

**ДООБУЧЕНИЕ РУССКОЯЗЫЧНОЙ BERT-LIKE МОДЕЛИ  
ПОД ЗАДАЧИ МУЛЬТИЛЕЙБЛА И ИЕРАРХИЧНОЙ КЛАССИФИКАЦИИ**

# АКТУАЛЬНОСТЬ

Для кого?



**Можно использовать модель чтобы:**

1. Видеть суть новости с первого взгляда.
2. Находить ту самую новость в информационном шуме.
3. Быстро принимать решения

## ПРОИЗВОДСТВЕННЫЕ ФАКТЫ

- **Обработка 1000 новостей вручную** (по 1 мин/новость = 5,5 часов) при ставке 500 руб/час стоила бы **27 500 ₽**
- **Наша модель** справляется с этим объемом за **33 минуты** с **минимальными затратами** на вычислительные ресурсы!

When I realize  
ChatGPT can do  
my job for me

When I realize  
ChatGPT can do  
my job for me

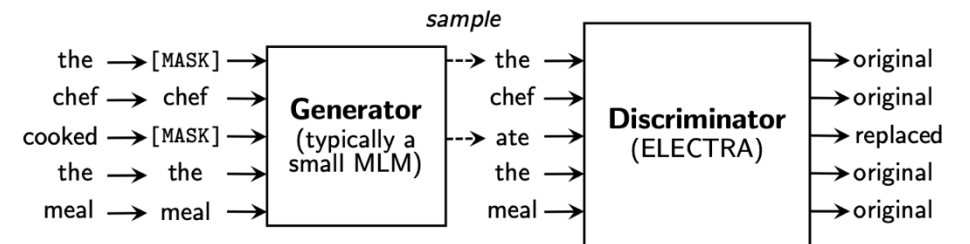
## ЦЕЛЬ

- Разработка и внедрение системы автоматической классификации текстовых новостных сообщений для Telegram-бота.

# ЗАДАЧИ

- 🤖 **Подготовка данных:** Использование LLM и промпт-инжиниринга для иерархической и многоклассовой разметки.
- 🧠 **Дообучение RuELECTRA:** Тонкая настройка модели под задачи проекта с кастомной лосс-функцией.
- 🖱️ **Интеграция модели:** Внедрение дообученной RuELECTRA для обработки потока новостей.
- 📰 **Автоматический сбор новостей:** Реализация парсера для получения свежих статей с телеграм канала РИА Новости.
- 🔍 **Интерпретируемость модели:** Визуализация значимости слов (SHAP values) для понимания решений модели.
- 🚀 **Доставка контента:** Своевременная отправка пользователю обработанных новостей с тегами и анализом.

