

ОТ ПОТОКОВ ДАННЫХ ДО ВЫСОКОКАЧЕСТВЕННЫХ ML-МОДЕЛЕЙ

Секрет конструирования и тестирования признаков

Обо мне



Юрий Гусев



[linkedin.com/in/ygusev](https://www.linkedin.com/in/ygusev)

- Окончил МГТУ им. Баумана
- Работал Principal Software Engineer в Oracle
- Старший руководитель разработки реалтайм систем кредитного скоринга
- Консультант
- Люблю Stream Processing и Machine Learning
- Работаю над LLM-системами
- Контрибьютор в Apache Flink

Продакшн в теории и в реальности

kaggle

Production

transfer_learning_with_hub.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Table of contents

- Copyright 2018 The TensorFlow Authors. Licensed under the Apache License, Version 2.0 (the "License");
- Transfer learning with TensorFlow Hub
 - Setup
 - An ImageNet classifier
 - Download the classifier
 - Run it on a single image
 - Decode the predictions
 - Simple transfer learning
 - Dataset
 - Run the classifier on a batch of images
 - Download the headless model
 - Attach a classification head
 - Train the model
 - Check the predictions
 - Export your model
 - Learn more
- Section

Transfer learning with TensorFlow Hub

TensorFlow Hub is a repository of pre-trained TensorFlow models.

This tutorial demonstrates how to:

- Use models from TensorFlow Hub with `tf.keras`
- Use an image classification model from TensorFlow Hub
- Do simple transfer learning to fine-tune a model for your own image classes

Setup

```
[53] import numpy as np
import time

import PIL.Image as Image
import matplotlib.pyplot as plt

import tensorflow as tf
import tensorflow_hub as hub
```

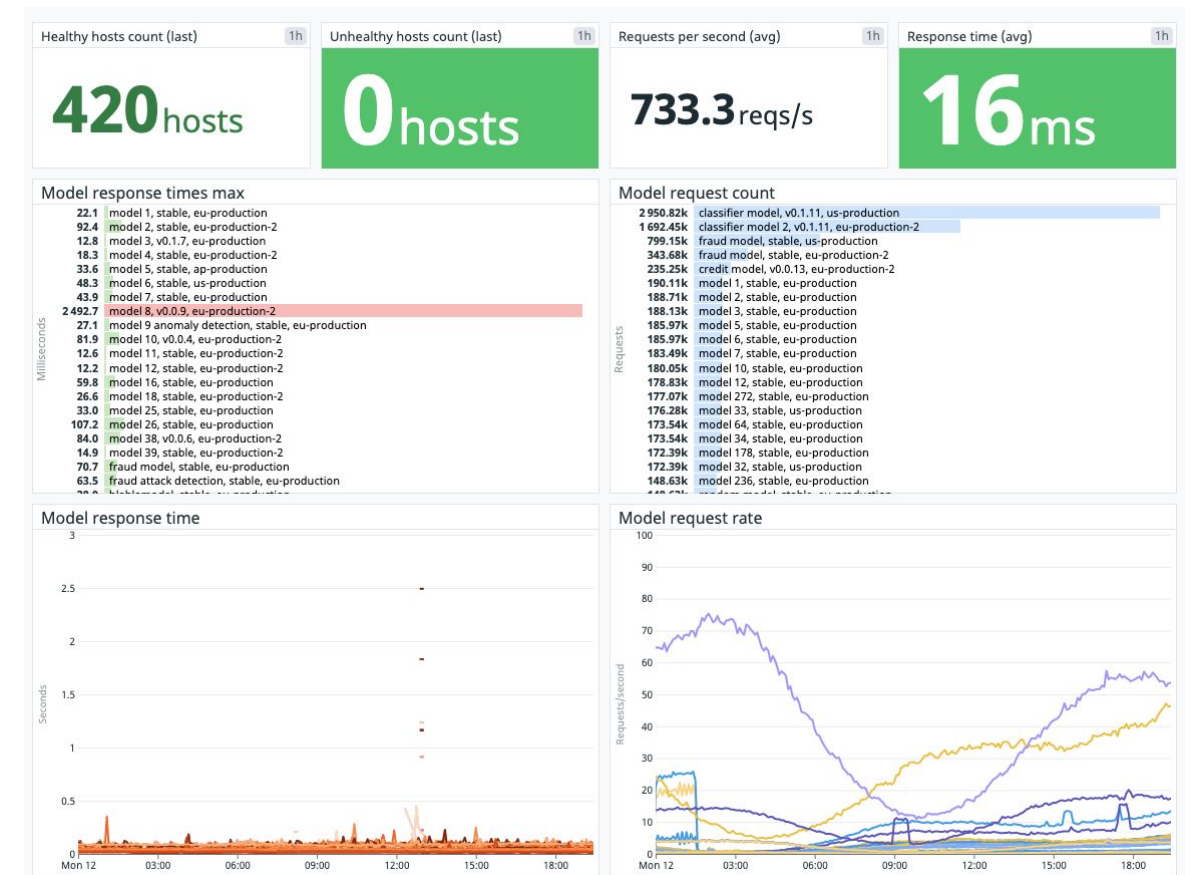
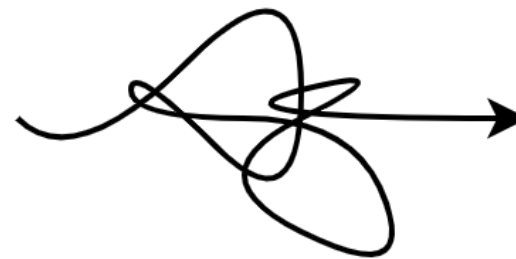
An ImageNet classifier

You'll start by using a pretrained classifier model to take an image and predict what it's an image of - no training required!

Download the classifier

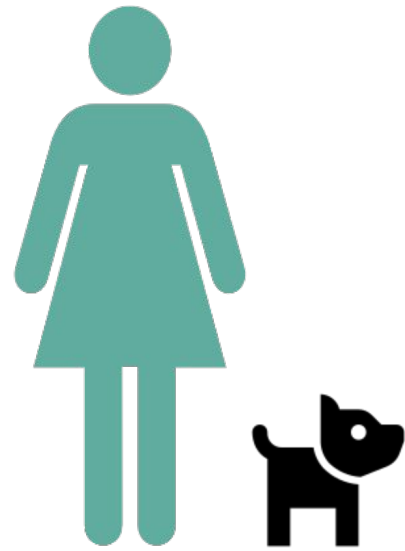
Use `hub.KerasLayer` to load a [MobileNetV2 model](#) from TensorFlow Hub. Any [compatible image classifier model](#) from [here](#).

```
[54] classifier_model = hub.KerasLayer("https://tfhub.dev/google/tf2-preview/mobilenet_v2/classification/4")
```



Задача

История



Мария



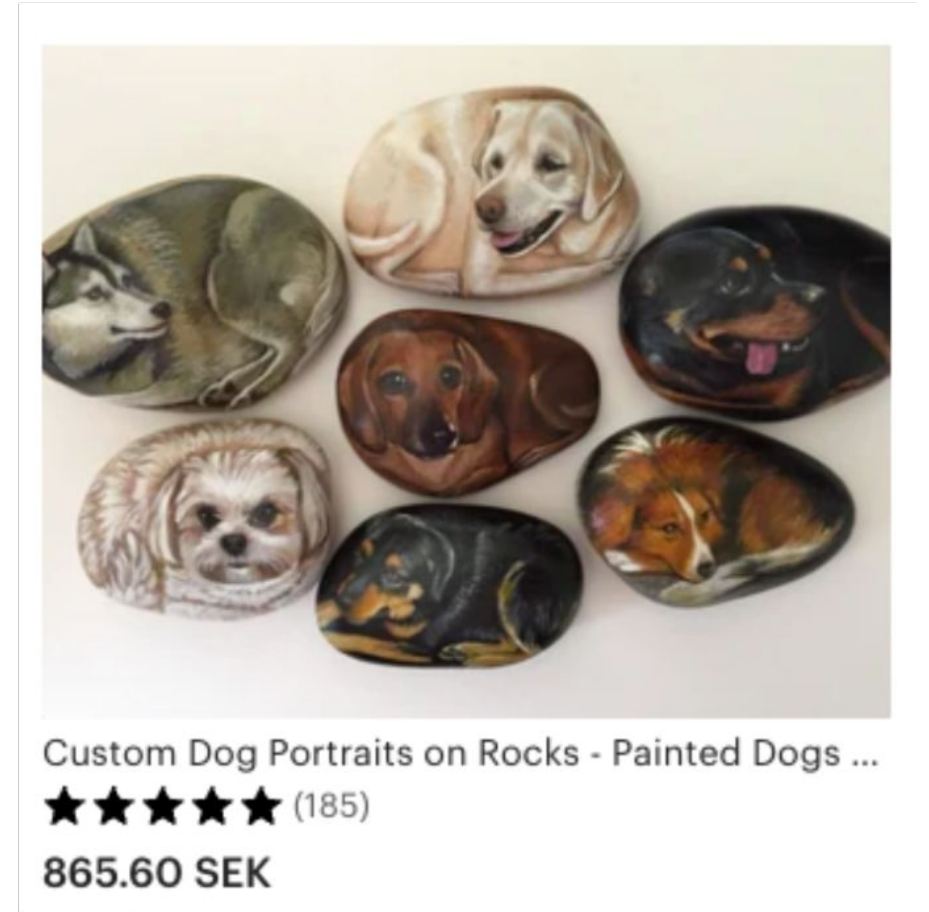
София



Одобрим?



