

# Feature Store

Семинар 11

# Определение и зачем надо

**Feature Store** – это централизованное хранилище для хранения, управления и предоставления признаков для моделей машинного обучения.

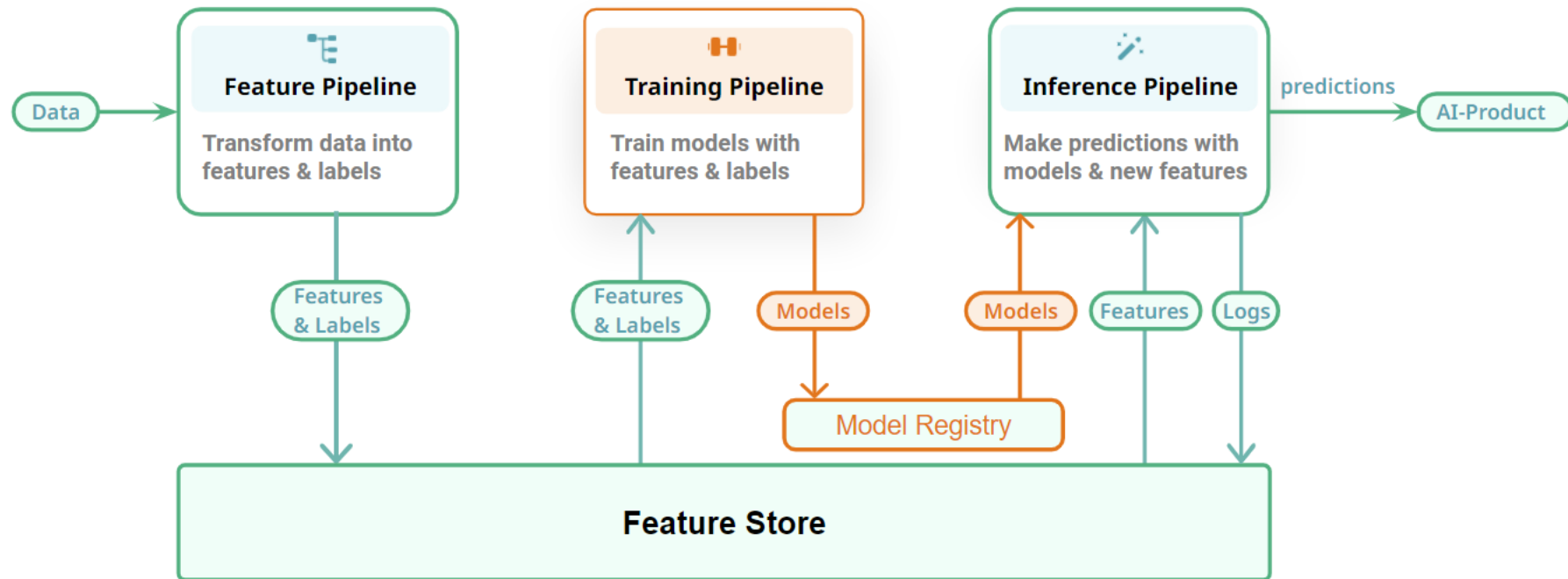
Нужен для:

1. Упрощения процесса извлечения и подготовки данных
2. Повторного использования признаков между проектами
3. Обеспечения согласованности и версионности данных

