

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ  
РОССИЙСКОЙ ФЕДЕРАЦИИ

федеральное государственное автономное  
образовательное учреждение высшего образования  
«Самарский национальный исследовательский университет  
имени академика С.П. Королева»  
(Самарский университет)

Институт информатики и кибернетики  
Кафедра суперкомпьютеров и общей информатики

**Отчет по лабораторной работе №1**

Дисциплина: «Инженерия данных»

Выполнил: Вечканова П.А.

Группа: 6233-010402D

Самара 2025

## ЗАДАНИЕ

### Цели работы

- На практике освоить полный цикл ETL: извлечение из публичного API, трансформация и загрузка в ClickHouse.
- Научиться собирать и отлаживать пайплайны в **Prefect**.
- Научиться настраивать окружение для работы с данными.

### Описание пайплайна

Заданием на лабораторную работу является реализация следующего пайплайна:

1. В качестве источника данных предлагается использовать [Free Weather API](#)
2. (Extract) Получить прогноз **на завтра** по переменным: *температура, осадки, скорость и направление ветра* для городов **Самара и Москва**. Сырые ответы API сохранить в объектном хранилище
3. (Transform)
  - Извлечь почасовые значения и нормализовать для таблицы `weather_hourly`
  - Посчитать дневную статистику (`min`, `max`, `avg` температуру и количество осадков) и подготовить для сохранения в таблице `weather_daily`
4. (Load) Загрузить преобразованные данные в соответствующие таблицы ClickHouse
5. Автоматически отправить уведомления в Telegram с кратким прогнозом на завтра и предупреждать о сильном ветре/осадках
6. **(Опционально)** Реализовать обработку различных ошибок.

## ВЫПОЛНЕНИЕ РАБОТЫ

### 1. АРХИТЕКТУРА

Пайплайн реализован по схеме ETL с использованием:

- Prefect (версия Prefect 2.10): для оркестрации задач, планирования (ежедневно в 08:00 по МСК), логирования и повторных запусков;
- MinIO (RELEASE.2024-12-05T23-40-55Z): для хранения сырых JSON-ответов API (обеспечивает воспроизводимость и аудит);
- ClickHouse (версия 24.8): для эффективного хранения и анализа структурированных метеоданных;
- Telegram Bot API: для отправки краткого прогноза и предупреждений.

Выбор обусловлен открытостью, лёгкостью развёртывания в Docker и соответствием задаче обработки временных рядов.

### 2. ИСТОЧНИКИ ДАННЫХ

Использован публичный Open-Meteo API (<https://api.open-meteo.com/v1/forecast>).  
Параметры запроса:

- latitude, longitude – координаты Москвы (55.7522, 37.6156) и Самары (53.1959, 50.1001);
- hourly=temperature\_2m,precipitation,wind\_speed\_10m,wind\_direction\_10m;
- timezone=Europe/Moscow;
- forecast\_days=2 (чтобы получить данные на завтра).

### 3. ОПИСАНИЕ РАБОТЫ ETL-ПАЙПЛАЙНА

*Извлечение данных (Extract):* Для каждого города выполняется HTTP-запрос к Open-Meteo. Полученный JSON сохраняется в MinIO в формате {город}/{дата}.json – это позволяет сохранить исходные данные для отладки.

*Обработка данных (Transform):* Из сырых данных извлекаются почасовые значения только на следующий день. На их основе формируются записи для таблиц weather\_hourly и агрегированные дневные статистики (min, max, avg температура, сумма осадков) для weather\_daily.

*Загрузка данных (Load):* Данные загружаются в ClickHouse через clickhouse-driver. В обе таблицы добавляется поле \_loaded\_at с меткой времени загрузки. Одновременно формируется и отправляется Telegram-уведомление.

### 4. КАЧЕСТВО ДАННЫХ И ОБРАБОТКА ОШИБОК

Реализованы базовые проверки: валидация HTTP-статуса, проверка существования бакета MinIO, корректный парсинг дат.

### 5. ВЫПОЛНЕНИЕ ПАЙПЛАЙНА И РЕЗУЛЬТАТЫ

Первым этапом работы была настройка телеграм бота с помощью @BotFather и получение ID чата для дальнейшей настройки работы.

