

# EXACTLY-ONCE В МИКРОСЕРВИСНОЙ СРЕДЕ

# О СЕБЕ

- 10 лет опыта на .NET
- Работаю в компании Mindbox
- Занимаюсь высоко-нагруженными сервисами



# ЧТО БУДЕТ В ДОКЛАДЕ

1. Чуть-чуть бизнес
2. Идемпотентность
3. Transactional outbox
4. Составление ключа идемпотентности на кейсах
5. Эволюция

# ПОЧЕМУ ТАКАЯ ТЕМА



**Mathias Verraes**

@mathiasverraes

 Follow

There are only two hard problems in distributed systems: 2. Exactly-once delivery 1. Guaranteed order of messages 2. Exactly-once delivery

RETWEETS 6,775 LIKES 4,727



10:40 AM - 14 Aug 2015

 69  6.8K  4.7K 

# ВСЁ-ТАКИ EXACTLY ONCE ИЛИ AT LEAST ONCE?

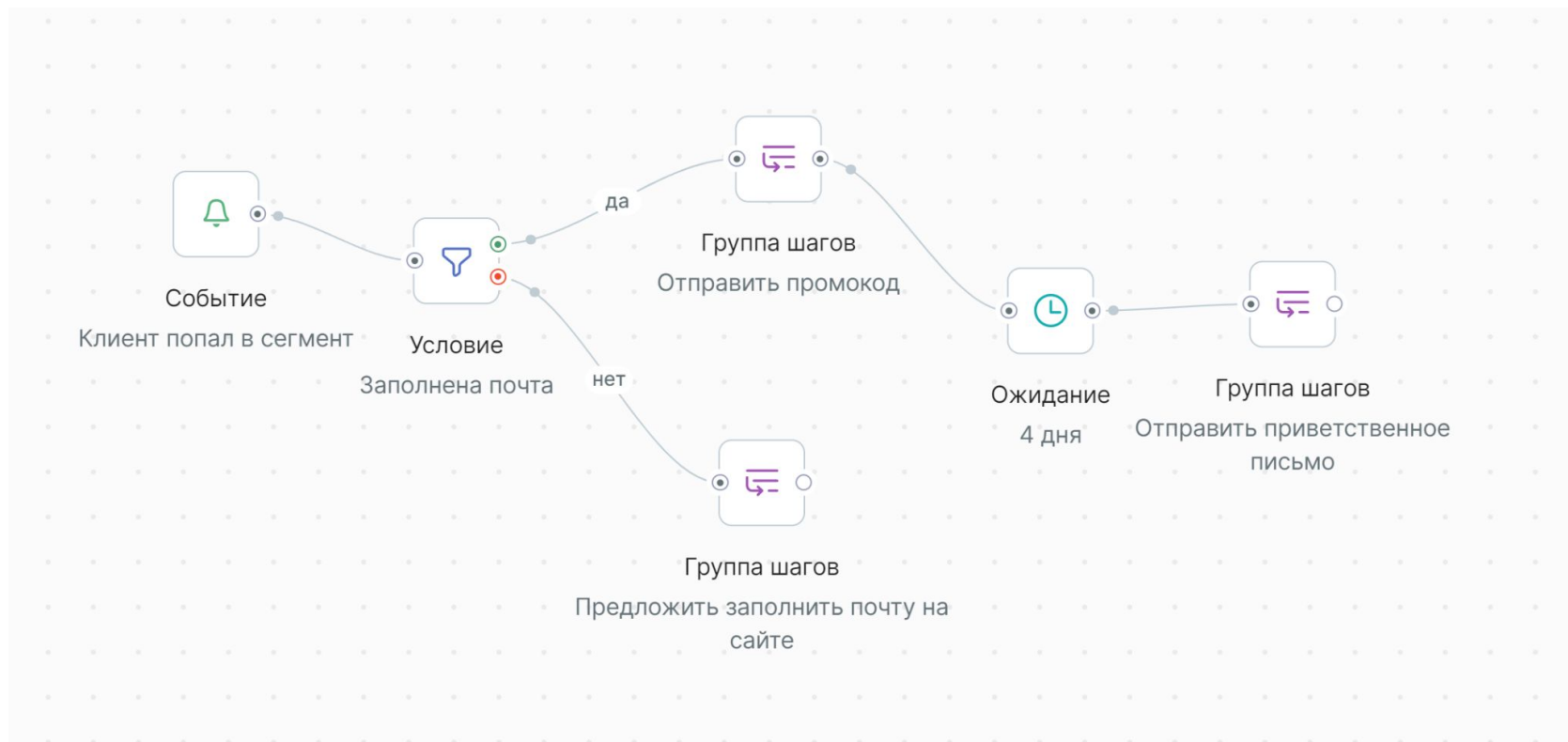
- не хотим распределенной транзакции
- идемпотентность решает проблемы