Наша компания



Задача

На основе данных о клиенте мы хотим разработать **статистическую модель** для предсказания **вероятности** того, что клиент не выполнит обязательства по кредиту.



Данные



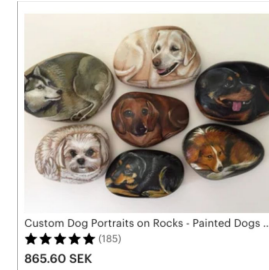
Наши данные

- Покупки
- Баланс
- Споры
- График оплаты
- Невозврат средств



Внешние данные

- Кредитный рейтинг



Данные продавца

- Сумма покупки
- Тип продукта
- Платежный инструмент
- Название продавца



Дизайн системы. Теория



Потоки данных

СЛОЙ РАСЧЕТА ПРИЗНАКОВ В РЕАЛЬНОМ ВРЕМЕНИ

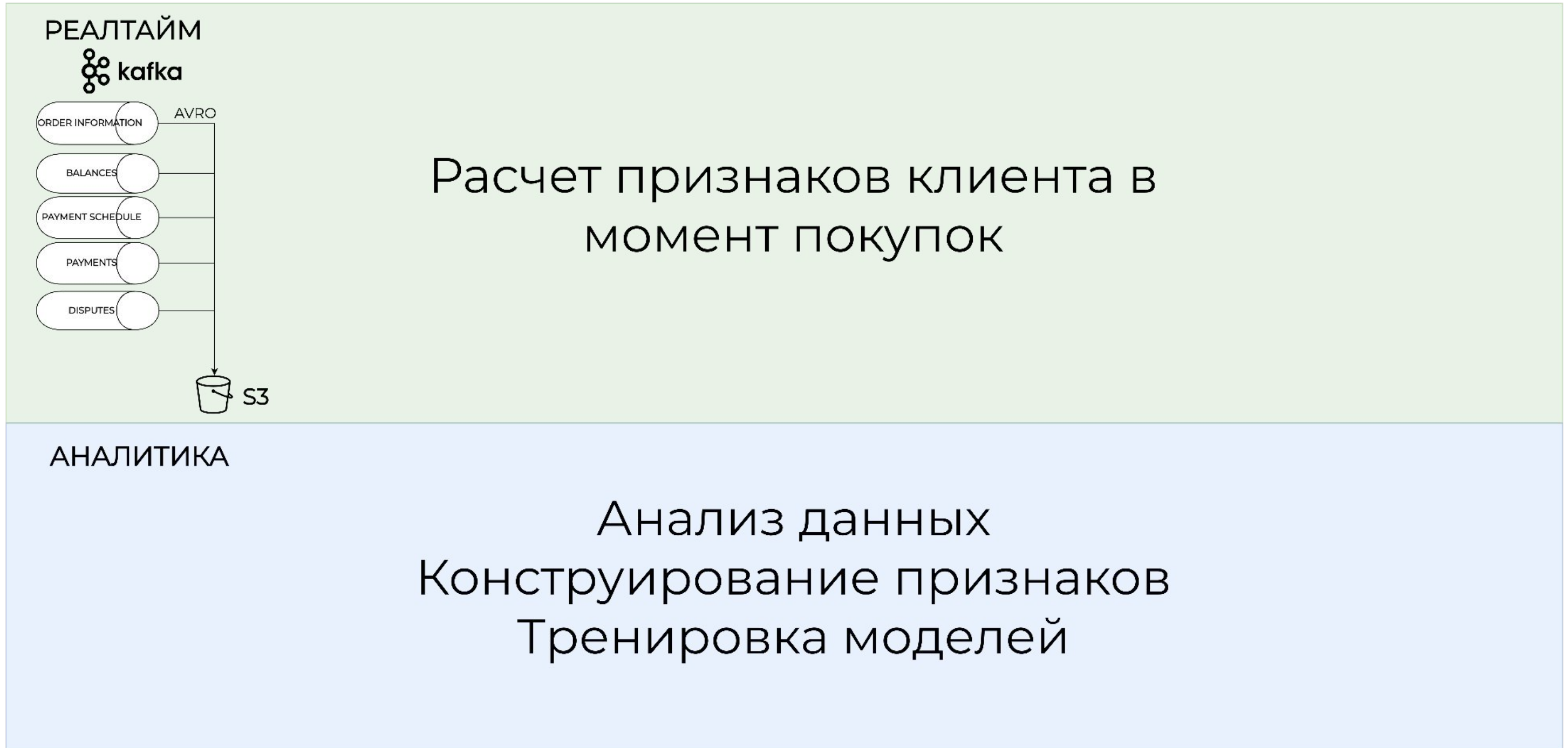
Расчет признаков клиента в
момент покупок

СЛОЙ АНАЛИТИКИ

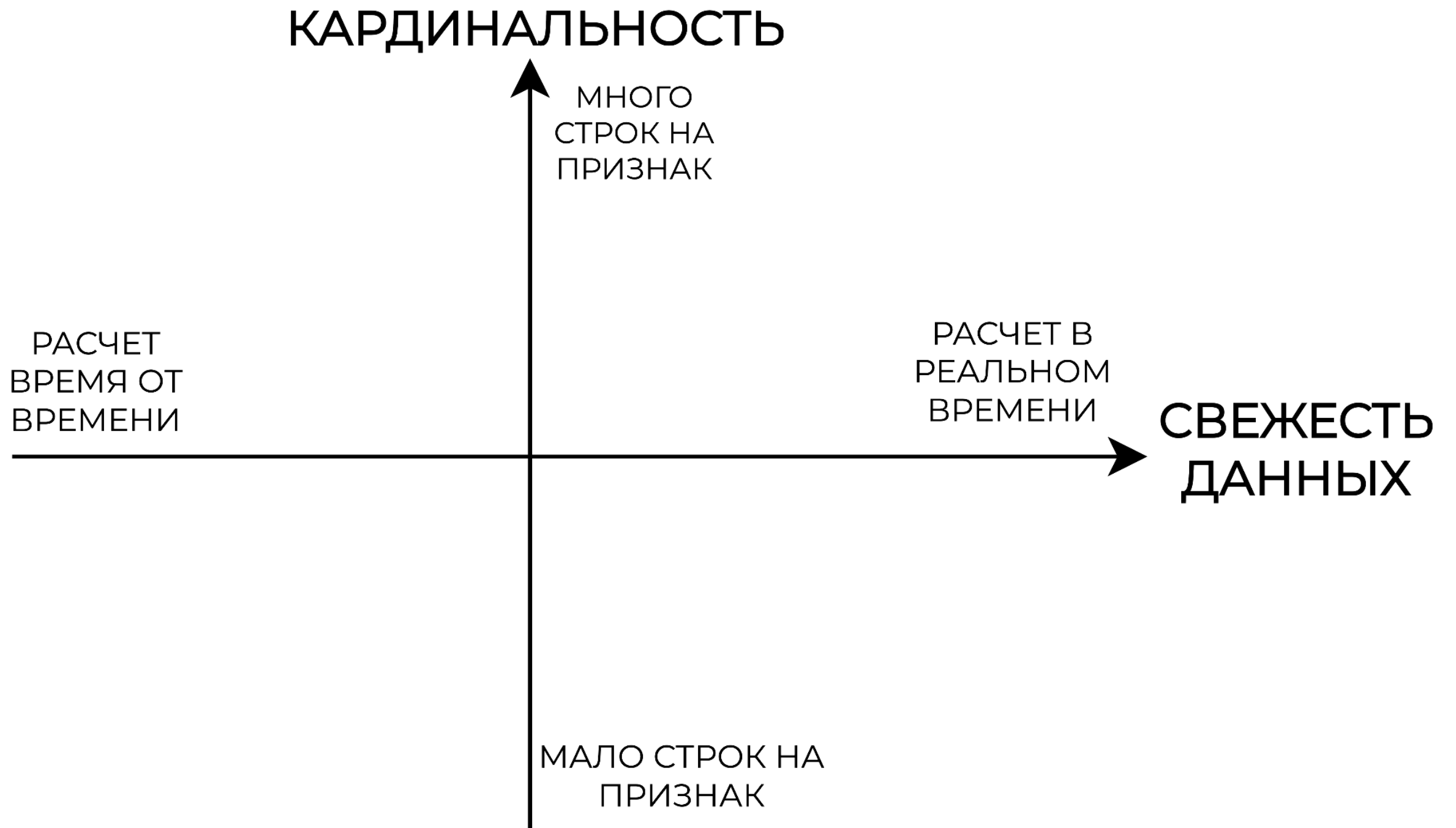
Анализ данных
Конструирование признаков
Тренировка моделей