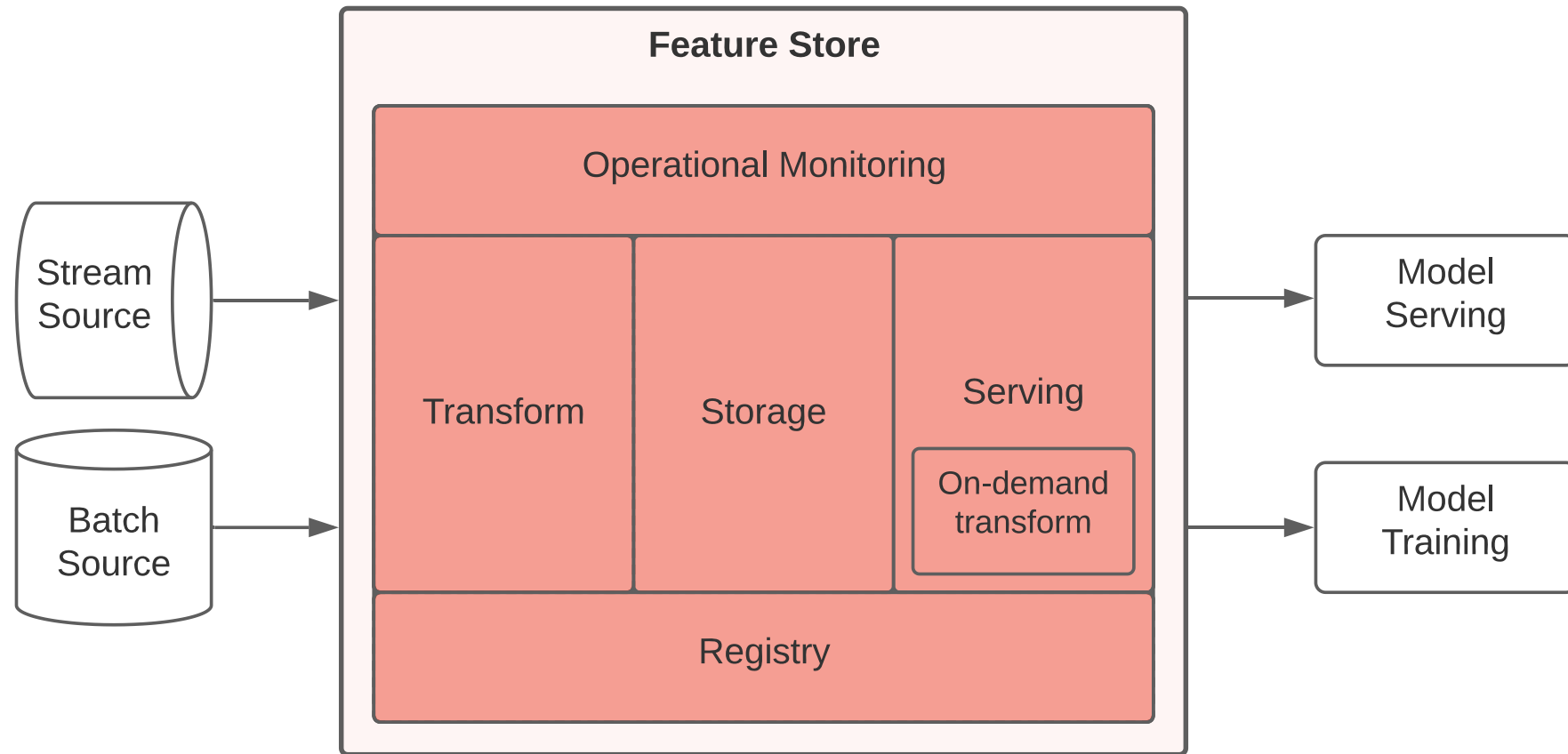
# Где можно почитать

- <https://www.featurestore.org>
- <https://www.hopsworks.ai/dictionary/feature-store>
- <https://habr.com/ru/companies/ glowbyte/articles/581458/>
- <https://feast.dev/blog/what-is-a-feature-store/>

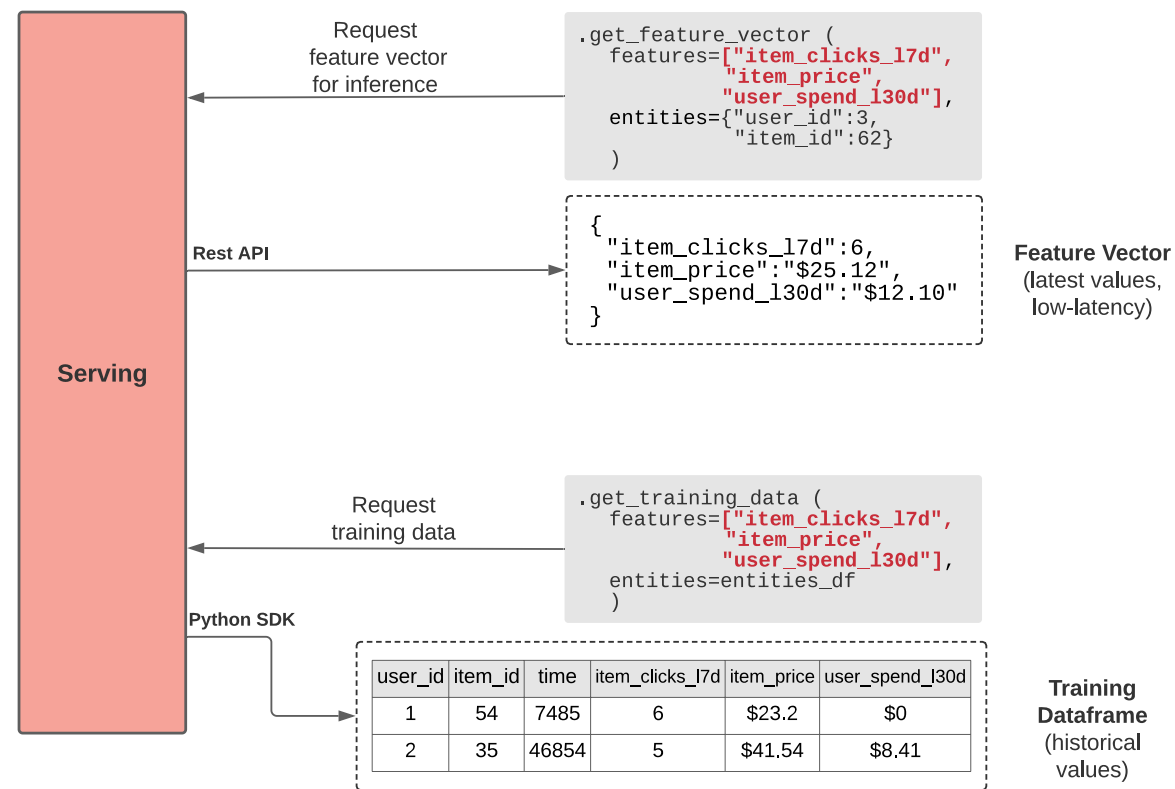
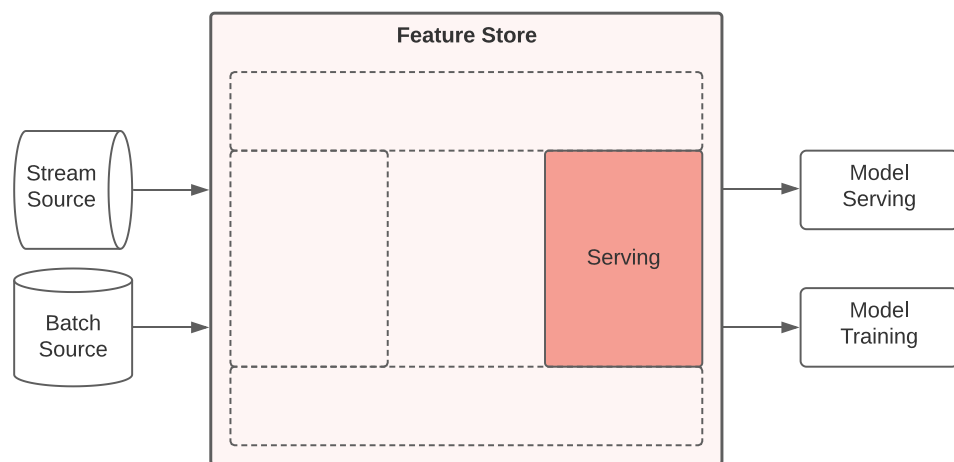
# Как выглядит в процессах