# КАКИЕ ВАШИ ДОКАЗАТЕЛЬСТВА?

- Кейсы
- Кейсы
- Кейсы



# БИЗНЕС



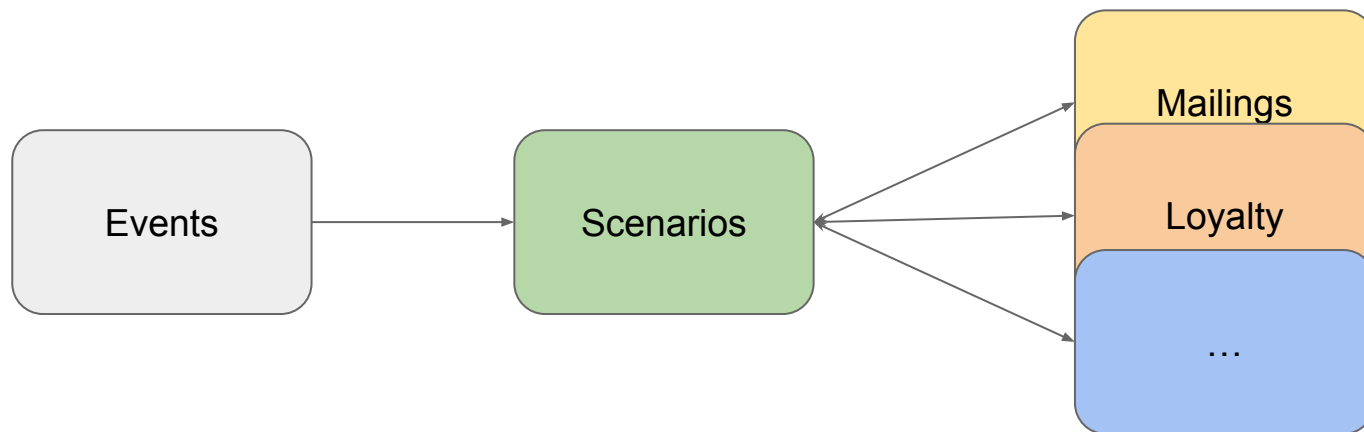
# МИКРОСЕРВИС

С монолитом будет тяжело