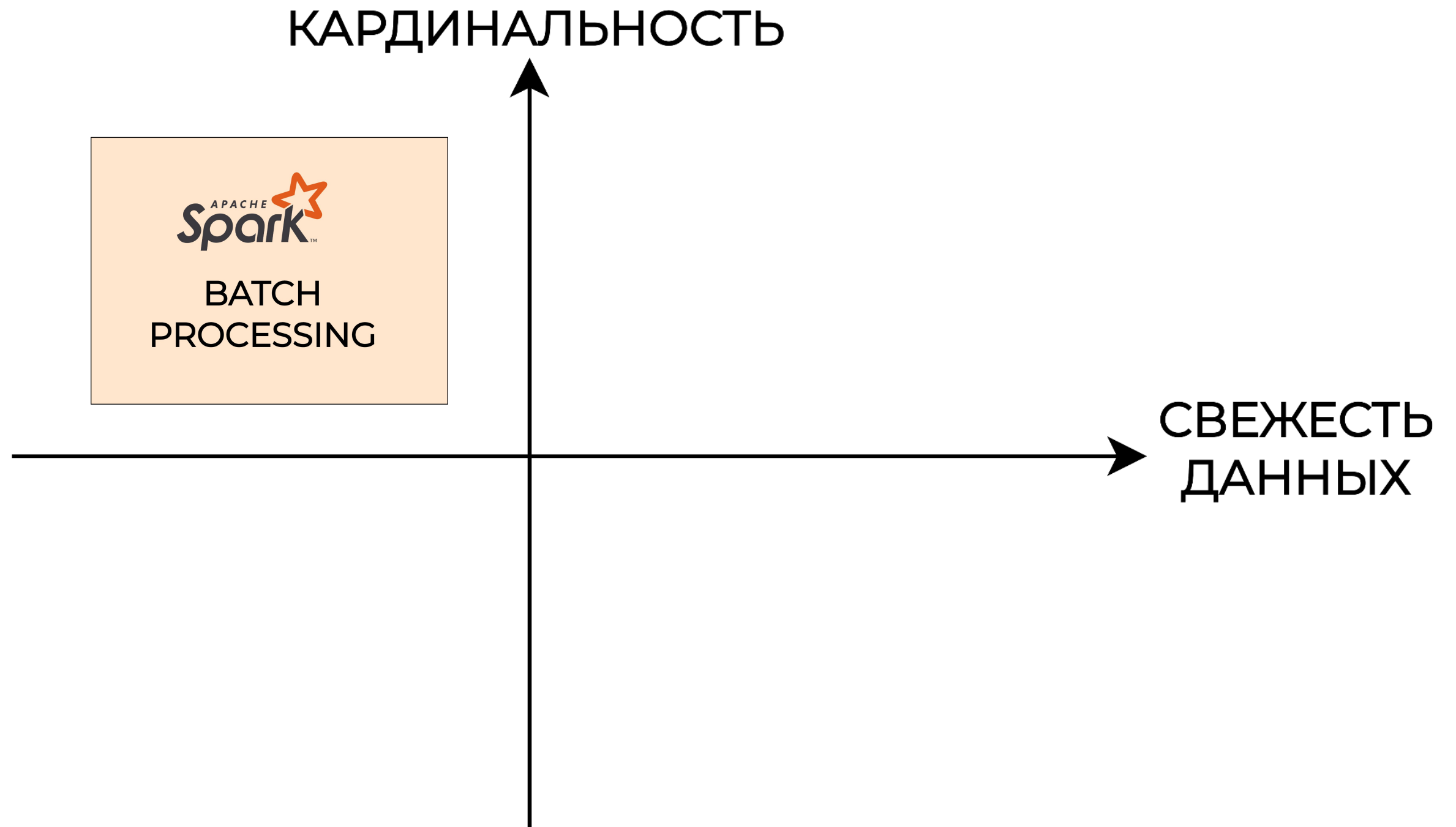
Потоки данных



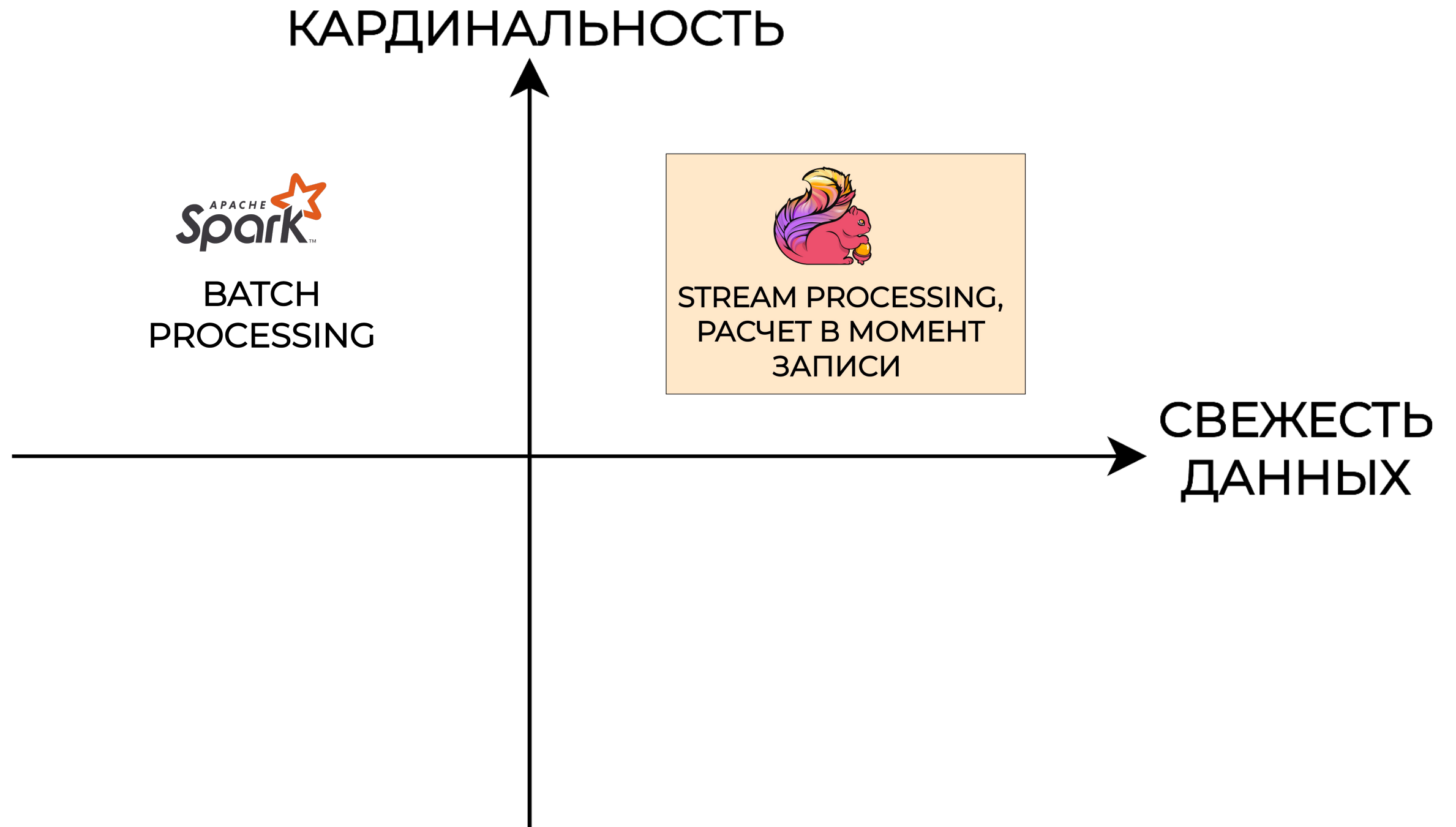
Признаки



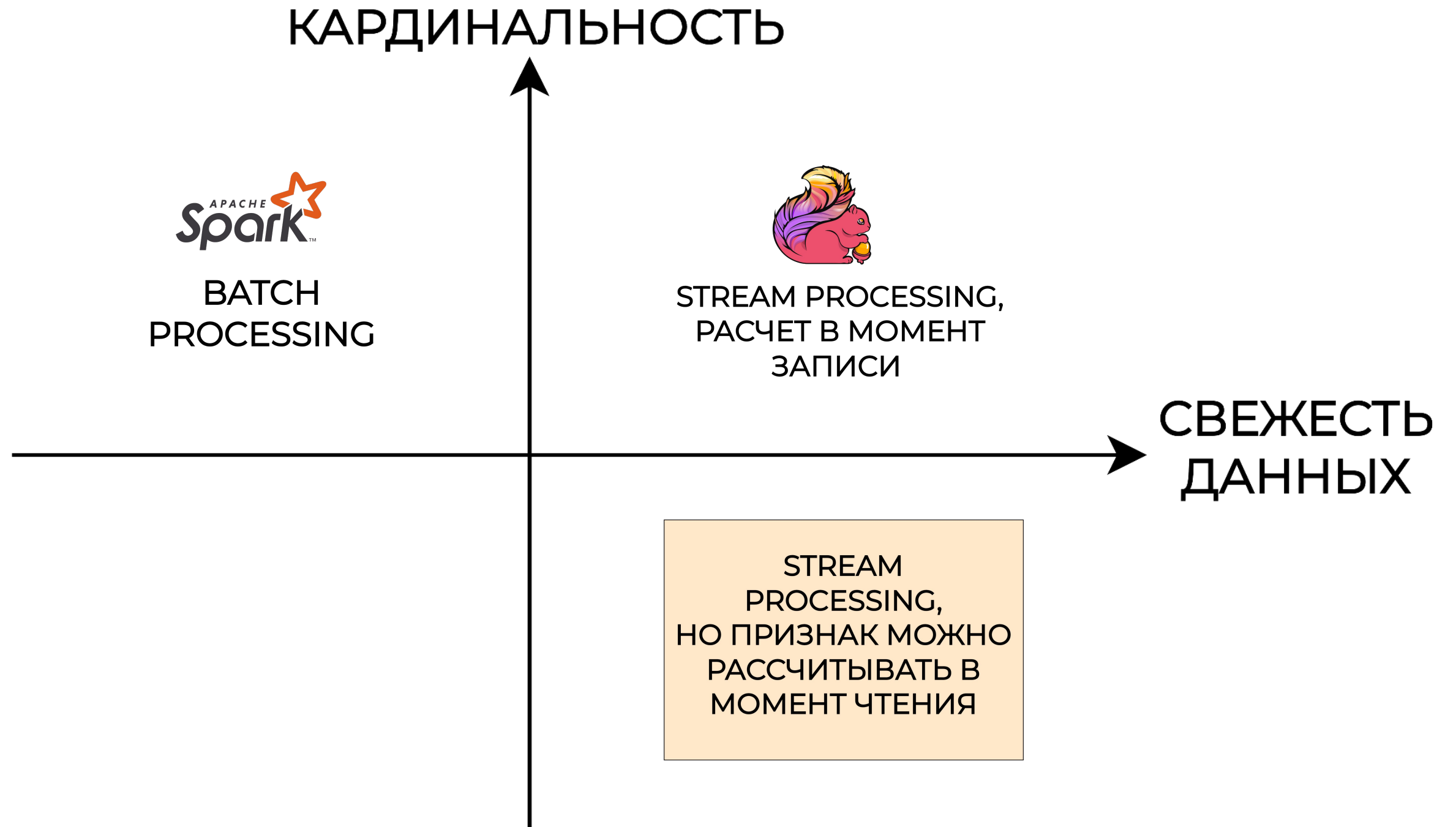
Признаки



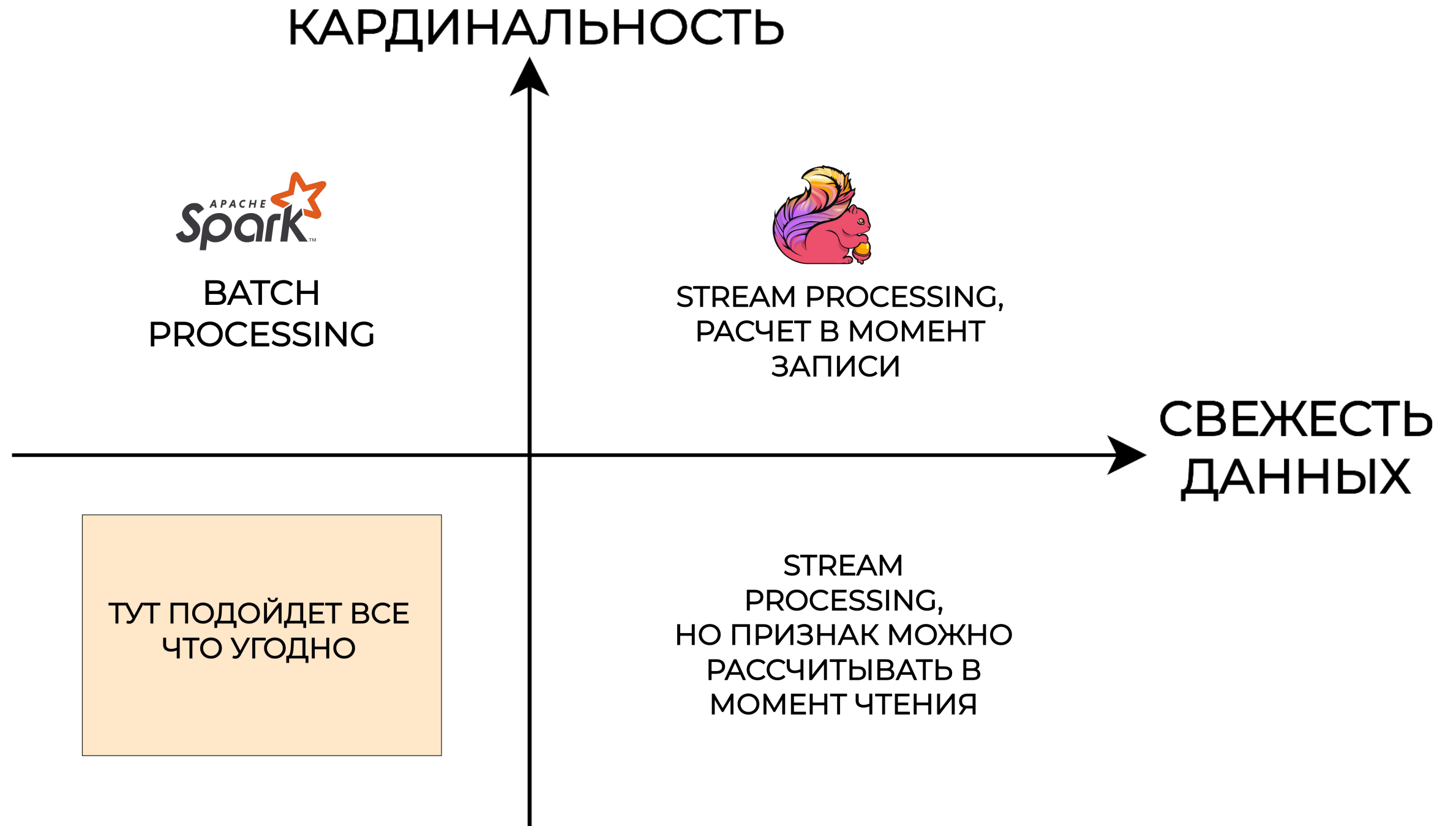
Признаки



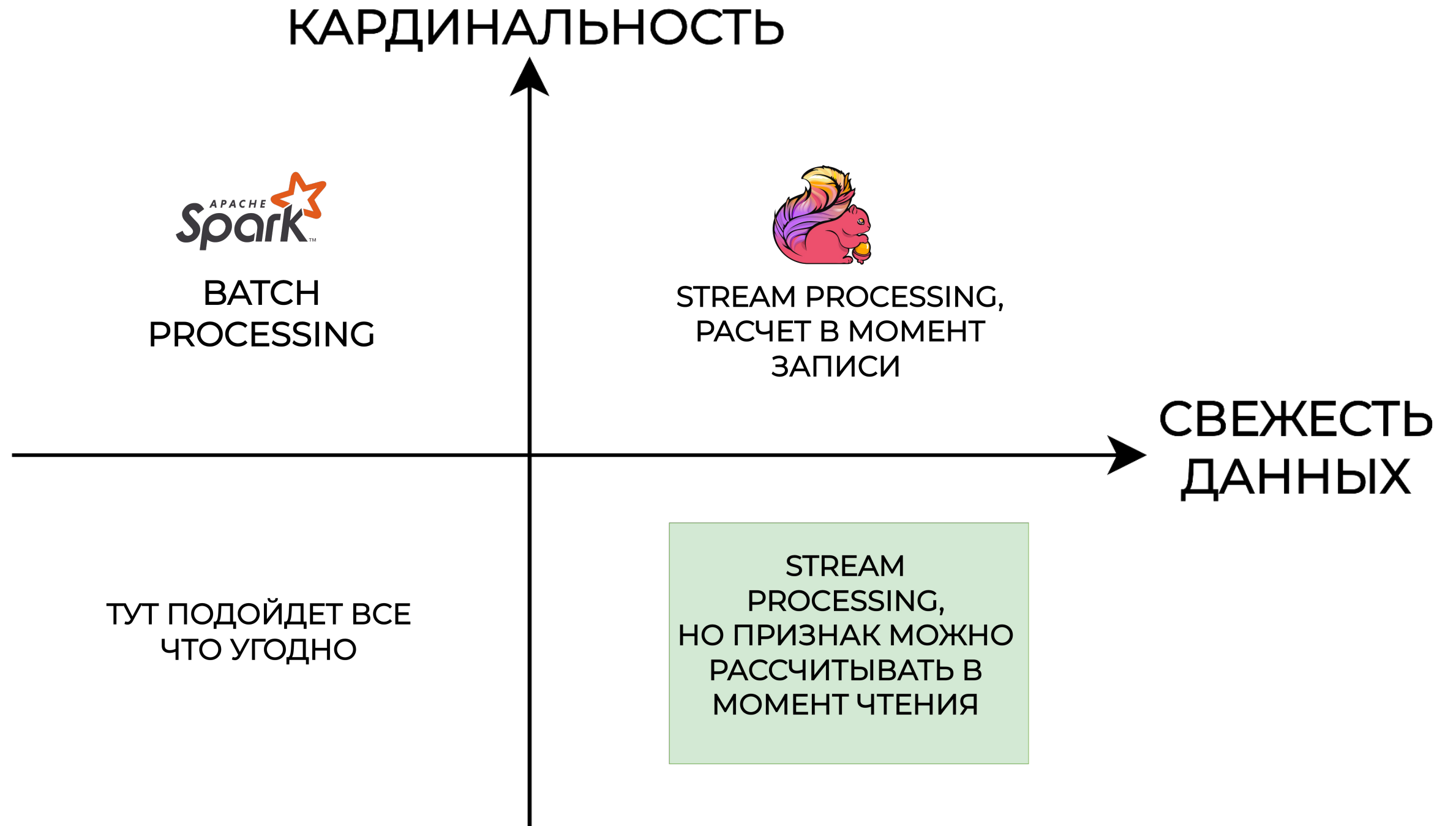
Признаки