Далее, написав код и настроив docker-compose, запускаем контейнер с помощью команды `docker-compose up -d` (пример запуска и вывод представлены на рисунке 1)

```
✓ Volume etl_weather_minio_data Removed
PS C:\Users\vechk\etl_weather> docker-compose up -d
[+] Running 7/7
✓ Network etl_weather_default Created
✓ Volume etl_weather_minio_data Created
✓ Volume etl_weather_clickhouse_data Created
✓ Container etl_weather-minio-1 Healthy
✓ Container etl_weather-prefect-server-1 Started
✓ Container etl_weather-clickhouse-1 Started
✓ Container etl_weather-prefect-agent-1 Started
```

Рисунок 1 – Запуск контейнера

Делаем деплой расписания для того, чтобы отправка была автоматической. Пример команды представлен на рисунке 2.

```
PS C:\Users\vechk\etl_weather> docker-compose exec prefect-agent python -m src.deploy
15:22:01.180 | WARNING | prefect.deployments - The field 'schedule' in 'Deployment' has been deprecated. It will not be available after Sep 2024. Define schedules in the 'schedules' list instead.
```

Рисунок 2 – Деплой расписания

Заходим в Prefect и проверяем работу пайплайна

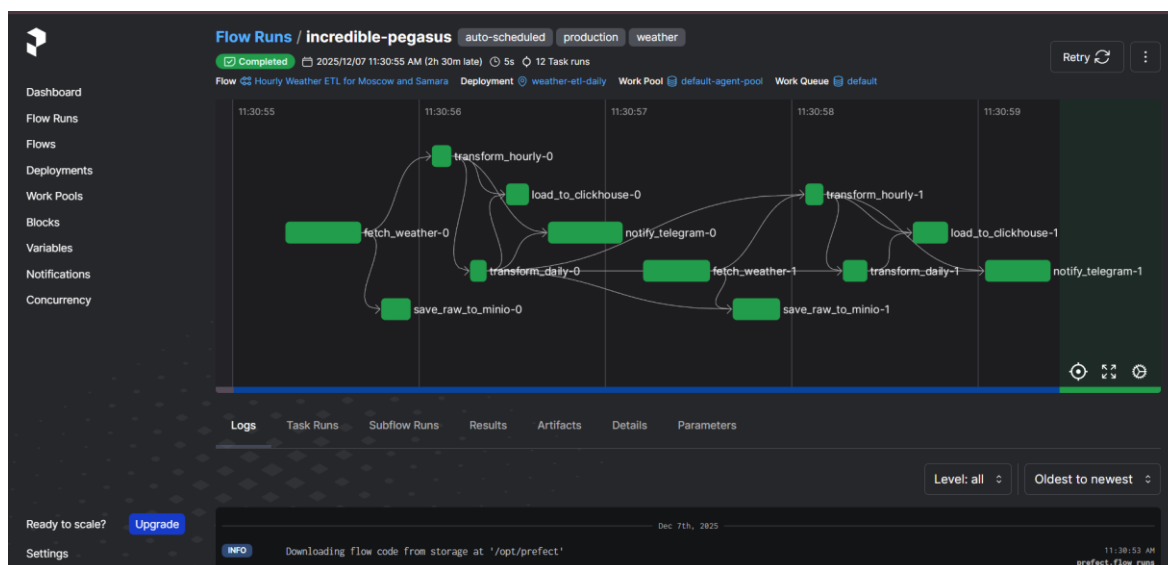


Рисунок 3 – Граф работы пайплайна

Зайдем в UI Clickhose (рисунок 4) и убедимся в том, что данные загружены с помощью команд

- `SELECT * FROM weather.weather_daily` (вывод представлен на рисунке 5)
- `SELECT * FROM weather.weather_hourly` (вывод представлен на рисунке 6)