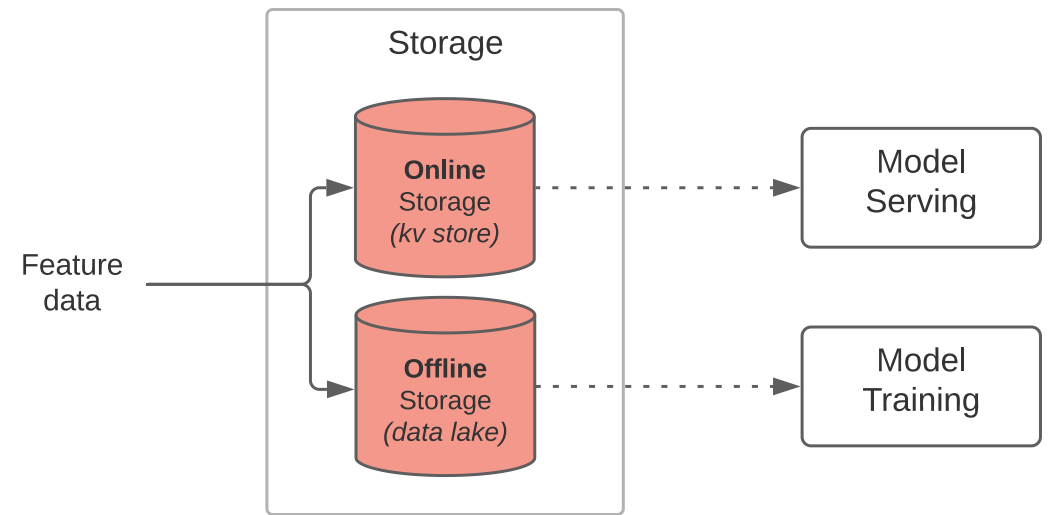
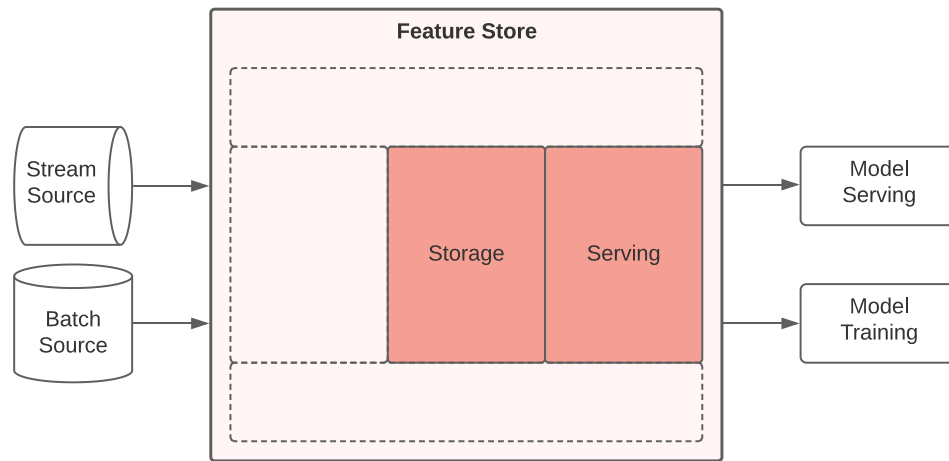
# Компоненты FS