# АРХИТЕКТУРА ELECTRA



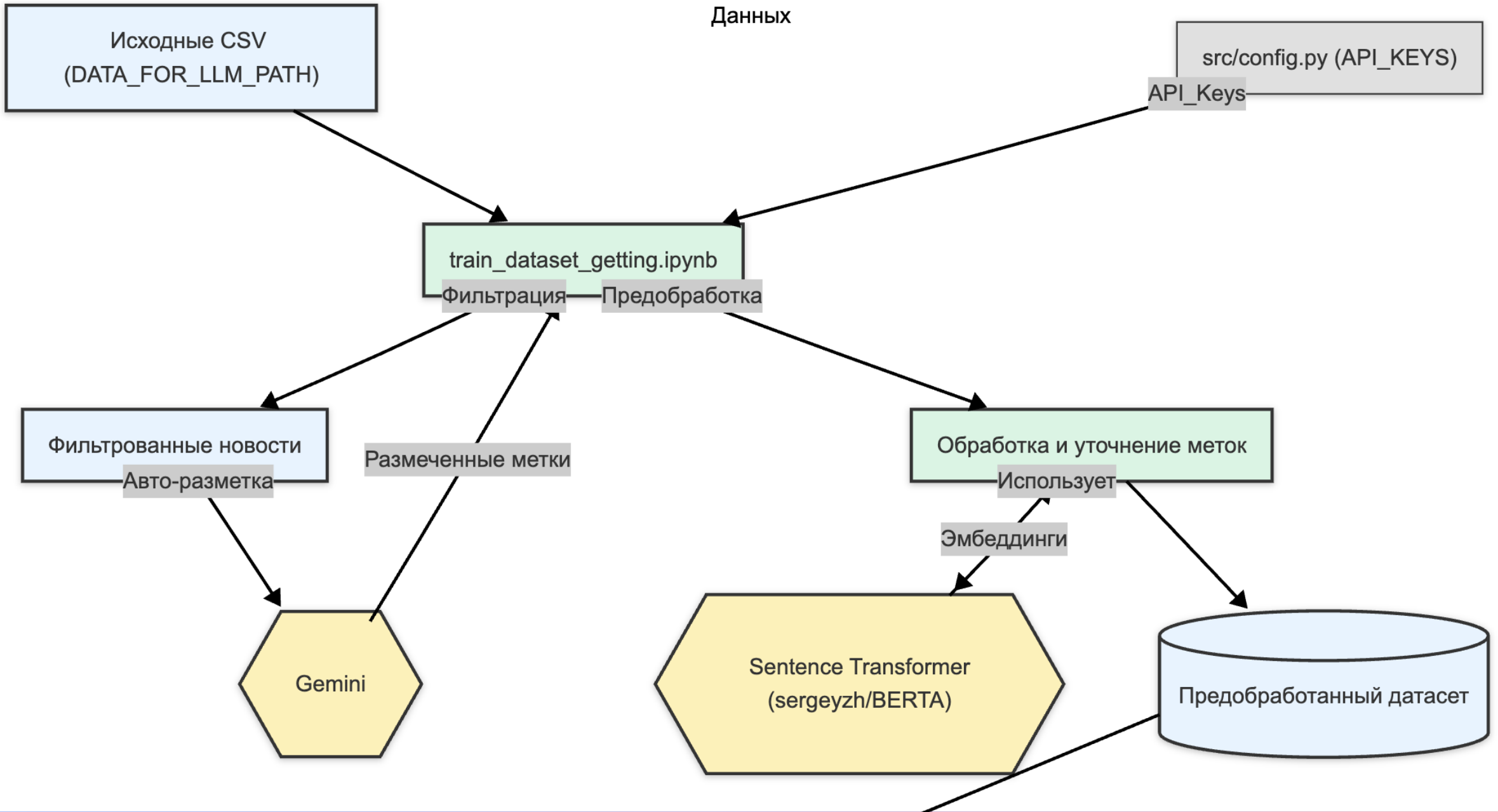
# ДАННЫЕ

- Телеграм канал РИА Новости 2023 – 2024 год
- Разметка лейблов с помощью **Gemini-2.5-flash-preview-05-20** с низкой температурой для уменьшения высокой вариативности в лейблах
- Благодаря этому, была получена разметка для:
  - **55198** текстов
  - **169** иерархических классов
  - **14** классов для мультитейблинга



