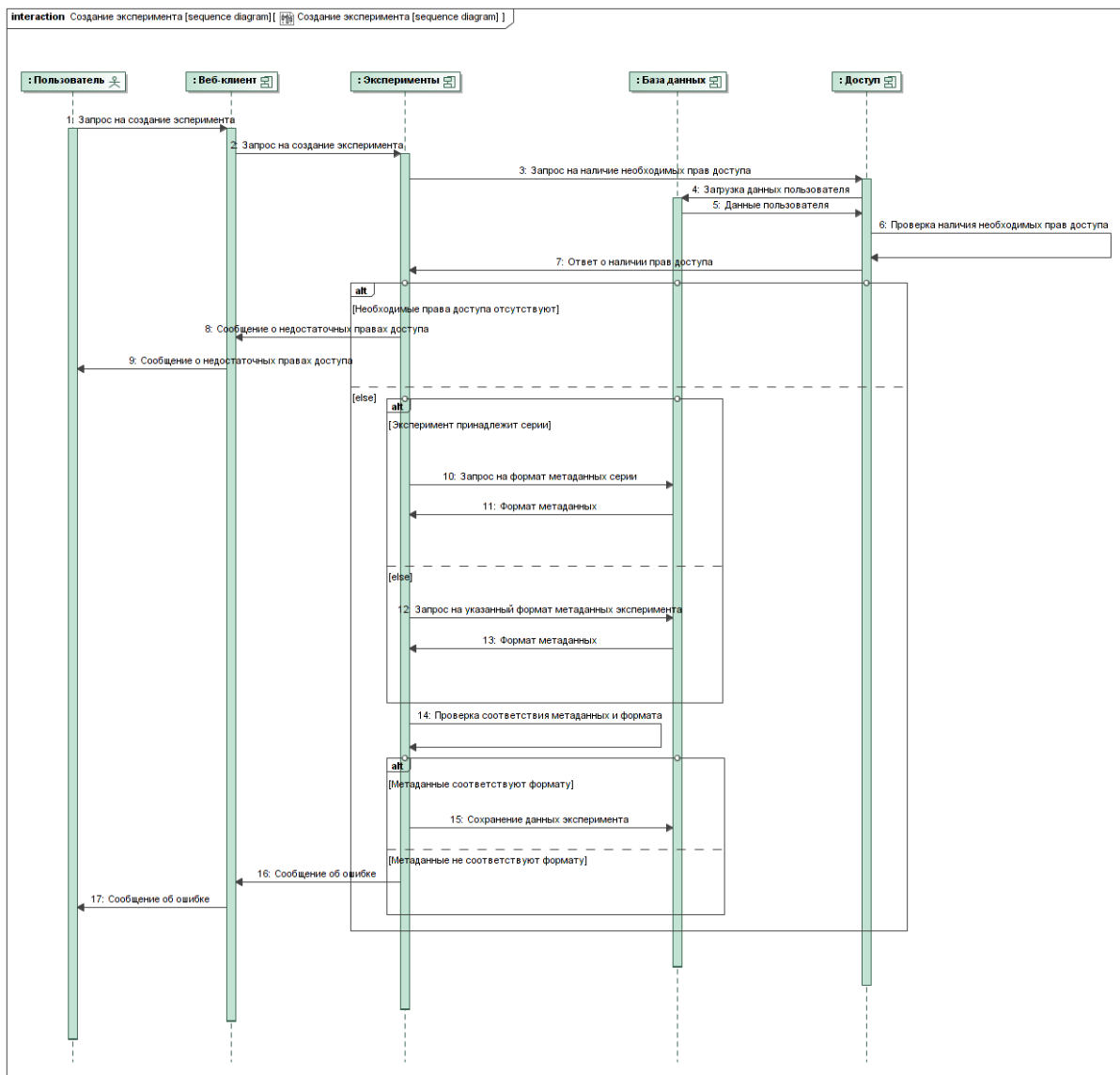


Рис. 5 Диаграмма последовательности процесса создания эксперимента

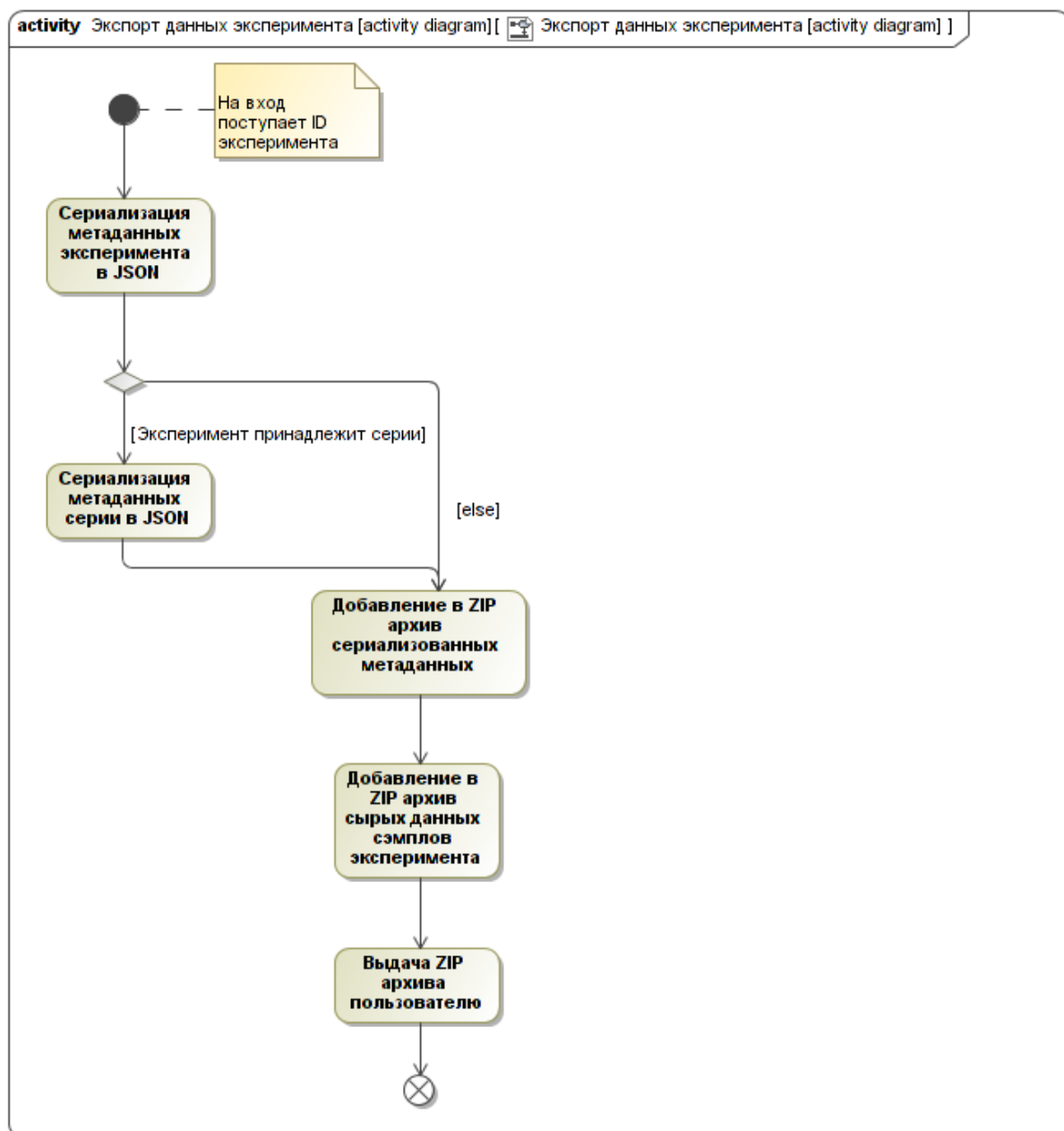


Рис. 6 Диаграмма деятельности процесса экспорта данных эксперимента

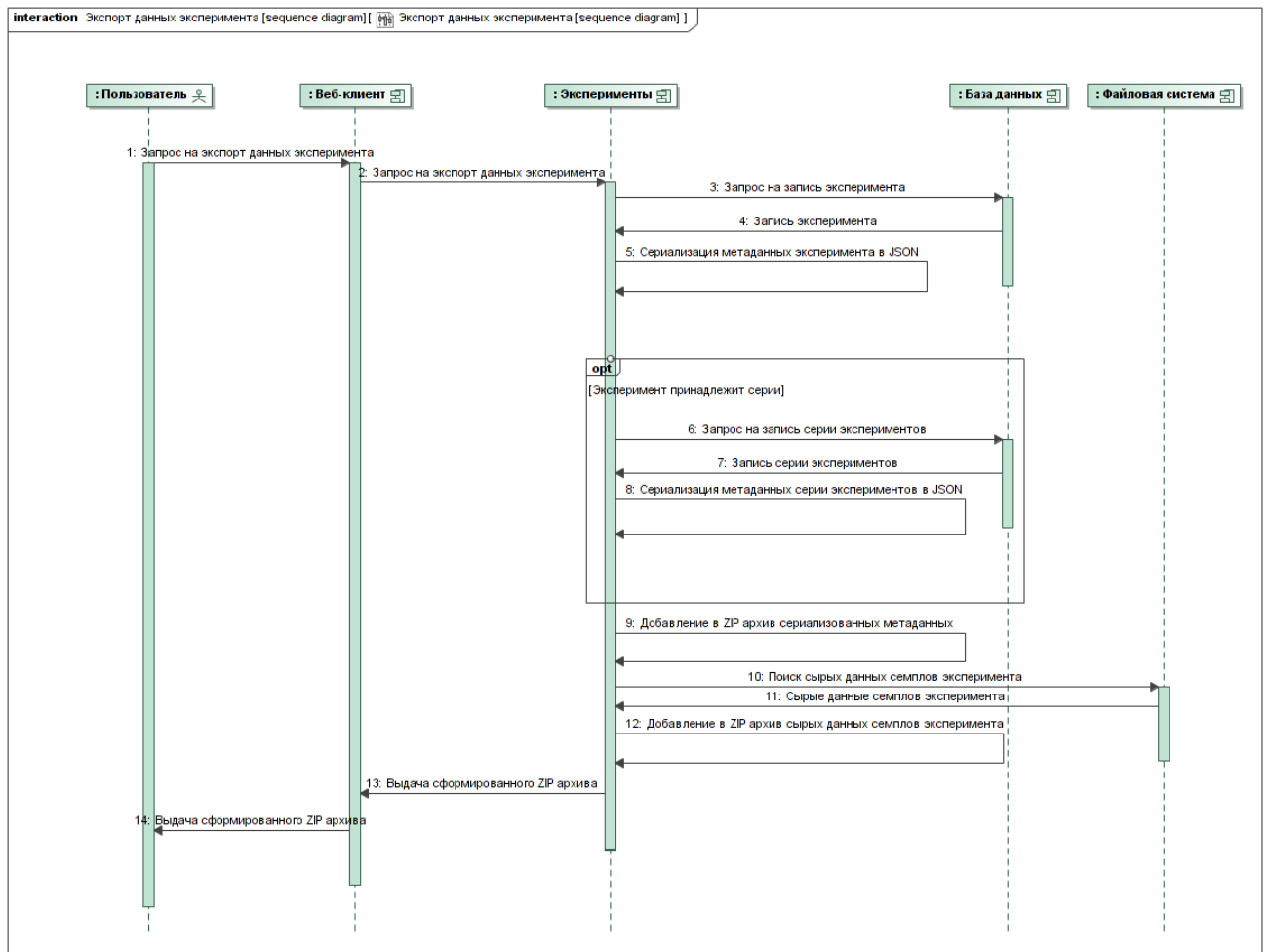


Рис. 7 Диаграмма последовательности процесса экспорта данных эксперимента