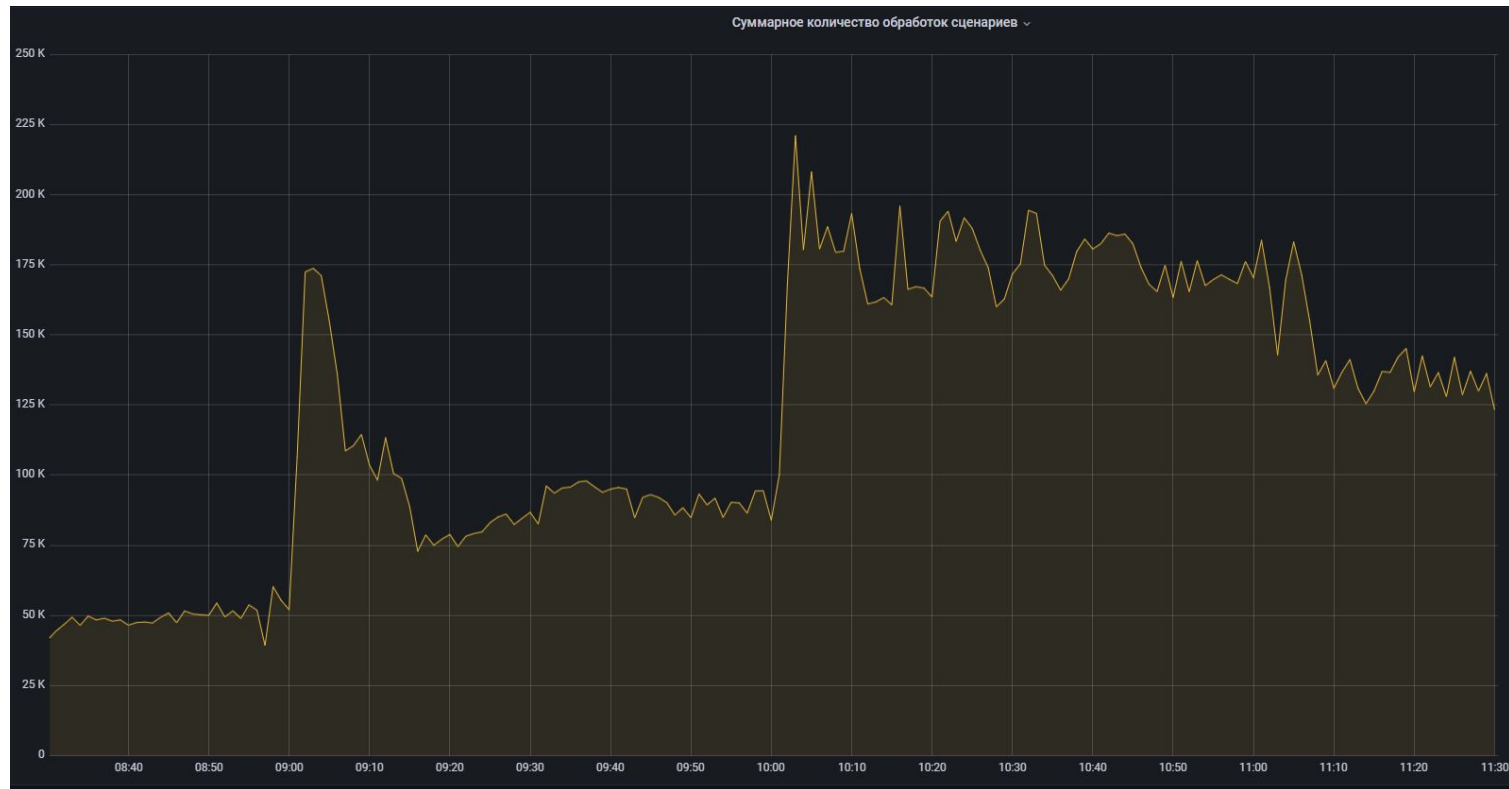
- Кросс-командная фича
- Много разработчиков
- Большая кодовая база