Признаки



Признаки

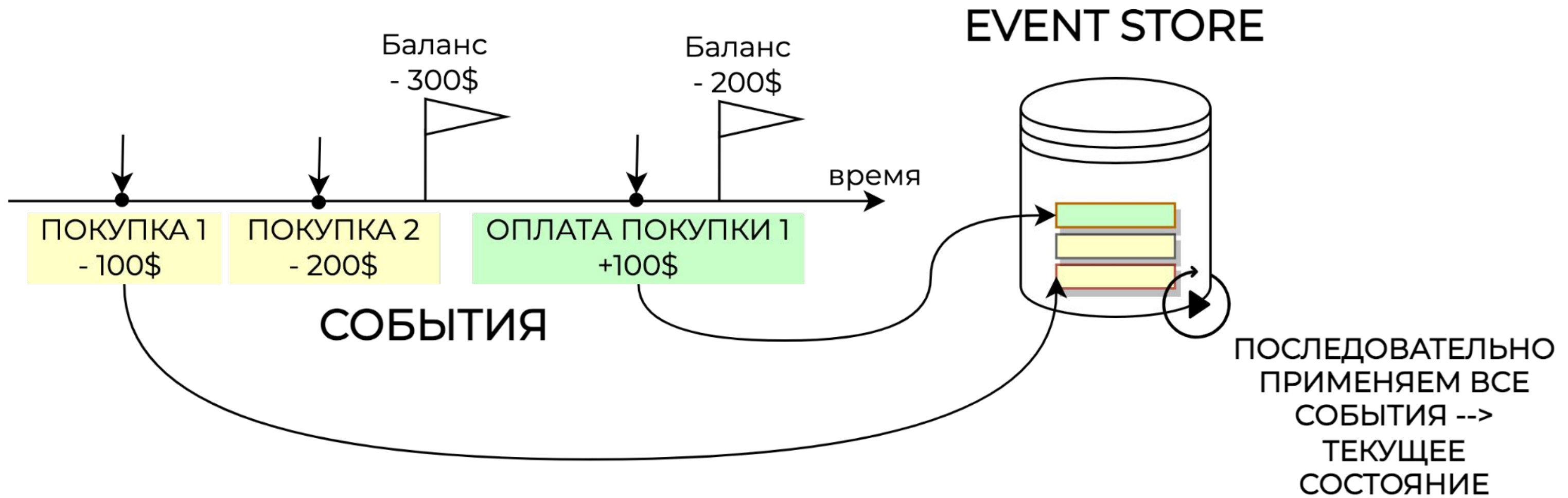


Event sourcing

- Модель, предсказывающая вероятность невыполнения обязательств по кредиту
- Модель надо тренировать на исторических данных
- Надо знать, как менялось наше знание о клиенте в течение времени