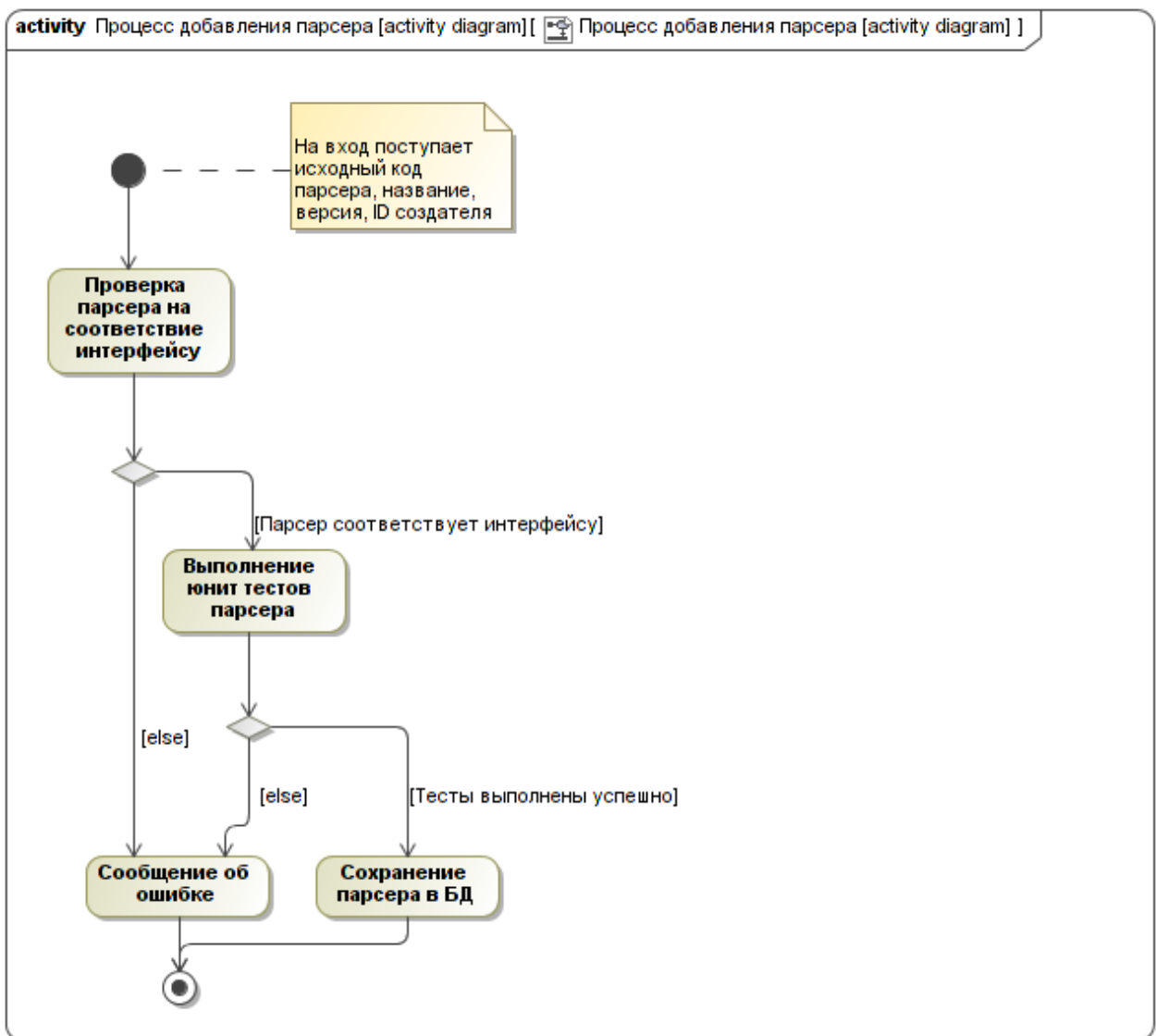


Рис. 8 Диаграмма деятельности процесса создания парсера

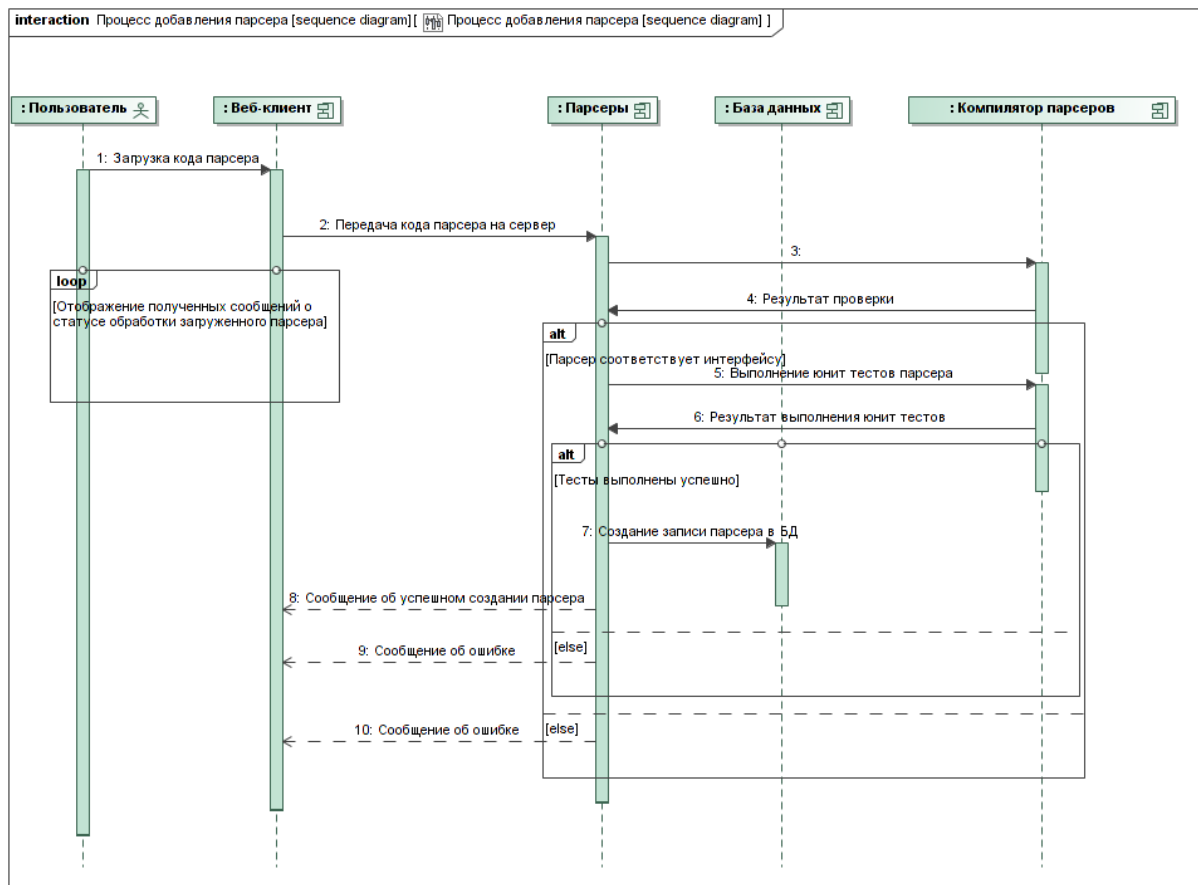


Рис. 9 Диаграмма последовательности процесса создания парсера

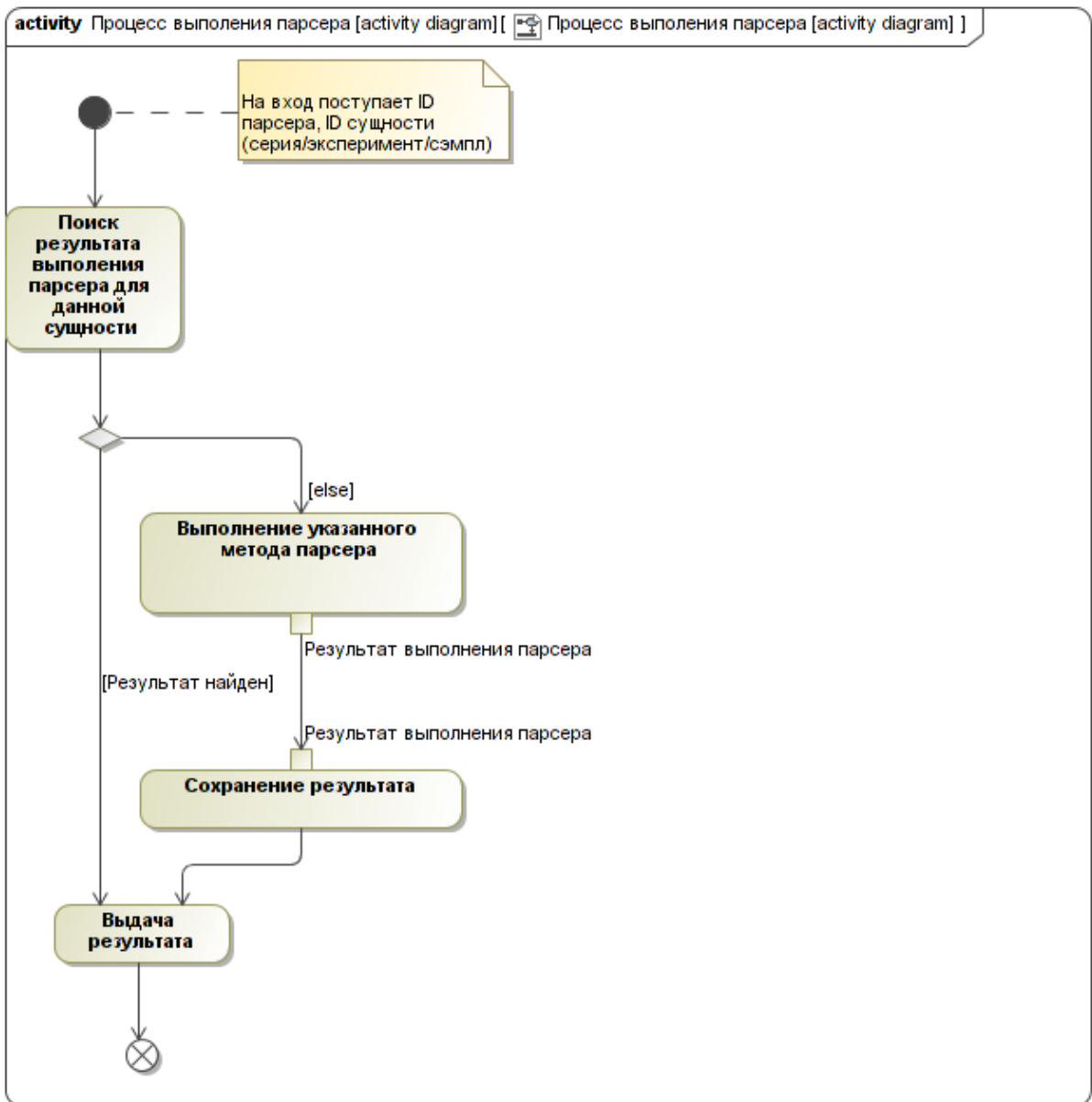


Рис. 10 Диаграмма деятельности процесса выполнения парсера

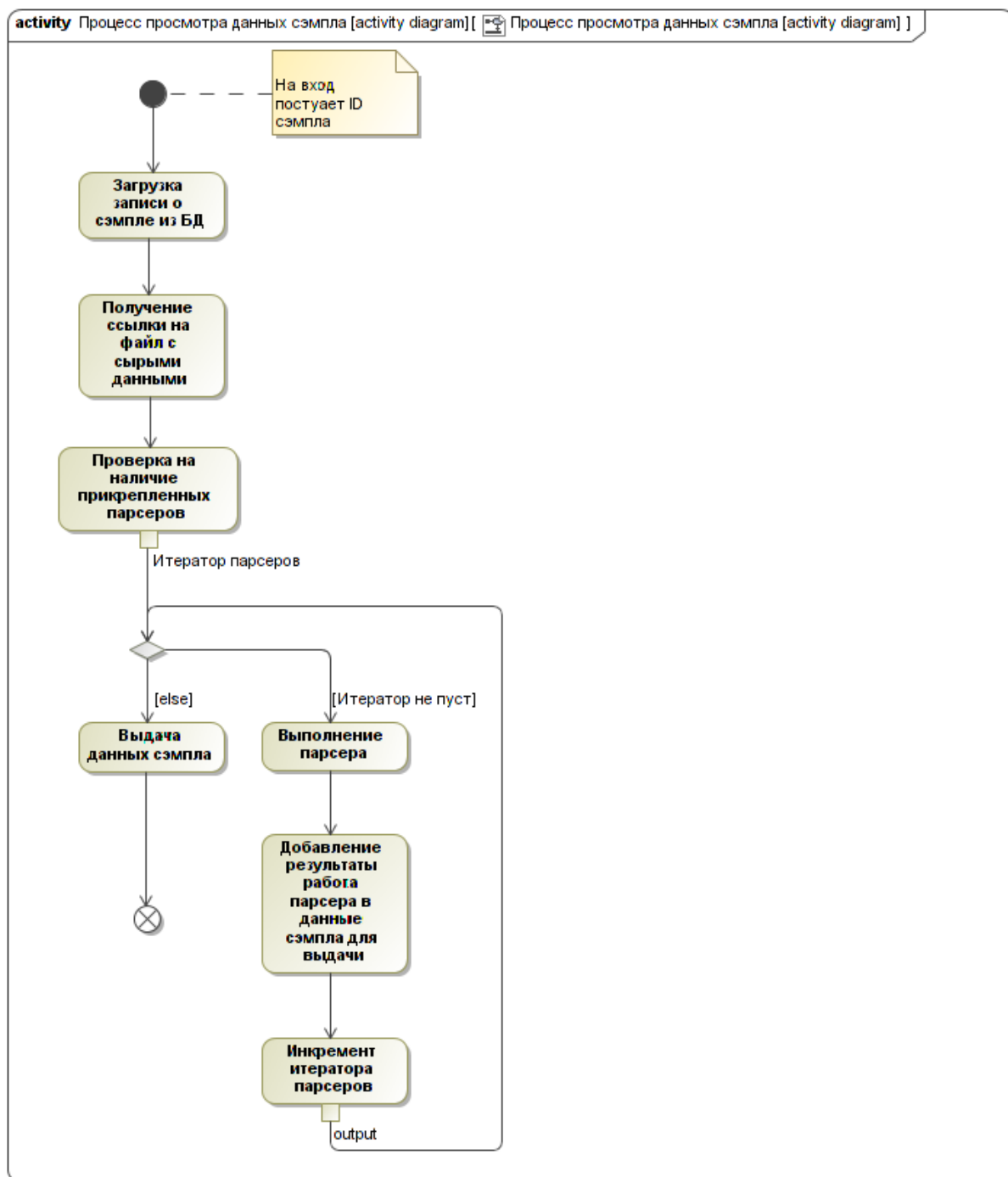