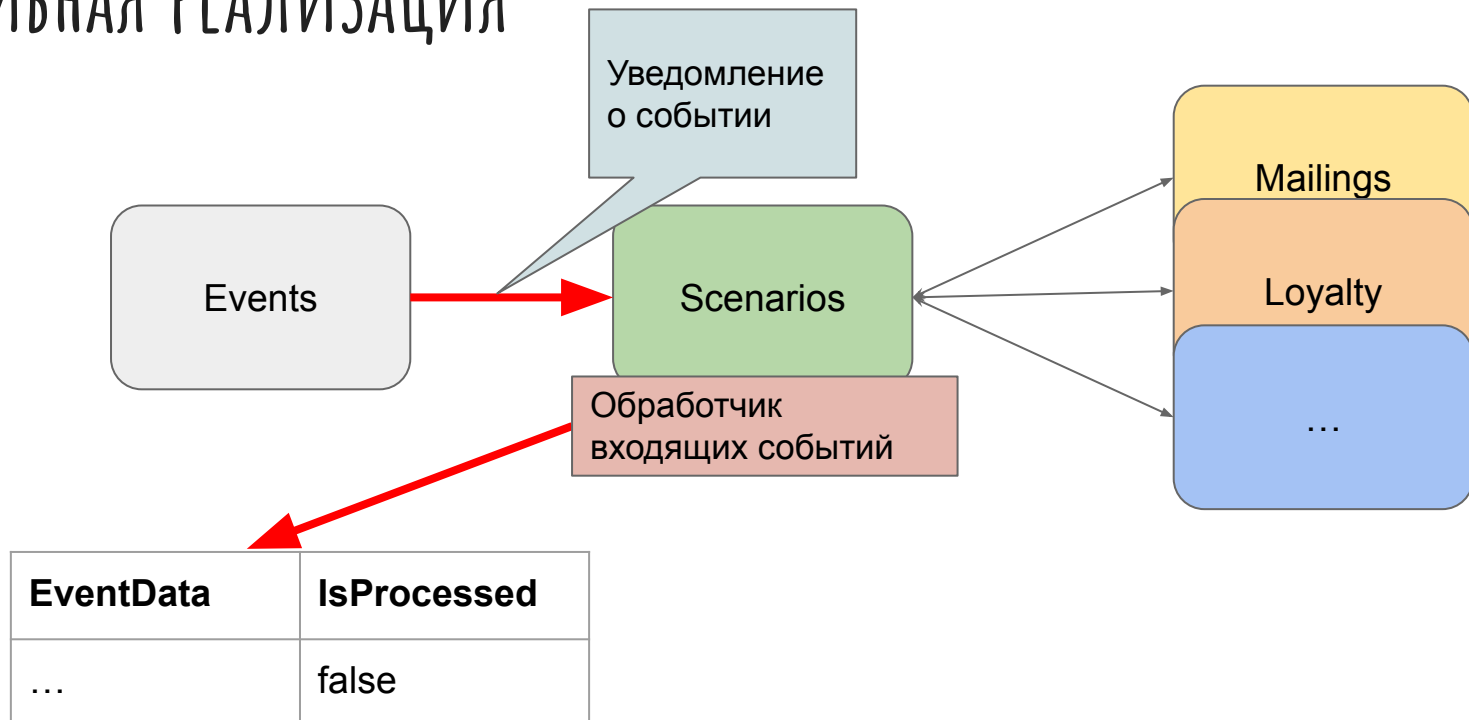
# КАКИЕ СЕРВИСЫ В ИТОГЕ ПОЛУЧИЛИСЬ



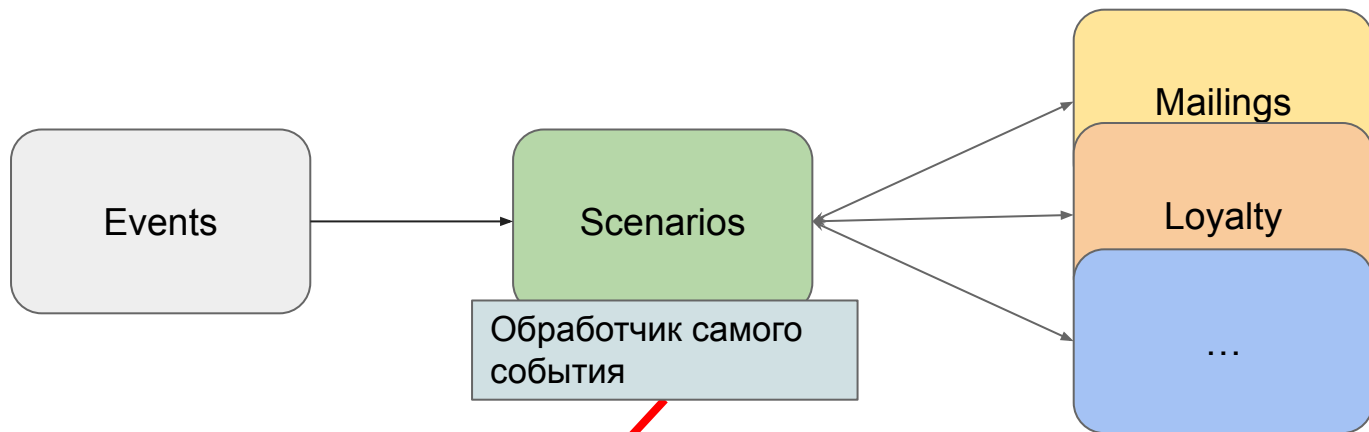
# ДОСТАТОЧНО ВЫСОКИЕ НАГРУЗКИ



# НАИВНАЯ РЕАЛИЗАЦИЯ

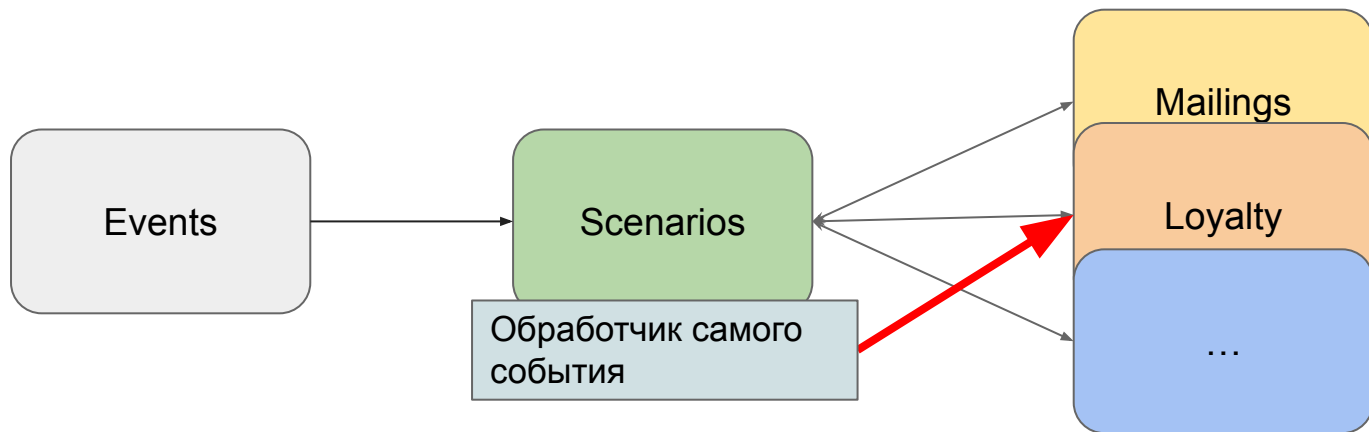