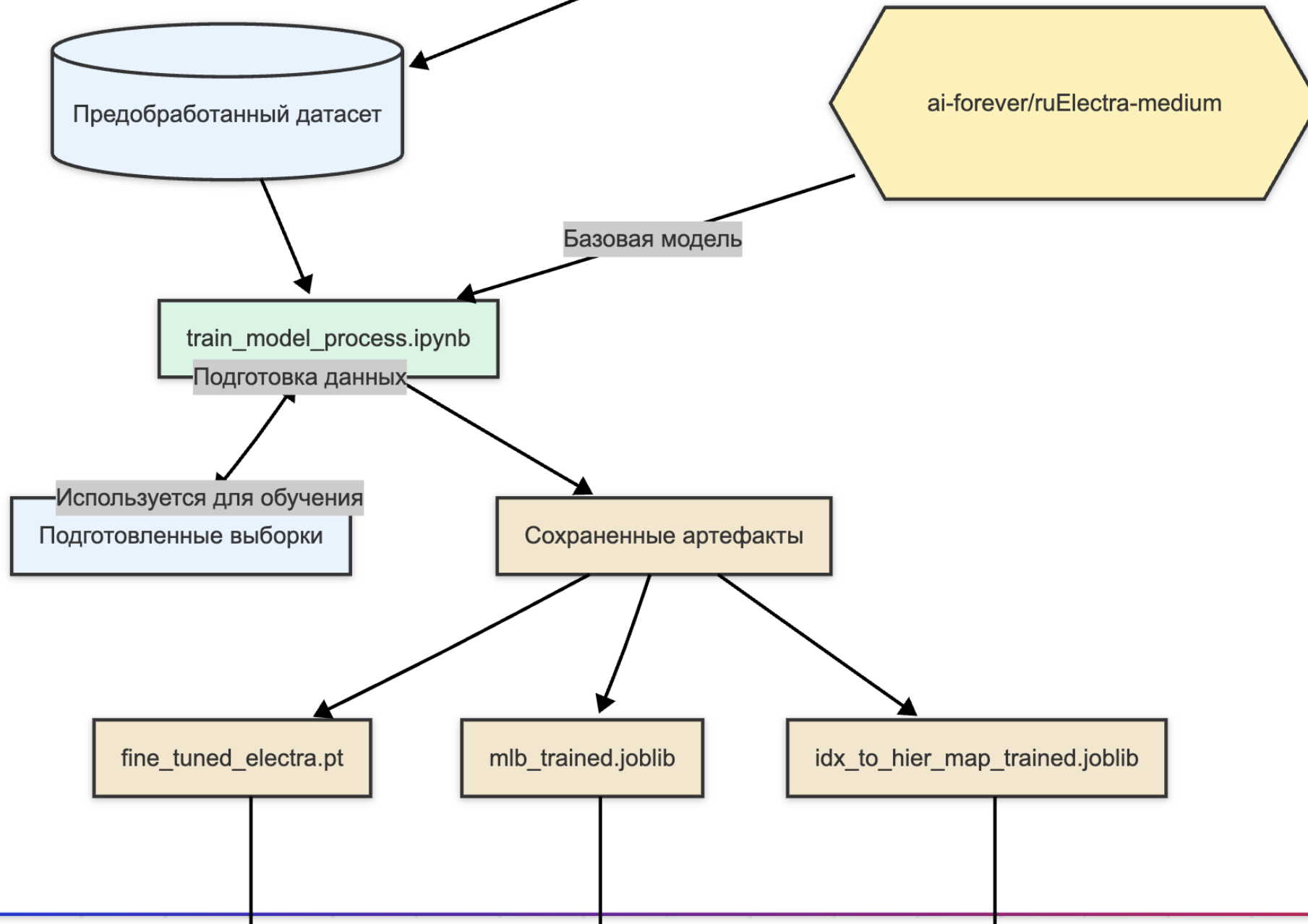
# АРХИТЕКТУРА ПРОЕКТА

Этап 1 и 2: Подготовка  
Данных

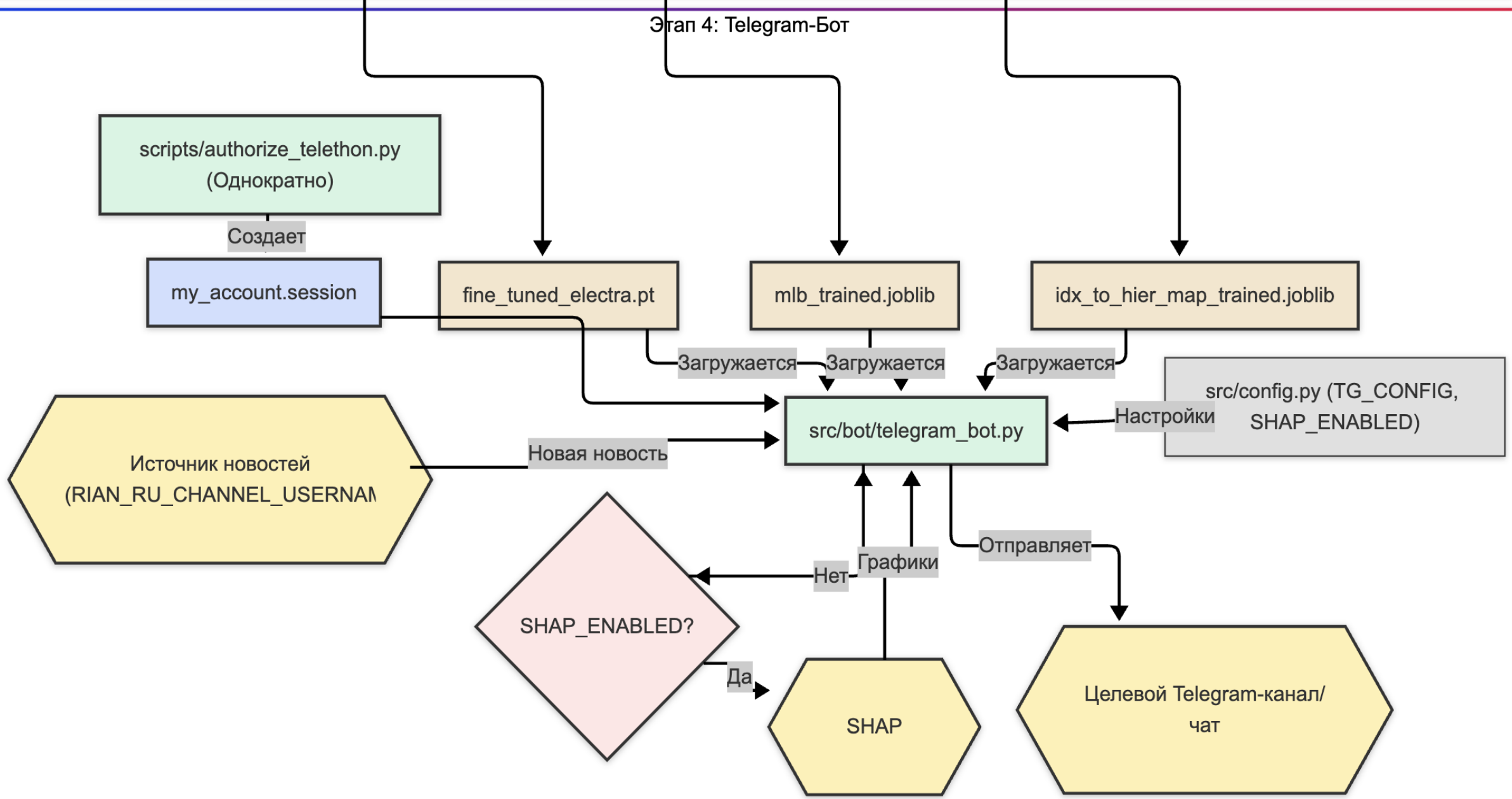




### Этап 3: Обучение Модели



#### Этап 4: Telegram-Бот



# ВОЗМОЖНОЕ ПРИМЕНЕНИЕ МОДЕЛИ

-  **Анализ инфополя:**
  - Выявление скрытых трендов.
  - Отслеживание пересечения тем.
  - Глубокое понимание повестки дня.
-  **Персонализация контента:**
  - Формирование лент и рекомендаций.
  - Учет интересов пользователя на разных уровнях детализации.
-  **Интеллектуальная каталогизация и архивация:**
  - Структурирование баз знаний.
  - Организация научных статей и юридических документов.
  - Навигация по сложным иерархиям данных.

## СПОНСОРЫ



СПАСИБО ЗА ВНИМАНИЕ