Рис. 11 Диаграмма деятельности процесса просмотра данных семпла

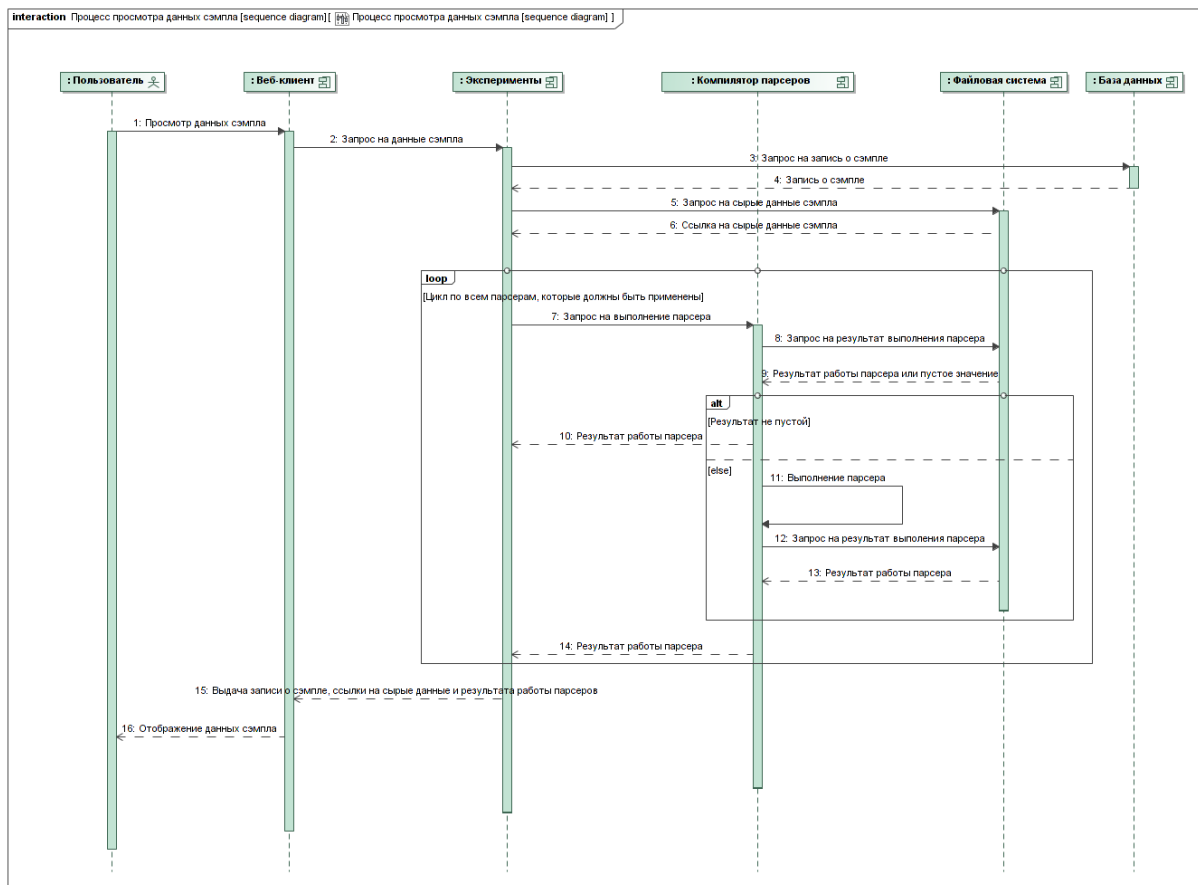


Рис. 12 Диаграмма последовательности процесса просмотра данных сэмпла

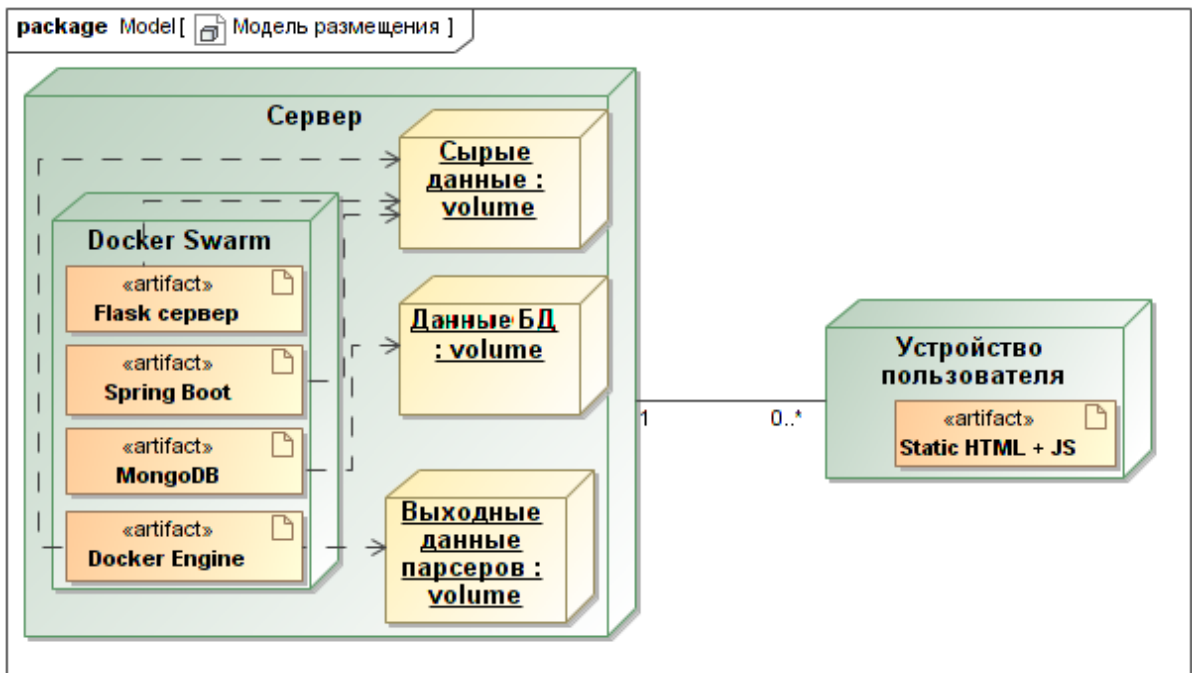


Рис. 13 Модель размещения

4. Используемый стек технологий

Для реализации системы, удовлетворяющей представленным требованиям, был использован следующий стек технологий:

- WEB интерфейс:
 - Фронтенд: Jinja2 + Twitter Bootstrap + JavaScript
 - Бэкенд: Python + Flask
- База данных: MongoDB
- Конвейер выполнения парсеров: Docker Engine + Python
- Компонент потоковой обработки “сырых” данных: Java
- Инфраструктура для развертывания приложения: Docker Engine, Docker Swarm, Bash скрипты

Исходный код системы доступен в Git репозитории:

<https://github.com/dev0x13/euclid>

5. Результаты

Разработанная система в целом удовлетворяет поставленной задаче: не был реализован только компонент генерации отчетов; компонент потокового API был реализован на уровне работающего прототипа.