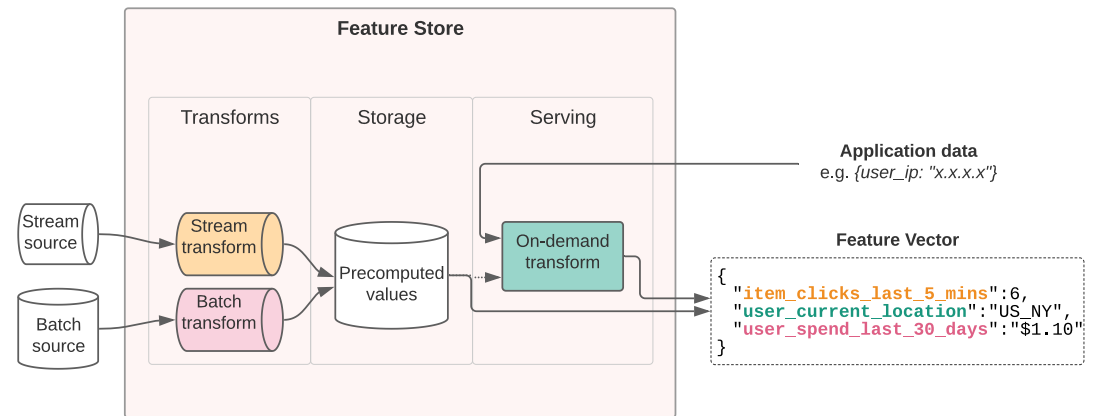
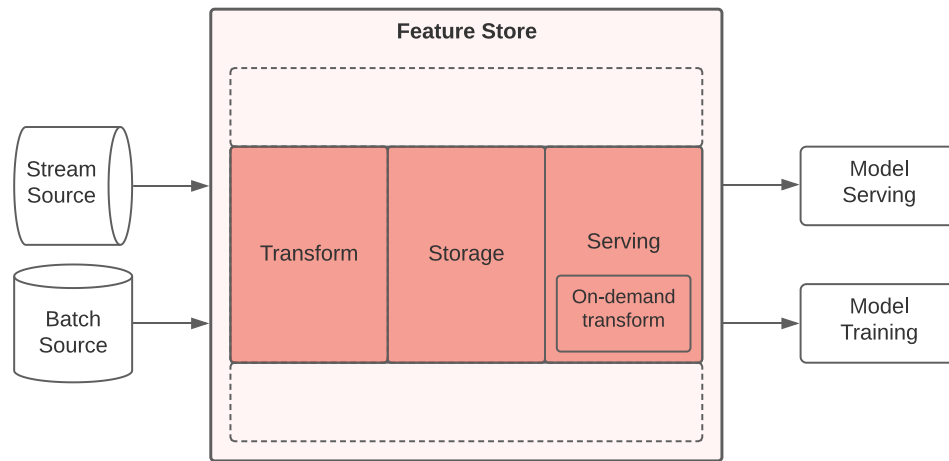
# 1. Serving



## 2. Storage

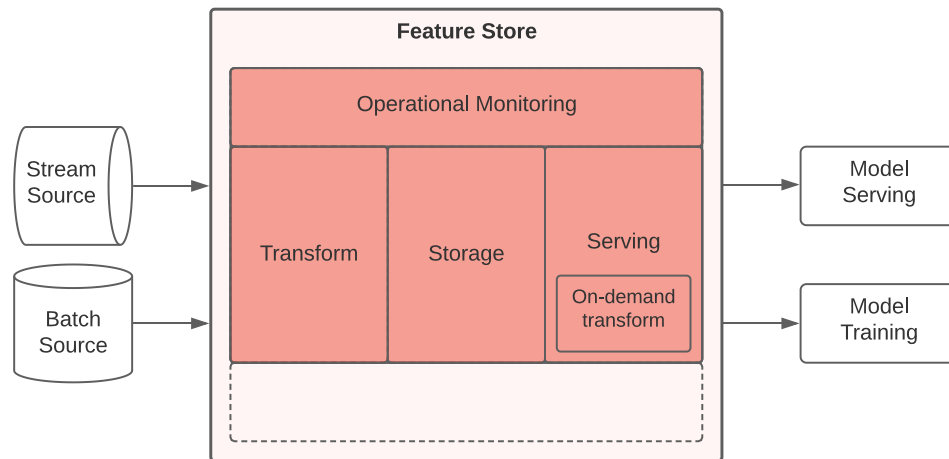


### 3. Transformation





## 4. Monitoring



### Dataset profile

Currently, Feast supports only [Great Expectation's ExpectationSuite](#) as dataset's profile. Hence, the user needs to define a function (profiler) that would receive a dataset and return an [ExpectationSuite](#).