Event sourcing



CQRS

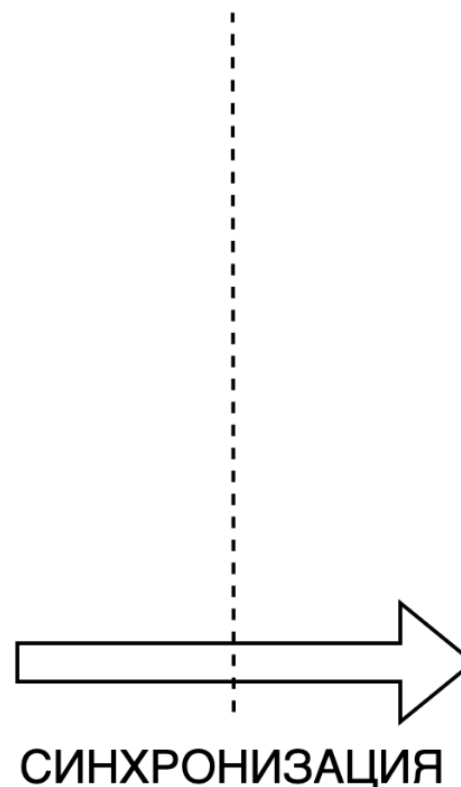
ЛЕГКО
ЗАПИСЫВАТЬ

EVENT STORE



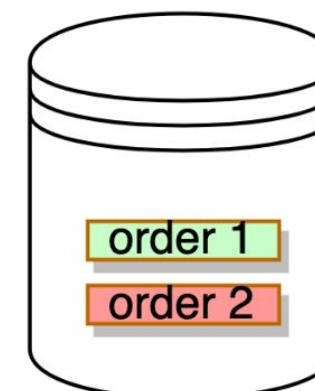
ИСТОРИЯ

Hash Key: Order ID



ЛЕГКО
ЧИТАТЬ

STATE STORE

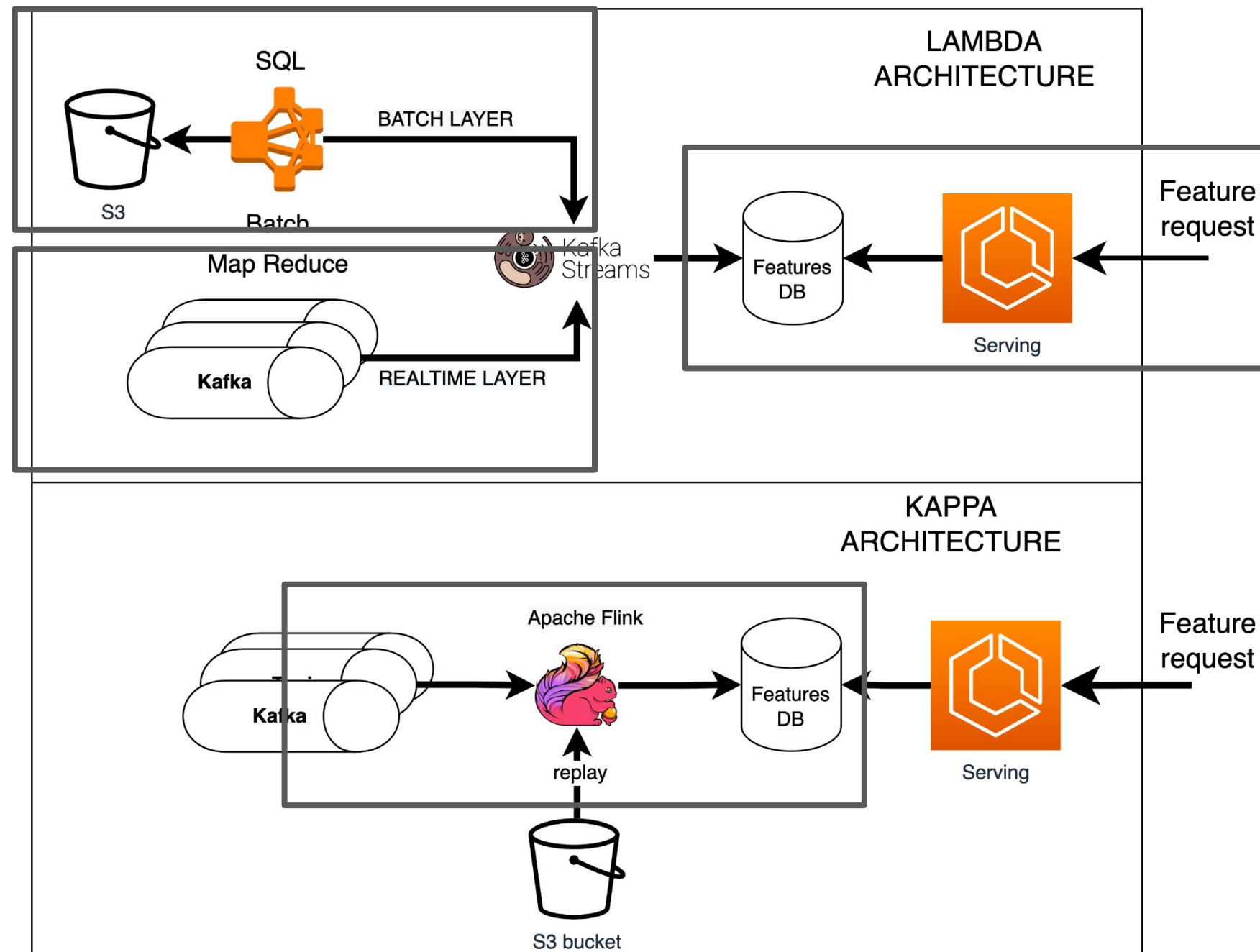


ТЕКУЩЕЕ СОСТОЯНИЕ
КАЖДОЙ ПОКУПКИ

Hash Key: Customer ID



Lambda и Карра архитектуры



Принятые решения

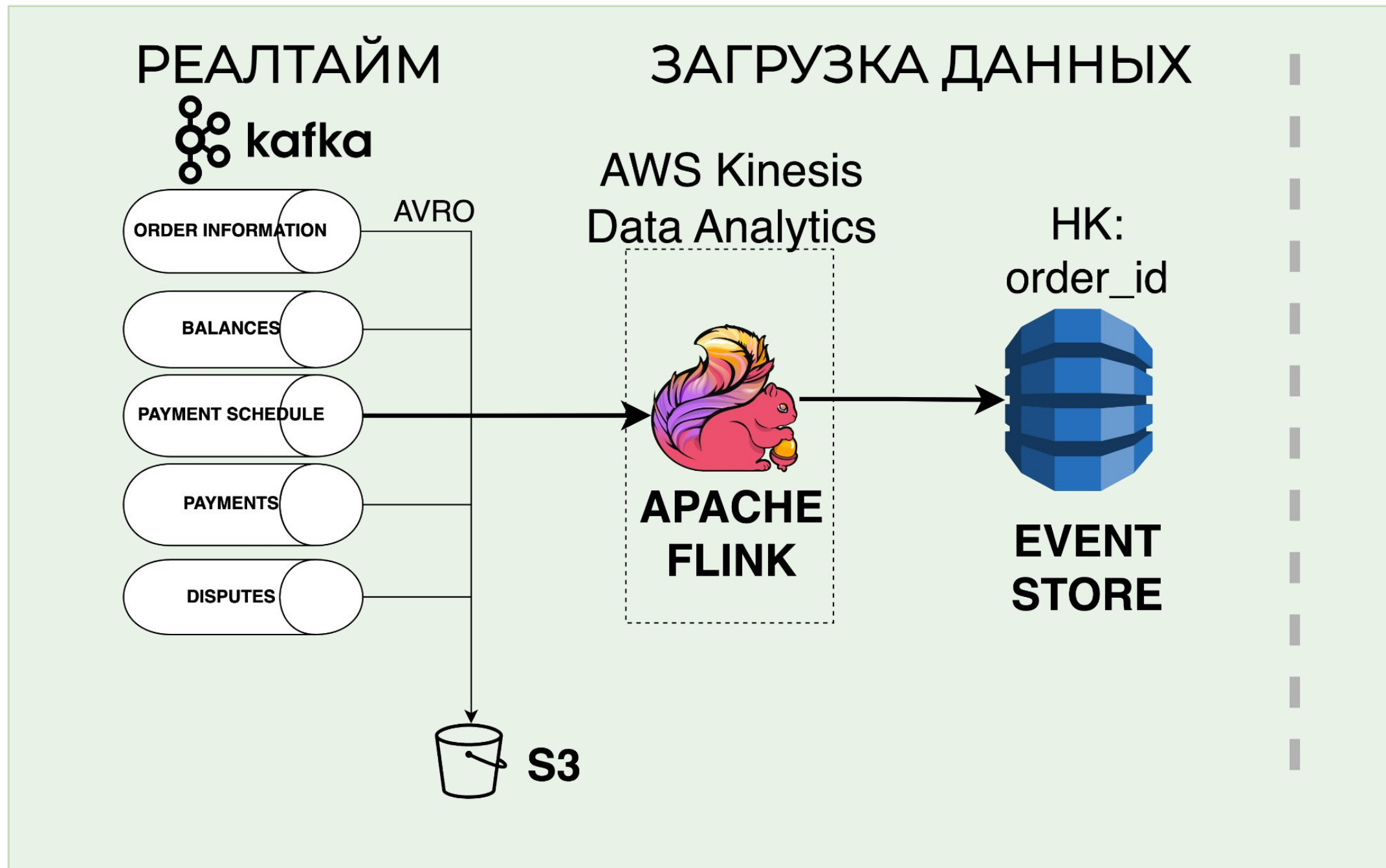
- Используем Event Sourcing и CQRS
- Выбираем архитектуру Карра



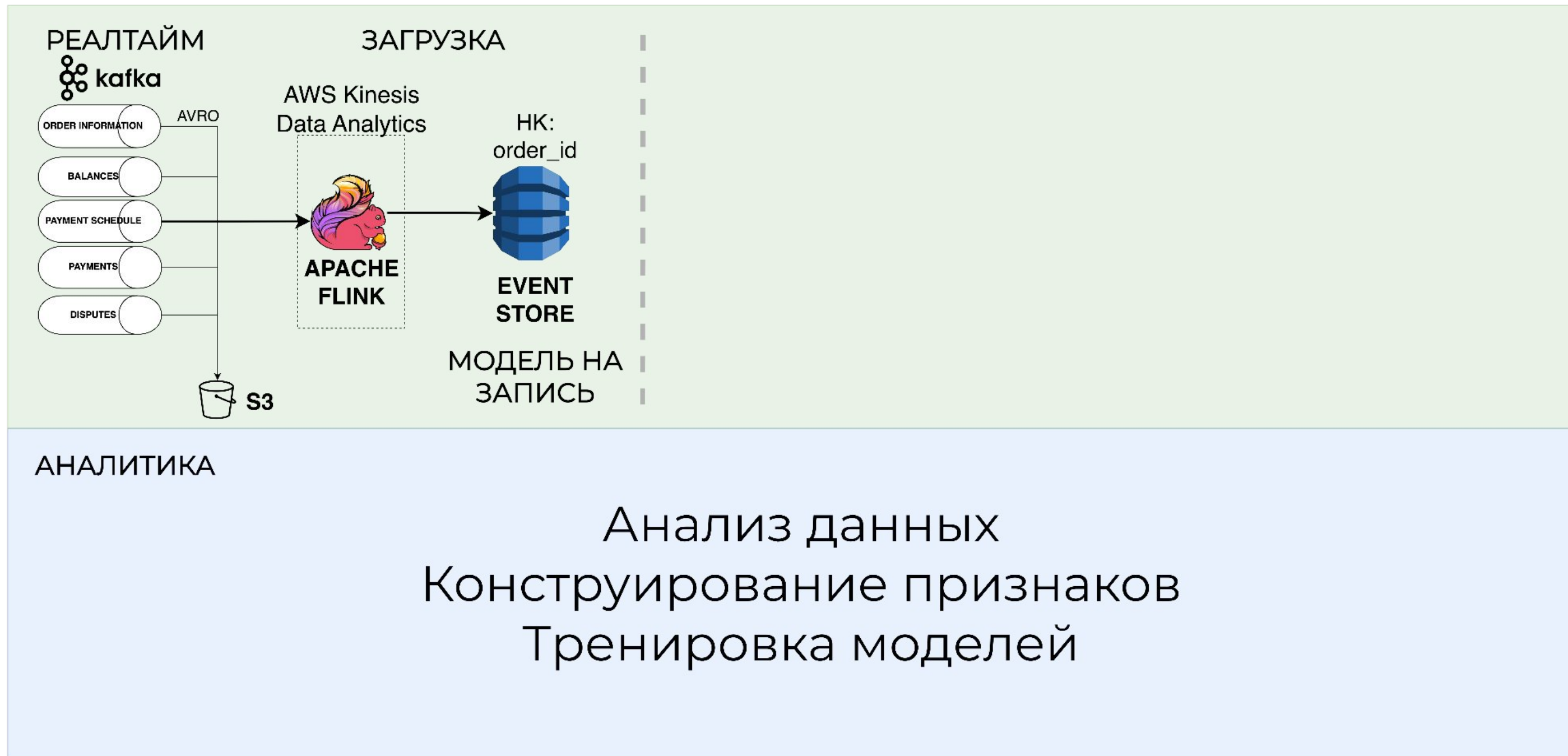
Дизайн системы. Практика



Модель на запись



Модель на запись



Модель на ЧТЕНИЕ

```
@deserialize
@dataclass
class OrderState(AvroModel):
    customer_id: None
    order_information: Optional[OrderInformation] = field(default=None)
    balances: Dict[str, Balance] = field(default_factory=dict)
    payments_schedule: Dict[str, PaymentsSchedule]
        = field(default_factory=dict)
    disputes: Dict[str, Dispute] = field(default_factory=dict)

@deserialize
@dataclass
class Dispute(AvroModel):
    id: str = field()
    balance_id: str = field()
    opened_at: Optional[int] = field(default=None)
    closed_at: Optional[int] = field(default=None)
    reason: Optional[str] = field(default=None)
```