Вся разработка выполнялась в рамках модели «Инфраструктура как код», то есть автоматизировано и отражено в коде скриптов все процессы, которые автоматизировать. Кроме того, средствами Docker Swarm все зависимости системы сведены к Docker Engine: достаточно поставить его на чистую операционную систему (Linux) и запустить скрипт, который развернет систему. На момент написания система развернута по адресу: <http://3.120.98.247:5000>

Реализована система аутентификации (средствами Flask-Login), а также система прав доступа и механизм управления ими. Пример конфигурации профиля пользователя при его создании показан на рис. 14.

Рис. 14 Окно создания и редактирования профиля пользователя

Реализован механизм управления доступом внешних приложений к потоковому API (рис. 15.). Сам потоковый API находится в стадии работающего прототипа (разработан на языке программирования Java).

Название	Ключ	Создатель	
Большая азотруба	f2979d32f6cc41c5b325ad6552124b2f	admin	
Малая азотруба	35e14936af5f437e922f19708c630f35	admin	

Рис. 15 Окно управления ключами потокового API

Реализована система управления форматами метаданных экспериментов и серий экспериментов, и их валидации. В качестве метаформата был выбран JSON. Валидация производится по структуре метаданных (рекурсивно), именам ключей, а также по типу значений. Таким образом, гарантируется наличие корректных метаданных для каждого набора данных. Пример созданных форматов метаданных показан на рис. 16. Пример создания формата в интерактивном редакторе и проверкой грамматики показан на рис. 17.

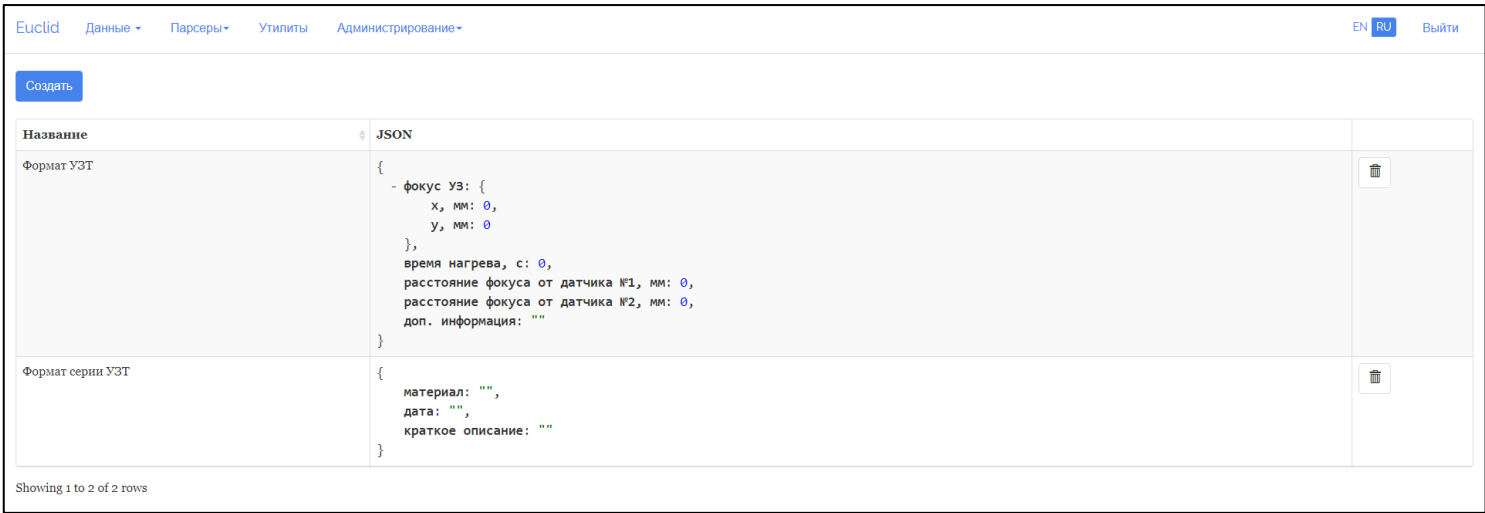


Рис. 16 Окно управления форматами метаданных

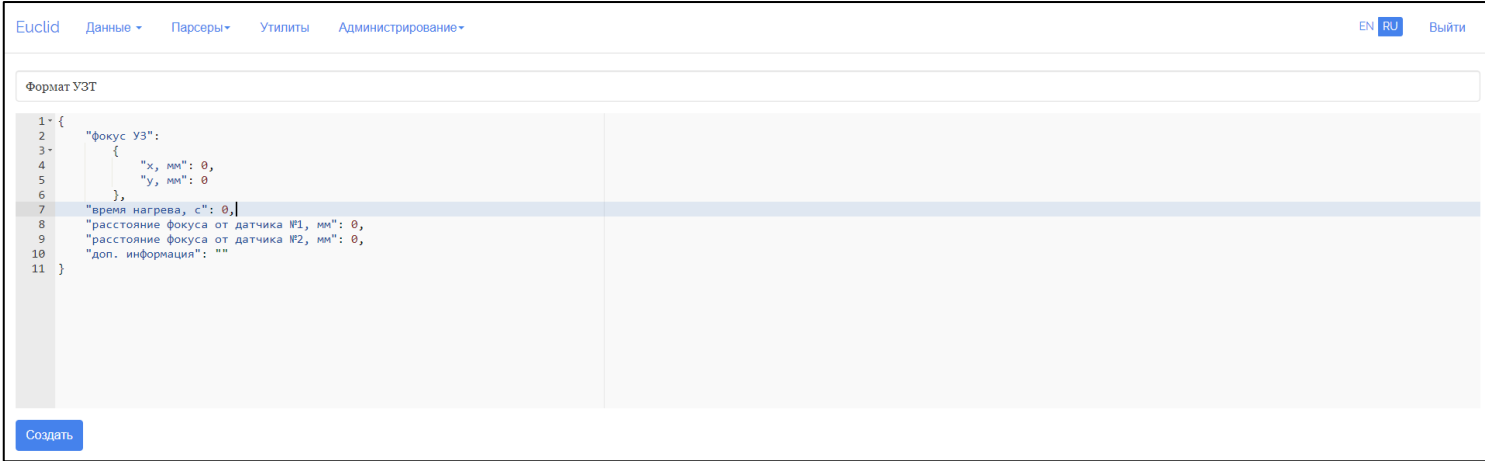


Рис. 17 Окно создания формата метаданных

Реализована система управления данными экспериментов. Для серий и экспериментов реализован механизм блокировки, который обеспечивает неизменяемость данных после их добавления. На рис. 18 показан внешний вид серии экспериментов с заполненными метаданными и добавленным экспериментом, а на рис. 19 – внешний вид эксперимента с добавленными семплами. Кроме того, реализована возможность экспорта данных на любом из уровней (серия, эксперимент, семпл). Запрошенные на экспорт данные упаковываются в ZIP архив совместно с метаданными и выдаются запросившему пользователю.