Great Expectations supports automatic profiling as well as manually specifying expectations:

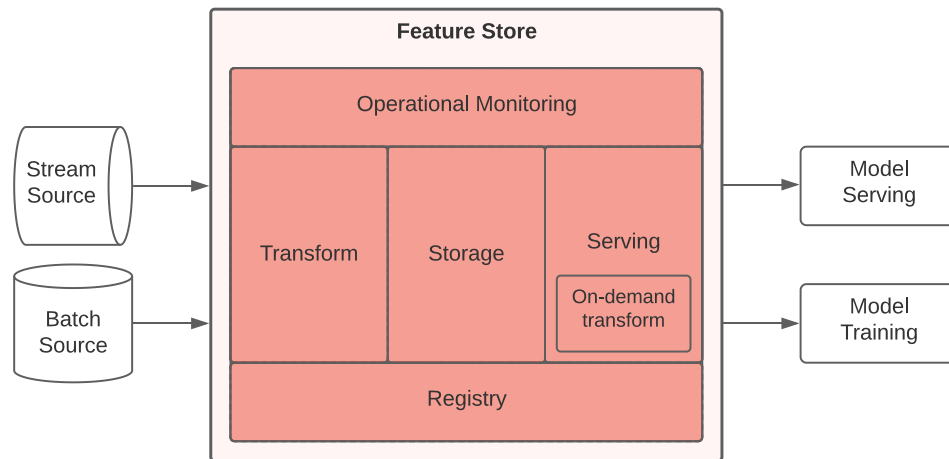
```
from great_expectations.dataset import Dataset
from great_expectations.core.expectation_suite import ExpectationSuite

from feast.dqm.profilers.ge_profiler import ge_profiler

@ge_profiler
def automatic_profiler(dataset: Dataset) -> ExpectationSuite:
    from great_expectations.profile.user_configurable_profiler import UserConfigurablePro

    return UserConfigurableProfiler(
        profile_dataset=dataset,
        ignored_columns=['conv_rate'],
        value_set_threshold='few'
    ).build_suite()
```

## 5. Registry



## Registries

Please see [Registry](#) for a conceptual explanation of registries.