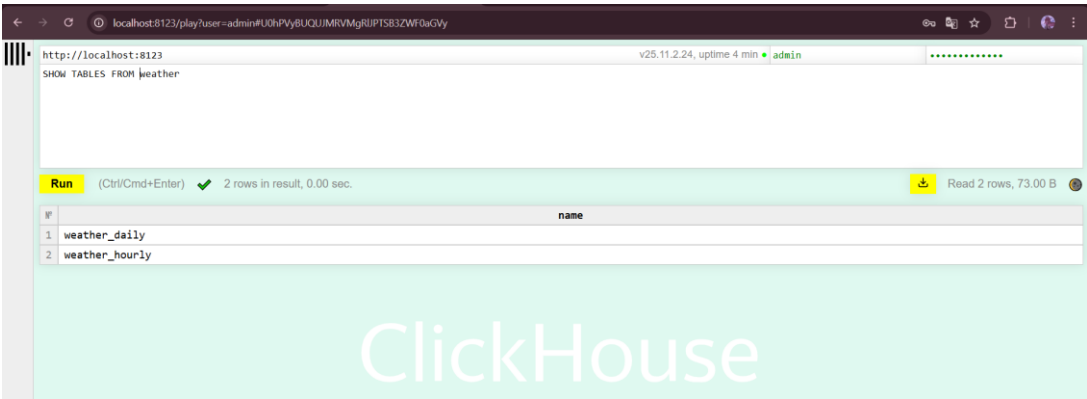


Рисунок 4 – UI Clickhose с созданными таблицами

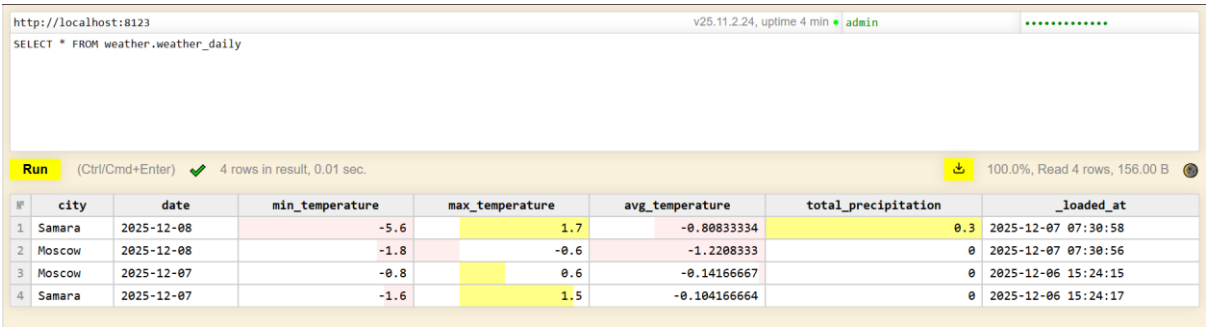


Рисунок 5 – SELECT запрос к таблице weather.weather\_daily

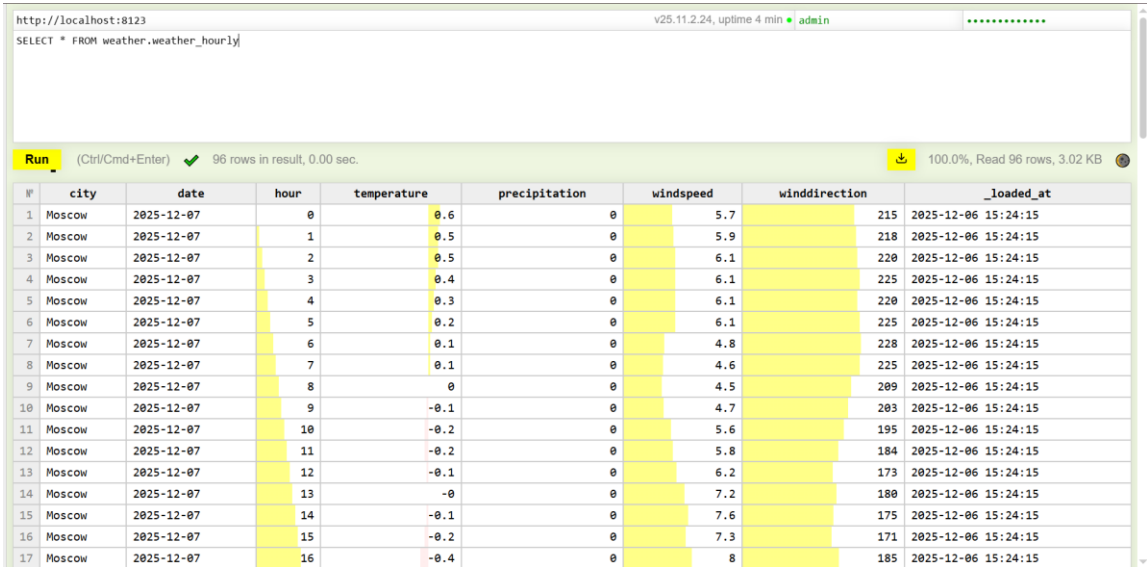


Рисунок 6 – SELECT запрос к таблице weather.weather\_hourly

Проверим, что json-файлы сохранены в MinIO (рисунки 7 и 8)