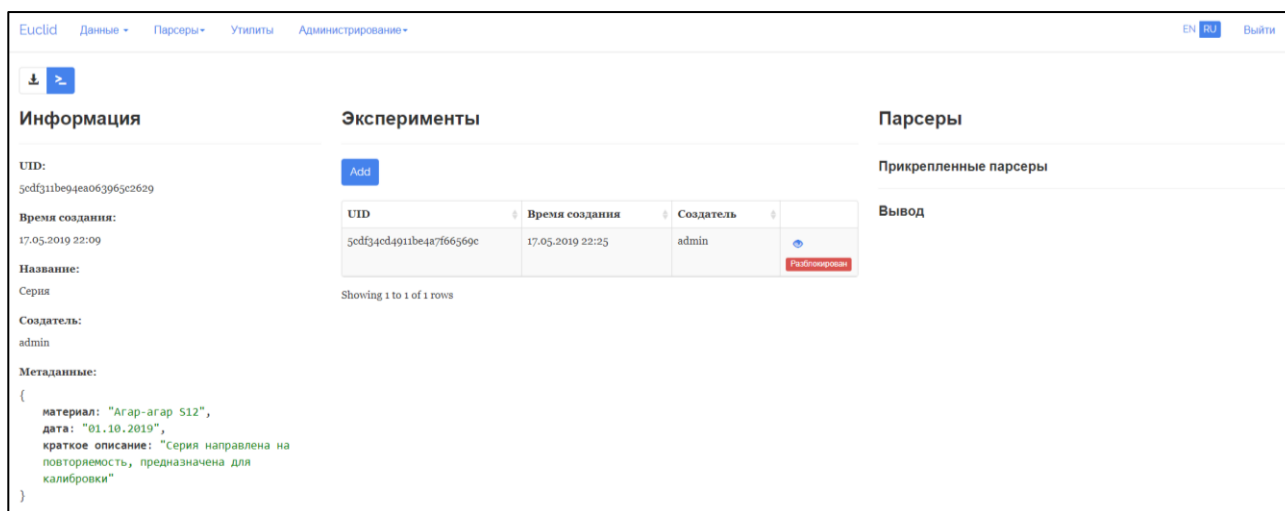


Рис. 18 Окно просмотра данных серии экспериментов

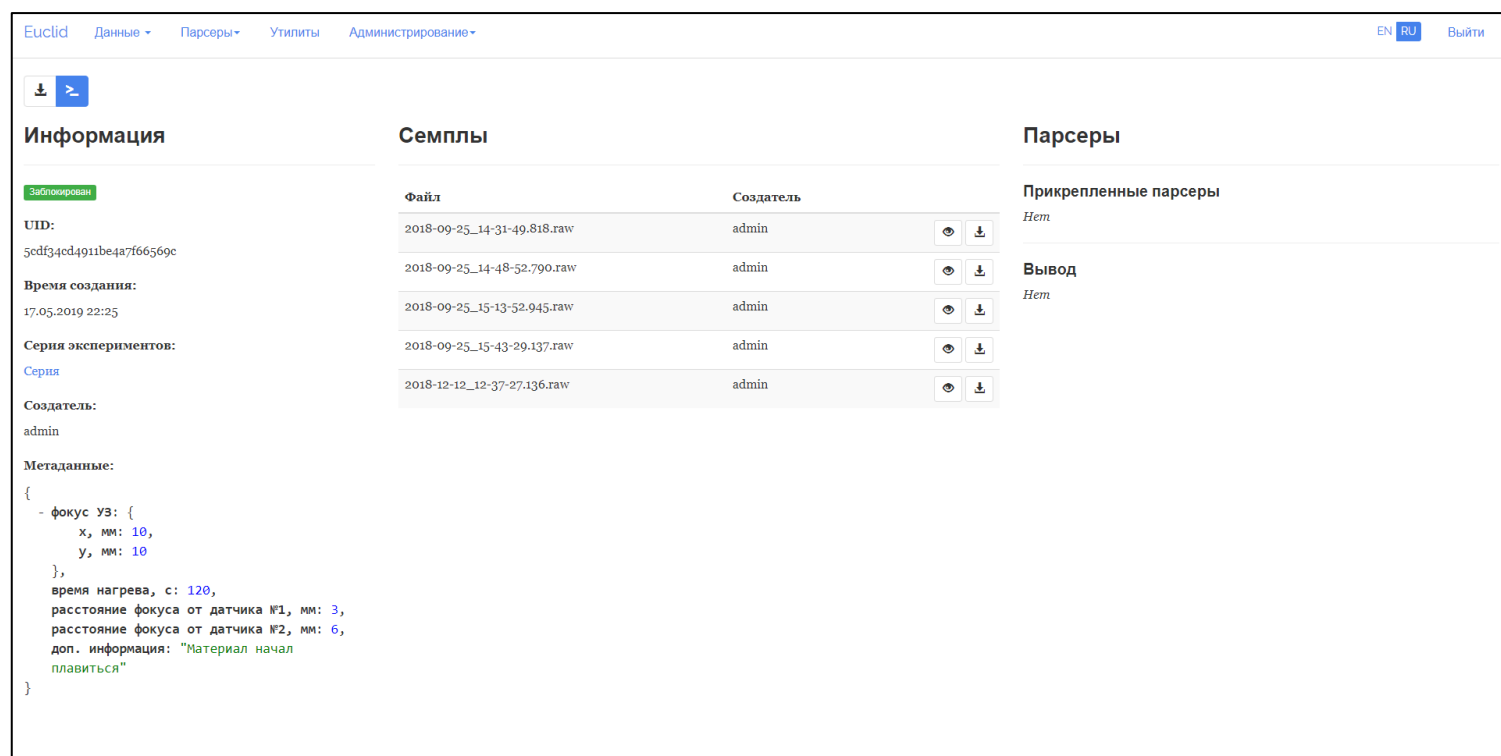


Рис. 19 Окно просмотра данных эксперимента

В соответствии с техническим заданием в качестве возможности первичного анализа данных экспериментов была разработана подсистема парсеров. Разработанная подсистема позволяет отлаживать и выполнять пользовательский код специального вида на языке программирования Python изолированно (и безопасно) от серверного окружения. Для разработки была создана интегрированная в браузер среда (рис. 20). Разработчику парсера предоставлены функции для вывода изображений и текста в форму веб-интерфейса сущности, к которой прикреплен парсер. Кроме того, для облегчения обработки и анализа данных в подсистеме установлены следующие пакеты Python: NumPy, SciPy, Pillow, Matplotlib.