Local



S3



GCS



SQL

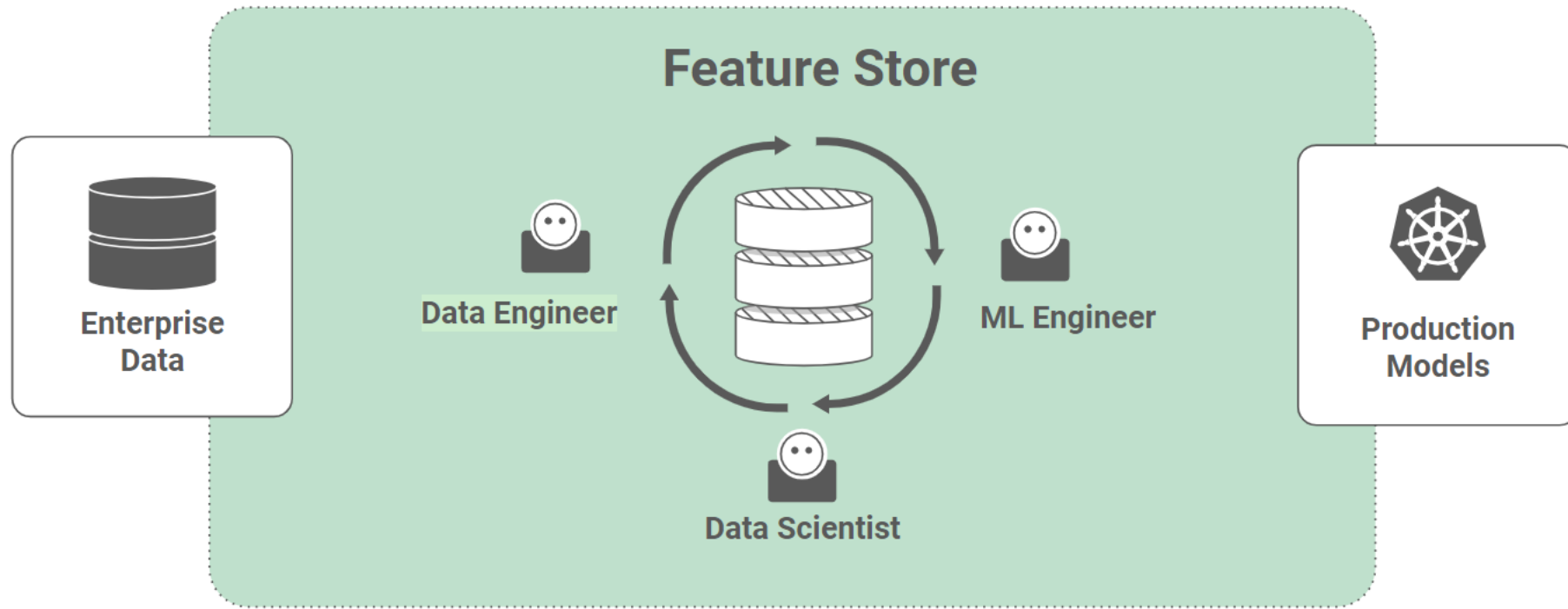


Snowflake



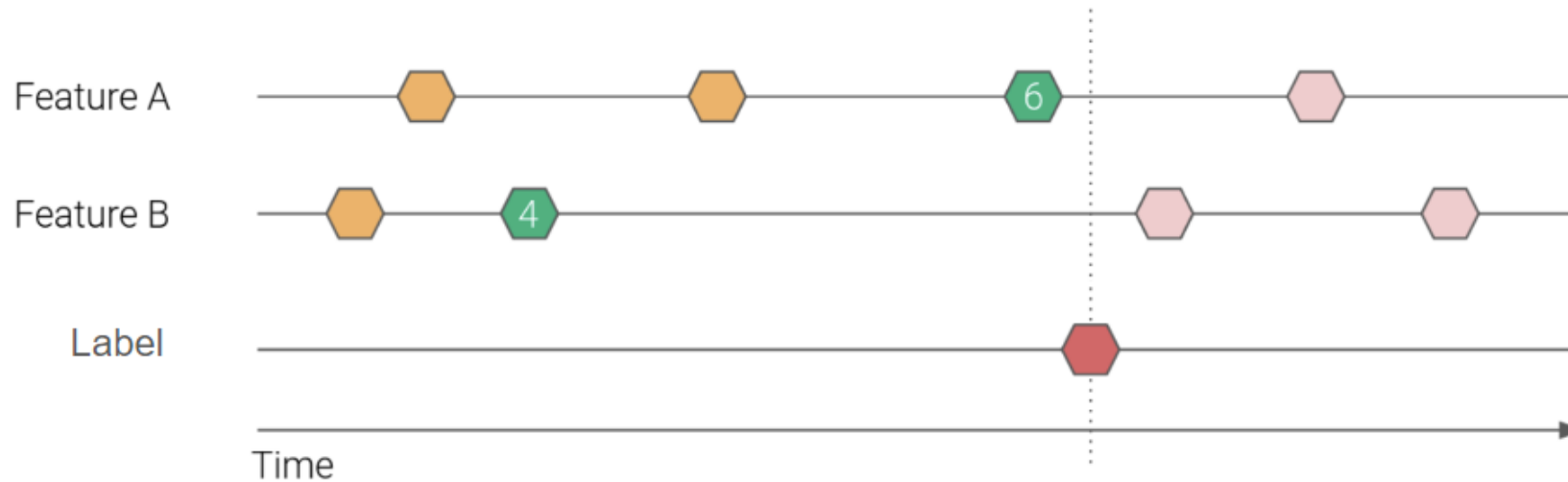
# Централизованный доступ

Cross-Team Collaboration with a Feature Store



# Сбор признаков по таргету

## Point-in-Time Correct Training Data



# Смотрим Feast

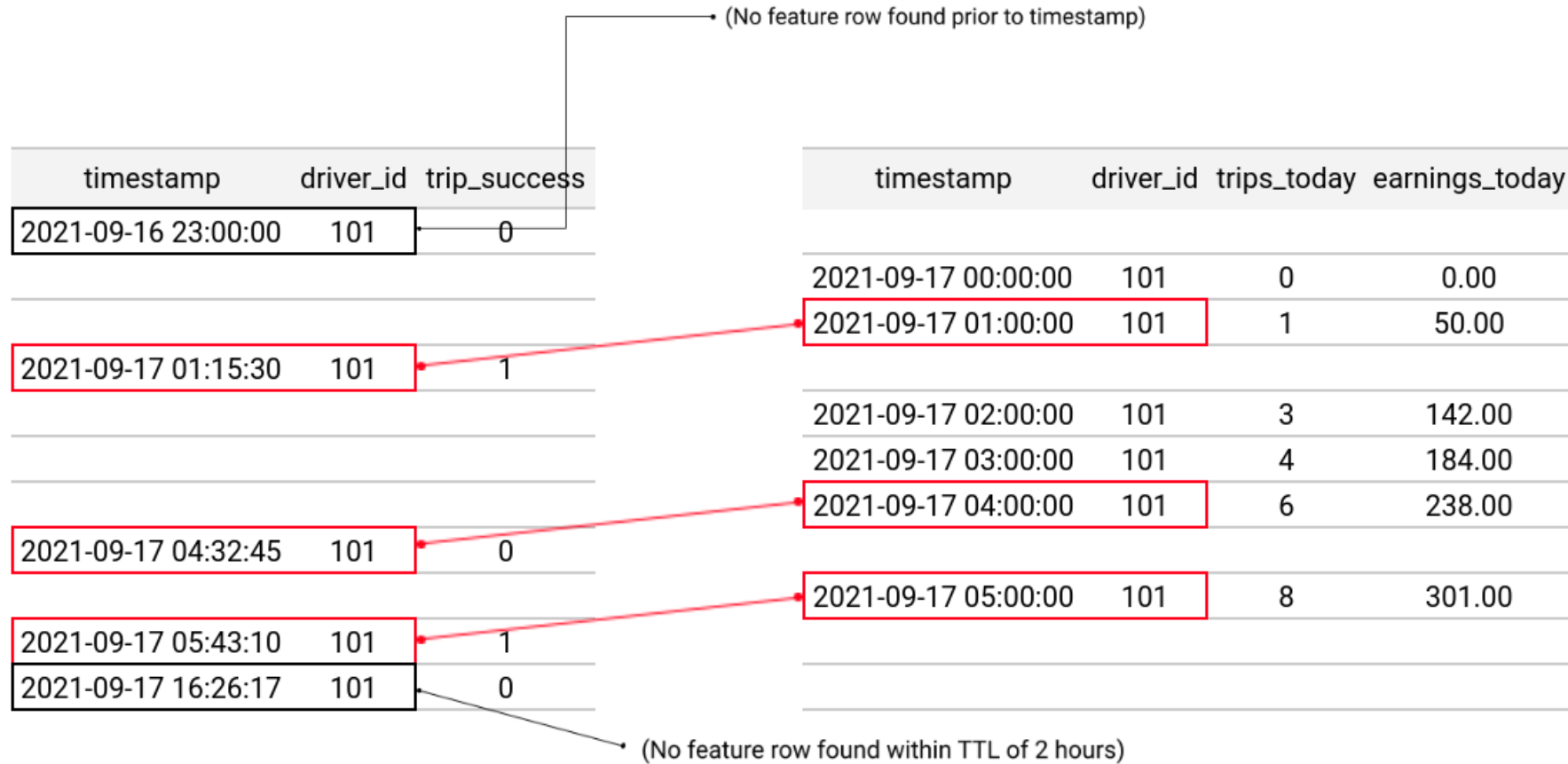
<https://docs.feast.dev/getting-started/quickstart>

# Point-in-time joins

| timestamp           | driver_id | trips_today | earnings_today |
|---------------------|-----------|-------------|----------------|
| 2021-09-17 00:00:00 | 101       | 0           | 0.00           |
| 2021-09-17 01:00:00 | 101       | 1           | 50.00          |
| 2021-09-17 02:00:00 | 101       | 3           | 142.00         |
| 2021-09-17 03:00:00 | 101       | 4           | 184.00         |
| 2021-09-17 04:00:00 | 101       | 6           | 238.00         |
| 2021-09-17 05:00:00 | 101       | 8           | 301.00         |

| timestamp           | driver_id | trip_success |
|---------------------|-----------|--------------|
| 2021-09-16 23:00:00 | 101       | 0            |
| 2021-09-17 01:15:30 | 101       | 1            |
| 2021-09-17 04:32:45 | 101       | 0            |