Модель данных. Обработчики сообщений

```
def handle_dispute_opened_event(arguments: Arguments) -> OrderState:
    dispute_id: str = arguments.event['dispute_id']
    balance_id: str = arguments.event['balance_id']
    opened_at: int = int(arguments.event['opened_at'])

    dispute: Dispute = Dispute(
        dispute_id=dispute_id,
        balance_id=balance_id,
        reason=arguments.event['reason'],
        opened_at=opened_at
    )

    arguments.state.disputes[dispute_id] = dispute

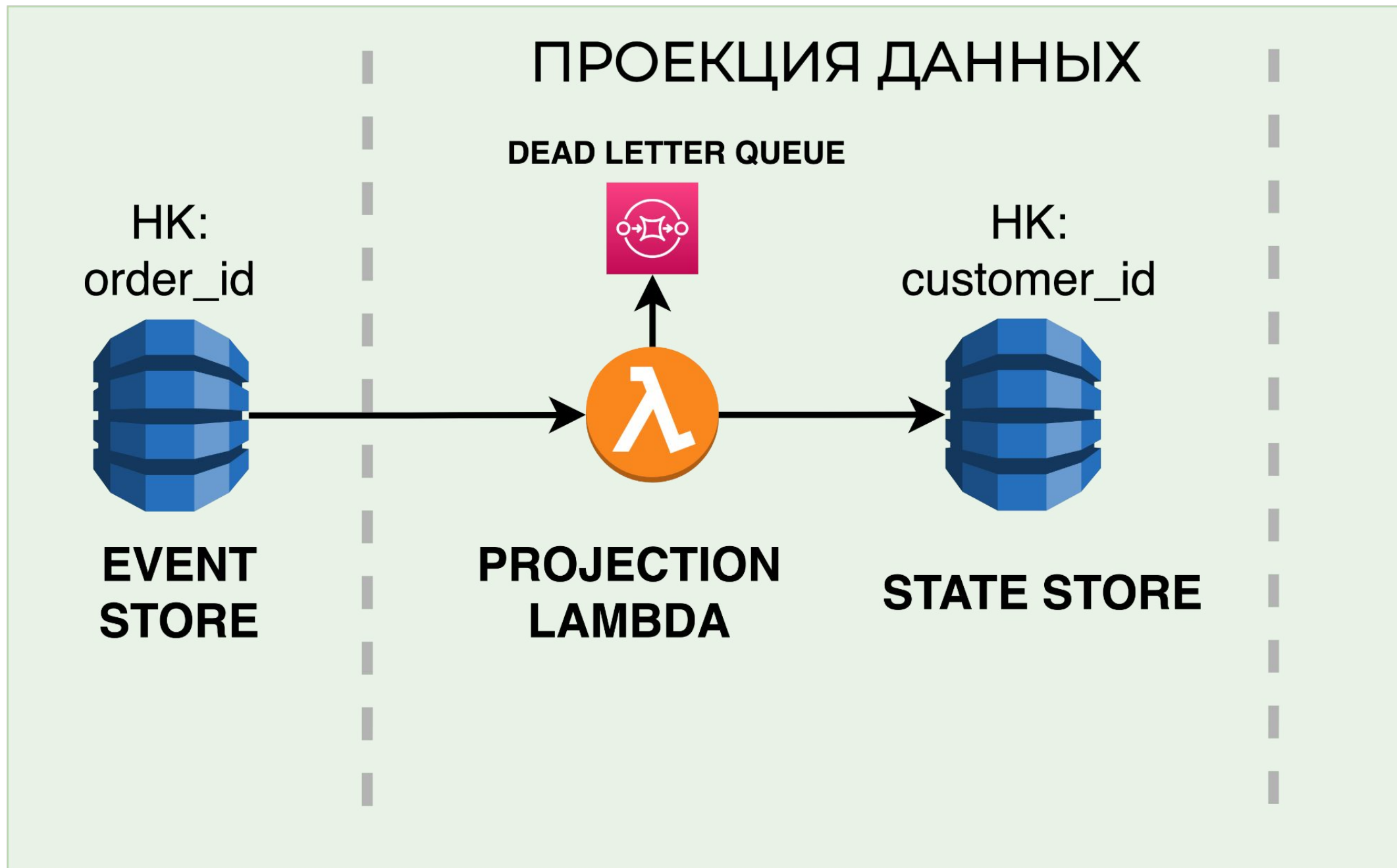
    return arguments.state

def handle_dispute_closed_event(arguments: Arguments) -> OrderState:
    dispute_id: str = arguments.event['dispute_id']
    dispute: Dispute = arguments.state.disputes.get(dispute_id)
    dispute.closed_at = int(arguments.event['closed_at'])
    arguments.state.disputes[dispute_id] = dispute

    return arguments.state
```



Синхронизация



Модель на чтение

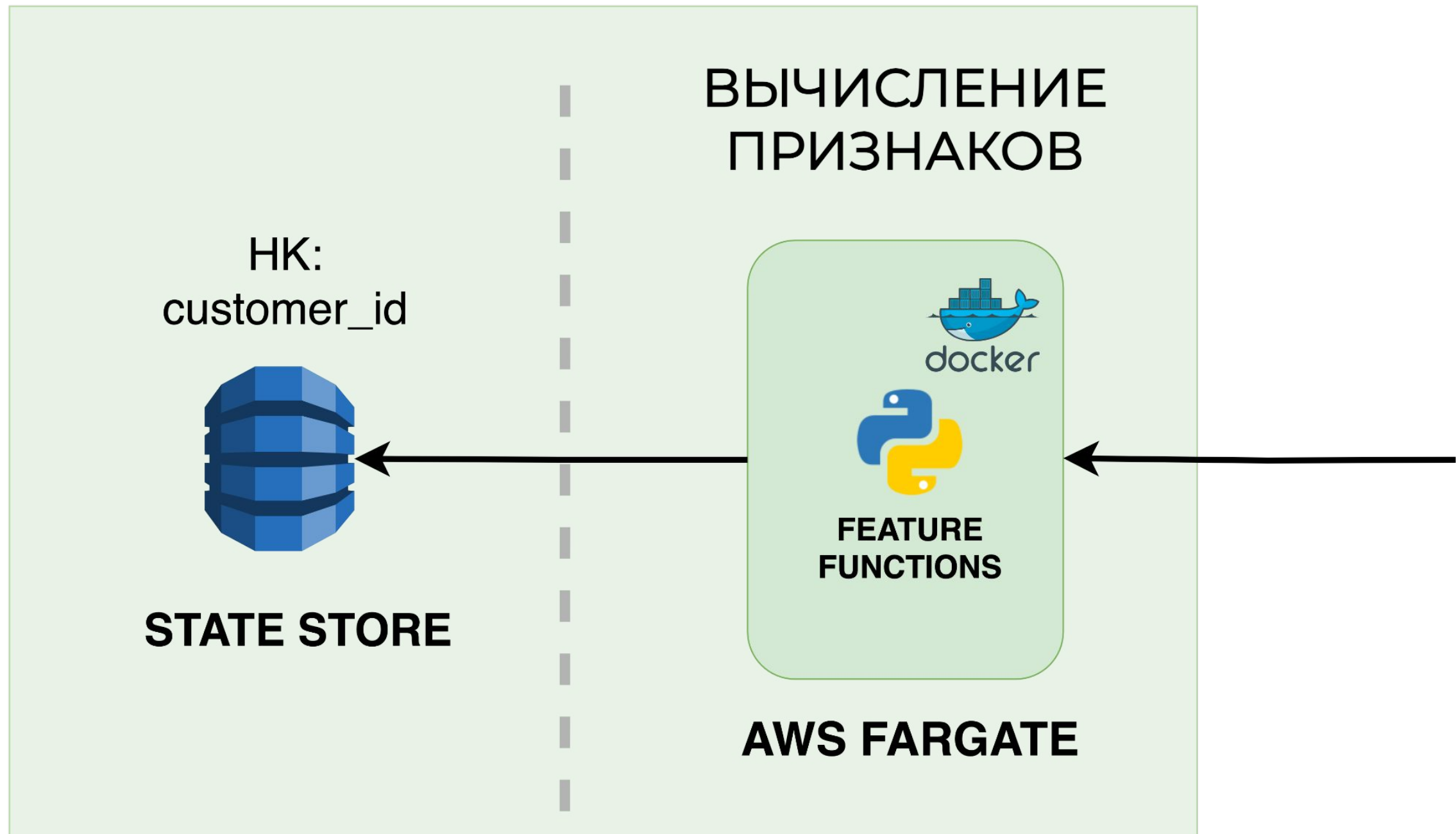


Вычисление признаков

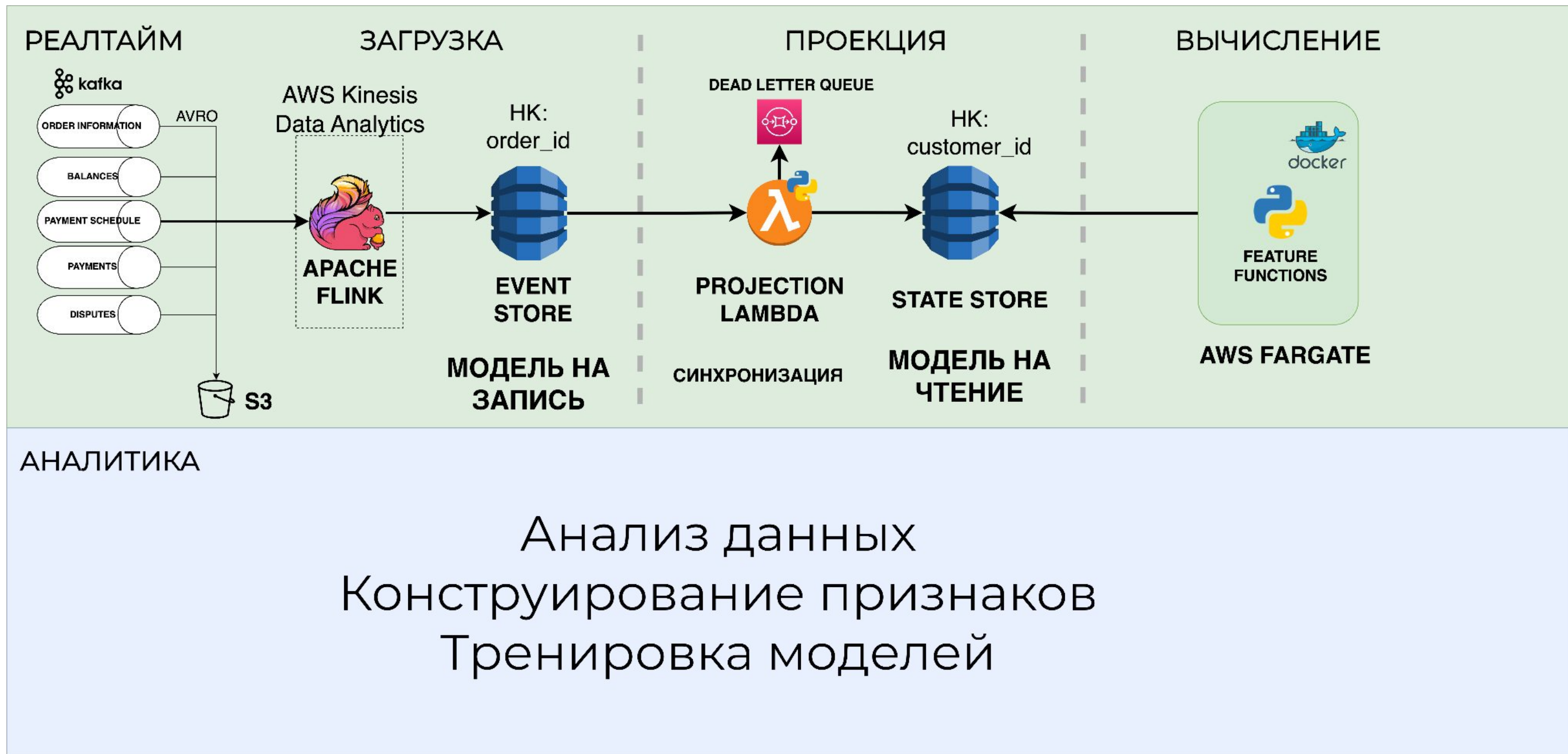
```
class CountOrdersPaid(IStatefulFeature):  
  
    @classmethod  
    def get_name(cls) -> str:  
        return "count_orders_paid"  
  
    @classmethod  
    def calculate(cls, projections: List[OrderState], decision_time: Optional[int] = None) -> int:  
        num_orders_paid: int = 0  
        for projection in projections:  
            if is_order_paid(projection):  
                num_orders_paid += 1  
  
        return num_orders_paid  
  
    @classmethod  
    def get_description(cls) -> str:  
        return """"  
            Count the number of orders that are fully paid  
        """"
```