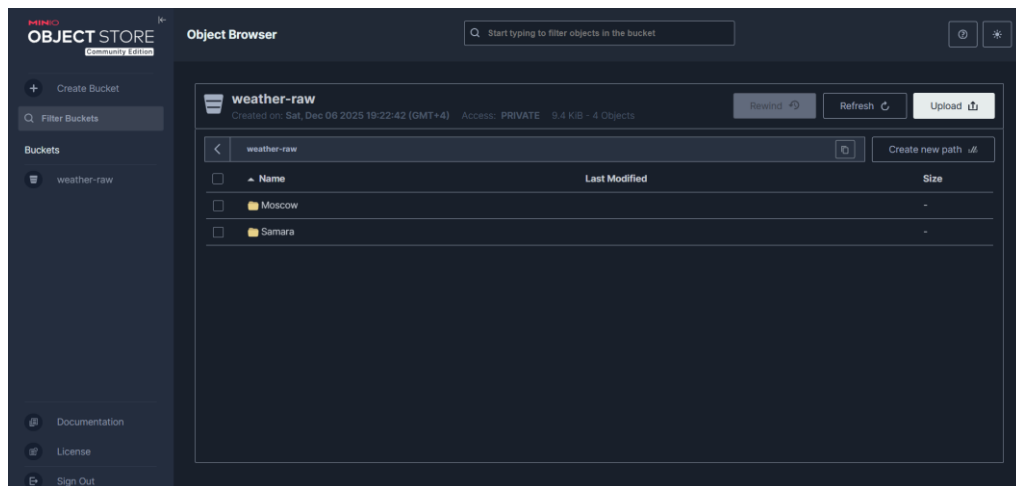


Рисунок 7 – Интерфейс MinIO с папками по городам

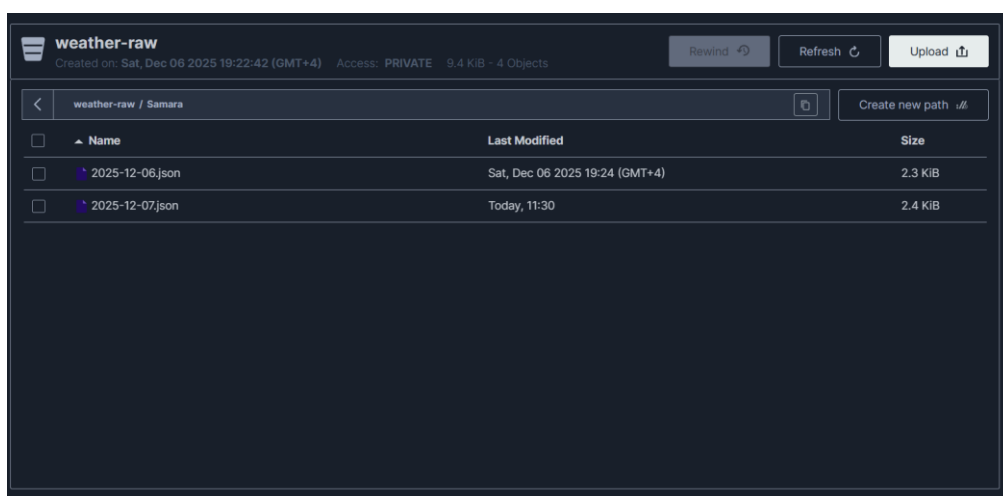


Рисунок 8 – Интерфейс MinIO с json-файлами для Самары

В телеграм также видим сообщения о прогнозе погоды, подтверждающие успех работы пайплайна (рисунок 8).

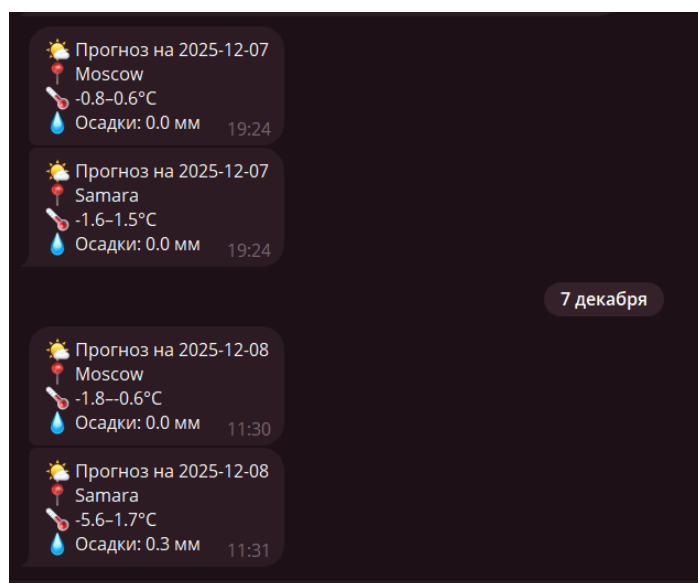


Рисунок 9 – Сообщения от бота в телеграм

## **ВЫВОДЫ**

В ходе выполнения лабораторной работы выполнены все поставленные цели: настроен полный цикл ETL (извлечение из внешнего ресурса, обработка и загрузка в Кликхаус), собран и отлажен пайплайн в Prefect, настроено окружение для работы с данными. Также настроена отправка сообщений через телеграмм бота.

Возникали трудности с настройкой окружения, так как некоторые версии инструментов конфликтовали друг с другом, но данная проблема была решена.

Из улучшений можно добавить выбор городов и определенных параметров в телеграмм боте, а также динамичное изменение времени отправки сообщений (выбирается пользователем во сколько будет приходить сообщение).