| 2021-09-17 05:43:10 | 101       | 1            |
| 2021-09-17 16:26:17 | 101       | 0            |

# Point-in-time joins



# Point-in-time joins

| timestamp           | driver_id | trip_success | trips_today | earnings_today |
|---------------------|-----------|--------------|-------------|----------------|
| 2021-09-16 23:00:00 | 101       | 0            | NULL        | NULL           |
| 2021-09-17 01:15:30 | 101       | 1            | 1           | 50.00          |
| 2021-09-17 04:32:45 | 101       | 0            | 6           | 238.00         |
| 2021-09-17 05:43:10 | 101       | 1            | 8           | 301.00         |
| 2021-09-17 16:26:17 | 101       | 0            | NULL        | NULL           |



# Домашняя работа №6 / Финальный проект

## 1 балл:

1. Поднять локально Feast (как в примере сегодня)
2. Собрать данные за 10 таймстемпов (любых, для всех driver\_id) и обучить модель предсказывать колонку **avg\_daily\_trips** по полям **conv\_rate** и **acc\_rate**

## 2 балла:

1. Поднять локально Feast + Airflow
2. Написать DAG для подтягивания актуальных фичей за 10 таймстемпов (любых) для всех driver\_id и обучение модели (как в задаче за 1 балл)