Вычисление признаков



Общая архитектура

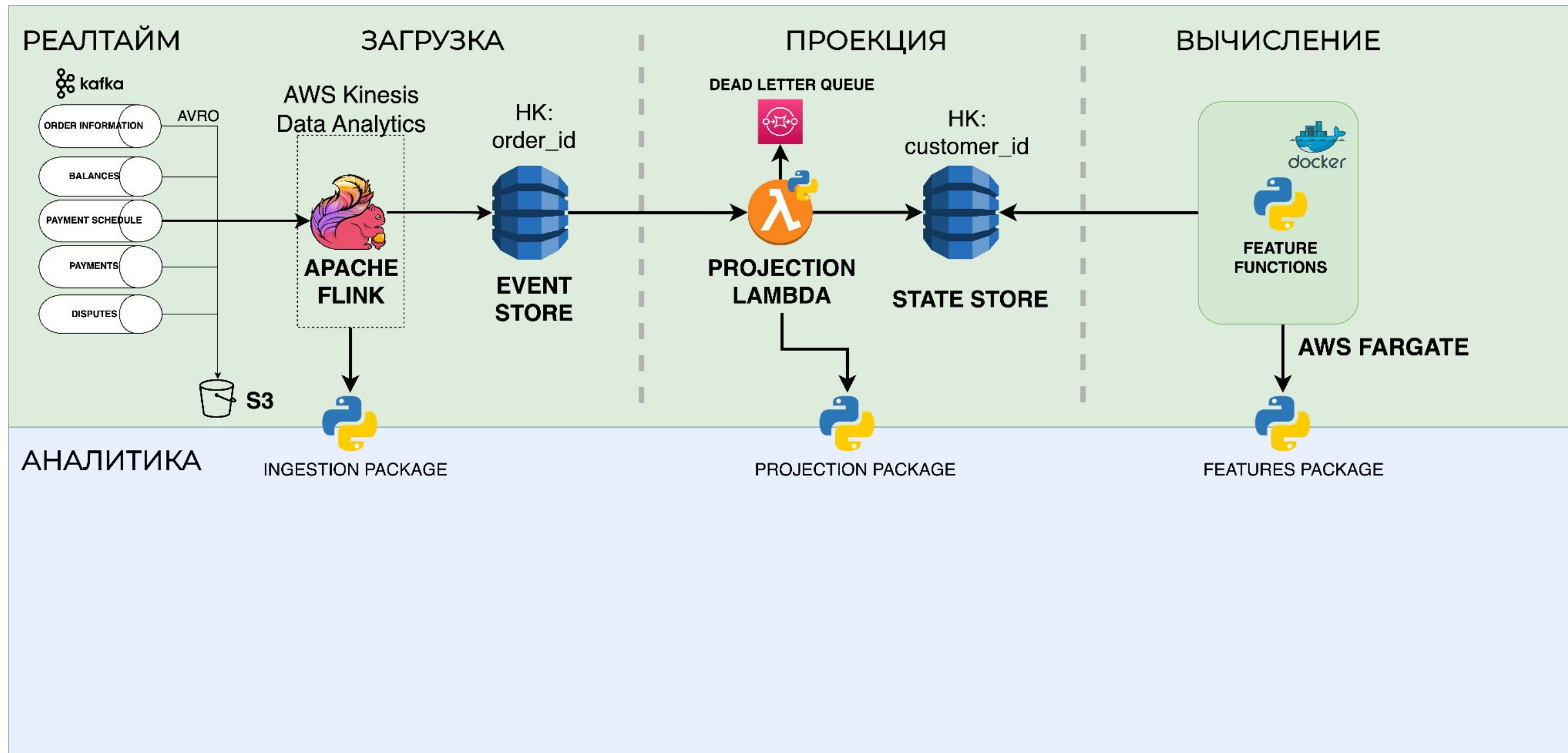


Слой аналитики

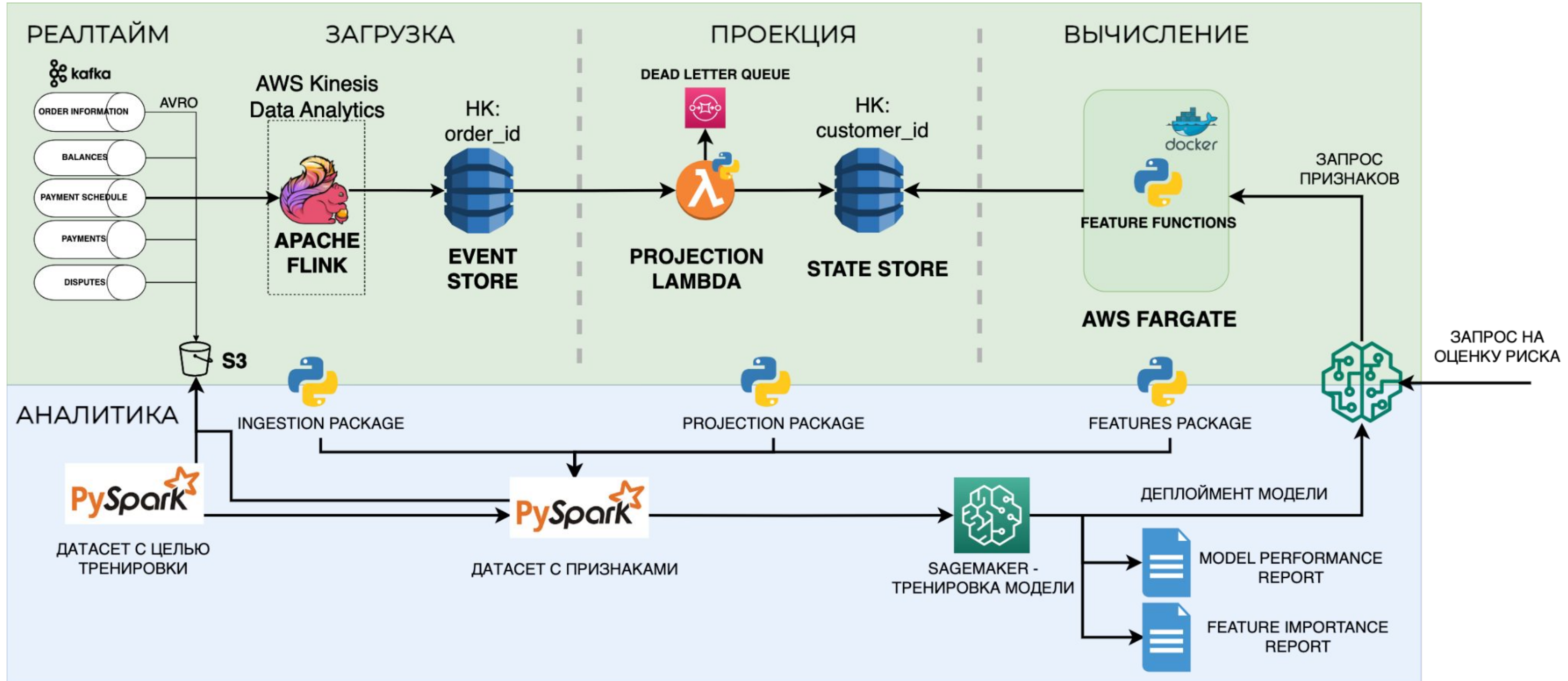
- Генерируем датасет для обучения модели. Были ли оплачены покупки?
- Вычисляем признаки для временных меток, в которые были сделаны покупки
- Запускаем обучение модели
- Оцениваем значимость признаков для модели



Слой аналитики



Слой аналитики



Так в чем секрет?

- Создаем и проверяем значимость признаков в офлайн
- Когда нас все устраивает, добавляем в реалтайм
- Новые признаки создавать просто
- Чаще всего не нужно пересчитывать данные всех клиентов
- Но при необходимости для этого можно использовать слой аналитики

Важность признака для модели (SHAP Value)