Пример работы парсера «Спектрограмма WAV файла» на тестовых экспериментальных данных показана на рис. 21. Неизменяемость данных как главный принцип системы позволяет выполнять парсеры всего один раз для каждой единицы данных: результат кэшируется и его актуальность в любой момент времени гарантирована.

EuclidДанныеПарсерыУтилитыАдминистрирование

ENRUВыйти

Спектрограмма WAV файла

```
1 import numpy as np
2 import math
3 import io
4 import wave
5 import matplotlib.pyplot as plt
6
7 def np_spectrogram(s, nfft, length, shift):
8     S = [np.fft.fft(s[n:n+length], n=nfft) for n in range(0, s.size - length + shift, shift)]
9     S = np.abs(np.asarray(S).T[:nfft // 2, :]).astype(np.complex64)
10    return S
11
12 nfft = 64
13 length = 480 # 30 ms for 16k
14 shift = length
15
16 class ParserImpl(Parser):
17     def process_sample(self, sample):
18         with wave.open("input/As" % sample["file"]) as f:
19             num_samples = f.getnframes()
20             sample_rate = f.getframerate()
21             samples = np.frombuffer(f.readframes(num_samples), dtype=np.int16)
22             s = np_spectrogram(samples, nfft, length, shift)
23             plt.figure()
24             plt.pcolormesh(np.squeeze(s))
25             buf = io.BytesIO()
26             plt.savefig(buf, format='png')
27             buf.seek(0)
28             self.print_image(buf.read())
29             buf.close()
30
31     def process_experiment(self, experiment):
32         pass
33
34     def process_batch(self, batch):
35         pass
36
```

ВалидацияКомпиляция

Вывод компилятора:

```
.....
STATIC: OK!
DYNAMIC: OK!
.....
STATIC: Error: Command 'python validator.py' in image 'euclid-parser-env' returned
non-zero exit status 1: b'missing 'process_experiment' method'
DYNAMIC: OK!
STATIC: OK!
DYNAMIC: OK!
```

Debug input:

Семплы50d3845491be4a766656a6

Отладочный вывод:

Рис. 2019 Интегрированная среда разработки парсеров

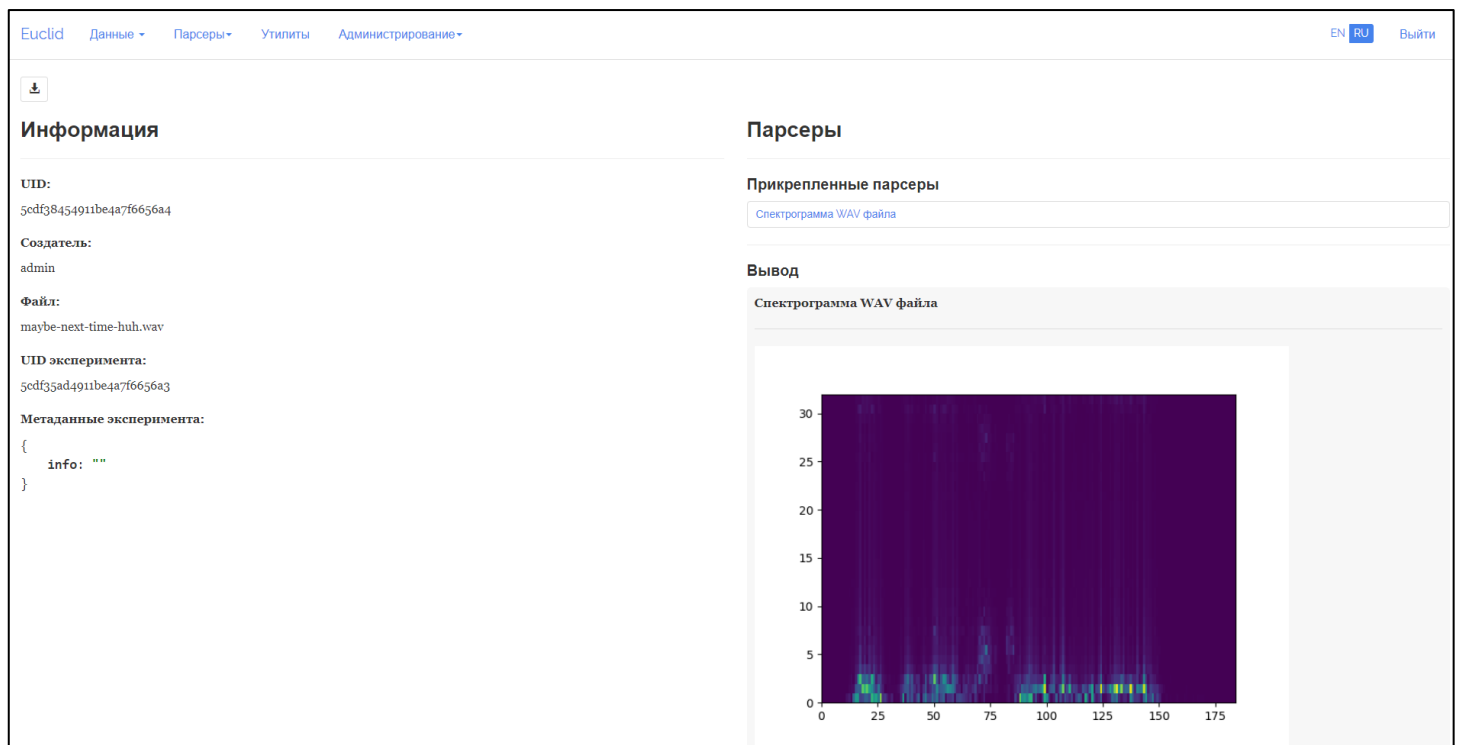


Рис. 201 Окно просмотра данных семпла с прикрепленным парсером

Кроме того, реализована поддержка английского и русского языков интерфейса с возможностью переключения пользователем. Языковая локализация реализована средствами Babel и имеет возможность легкого расширения для любого набора языков.

6. Дальнейшее развитие проекта

Дальнейшие шаги по развитию проекта:

- Доработка потокового API
- Модернизация подсистемы парсеров: добавление возможности установки целых пользовательских модулей (по аналогии с AWS Lambda)
- Написание подробной документации
- Разработка CLI (Command Line Interface) для возможности автоматизации работы с данными экспериментов