# Домашняя работа №6 / Финальный проект

## **3 балла:**

1. Все то, что было в задании на 2 балла.
2. Написать DAG для батч инференса модели (должен уметь делать инференс как на актуальную дату, так и на любую дату в прошлом). Модель сохранять в MLflow

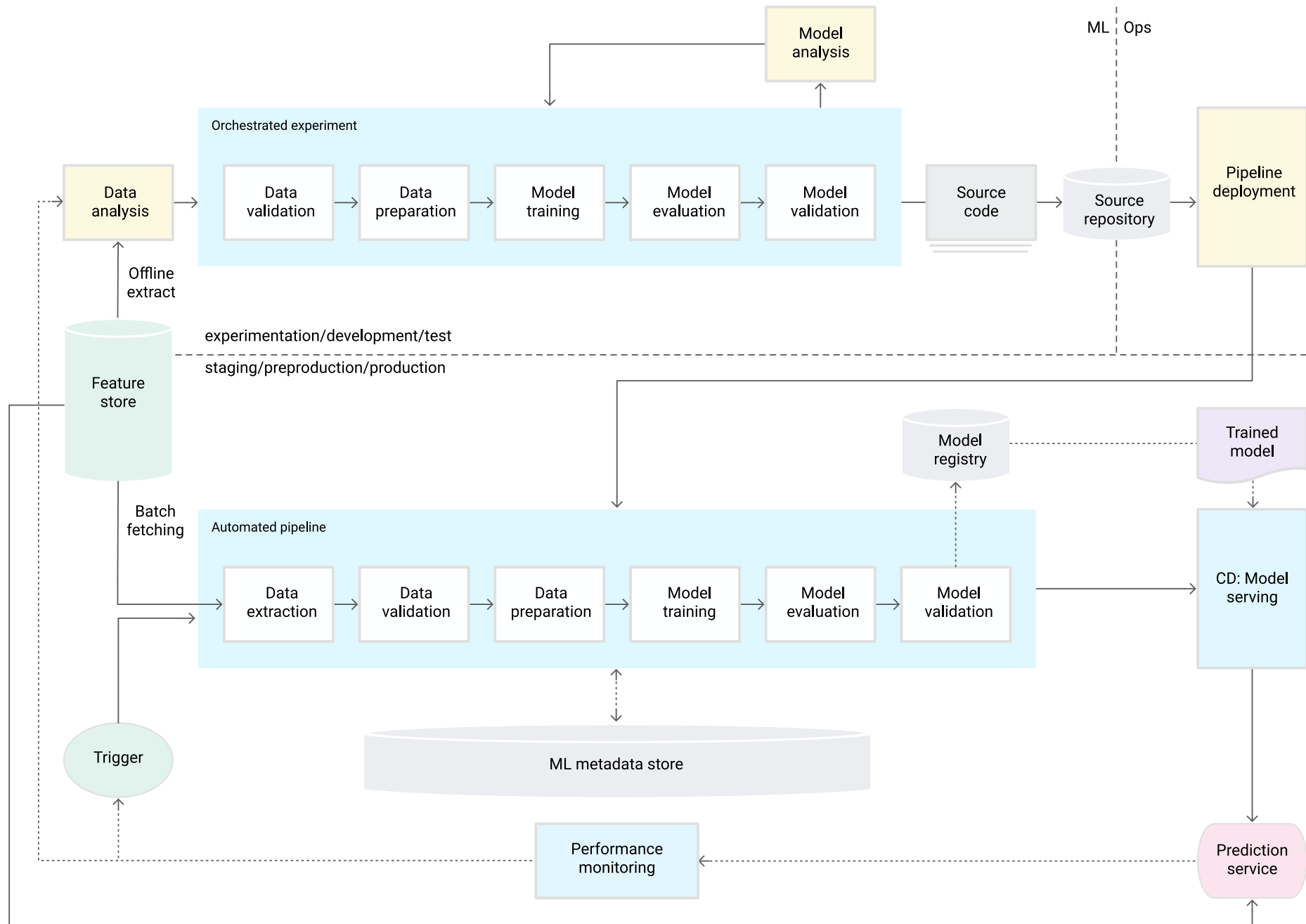
## **4 балла:**

1. Все то, что было в задании на 3 балла.
2. Поднять сервис с онлайн инференсом модели
3. Забирать онлайн фичи из FS и отдавать на вход сервису с онлайн инференсом
4. \* (опционально) прикрутить мониторинг входных и выходных данных

# Домашняя работа №6 / Финальный проект

[https://github.com/feast-dev/feast-workshop/tree/main/module\\_1](https://github.com/feast-dev/feast-workshop/tree/main/module_1)

Для вдохновения



# Итоговые баллы

В итоге 14 баллов за семестр

Для автомата надо сдать **все** лабы хотя бы на один балл

«5» = 11+ баллов

«4» = 8+ баллов

«3» = 6+ баллов