# НАИВНАЯ РЕАЛИЗАЦИЯ



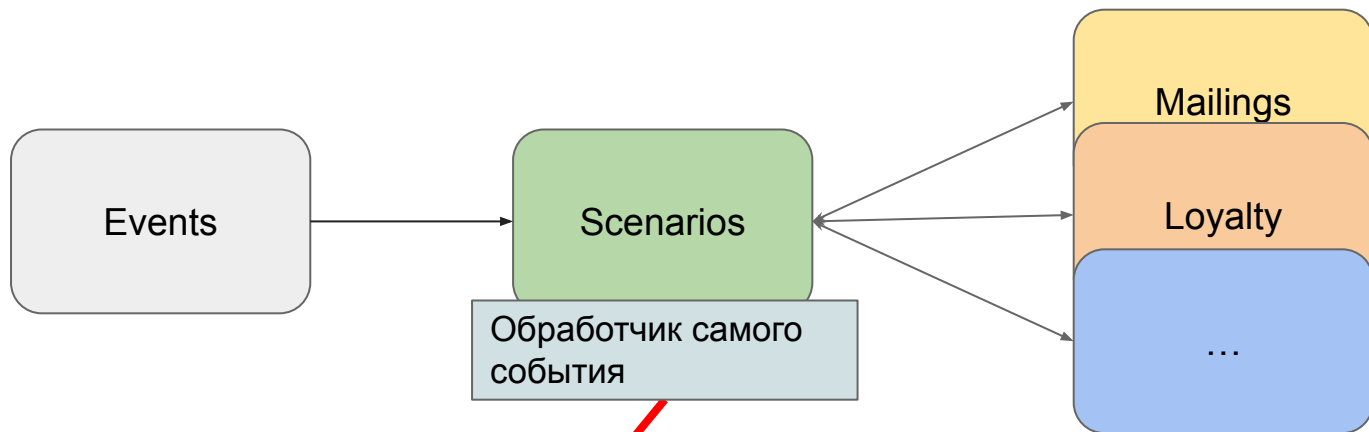
EventData	IsProcessed
...	<b>true</b>

# НАИВНАЯ РЕАЛИЗАЦИЯ



EventData	IsProcessed
...	<b>true</b>

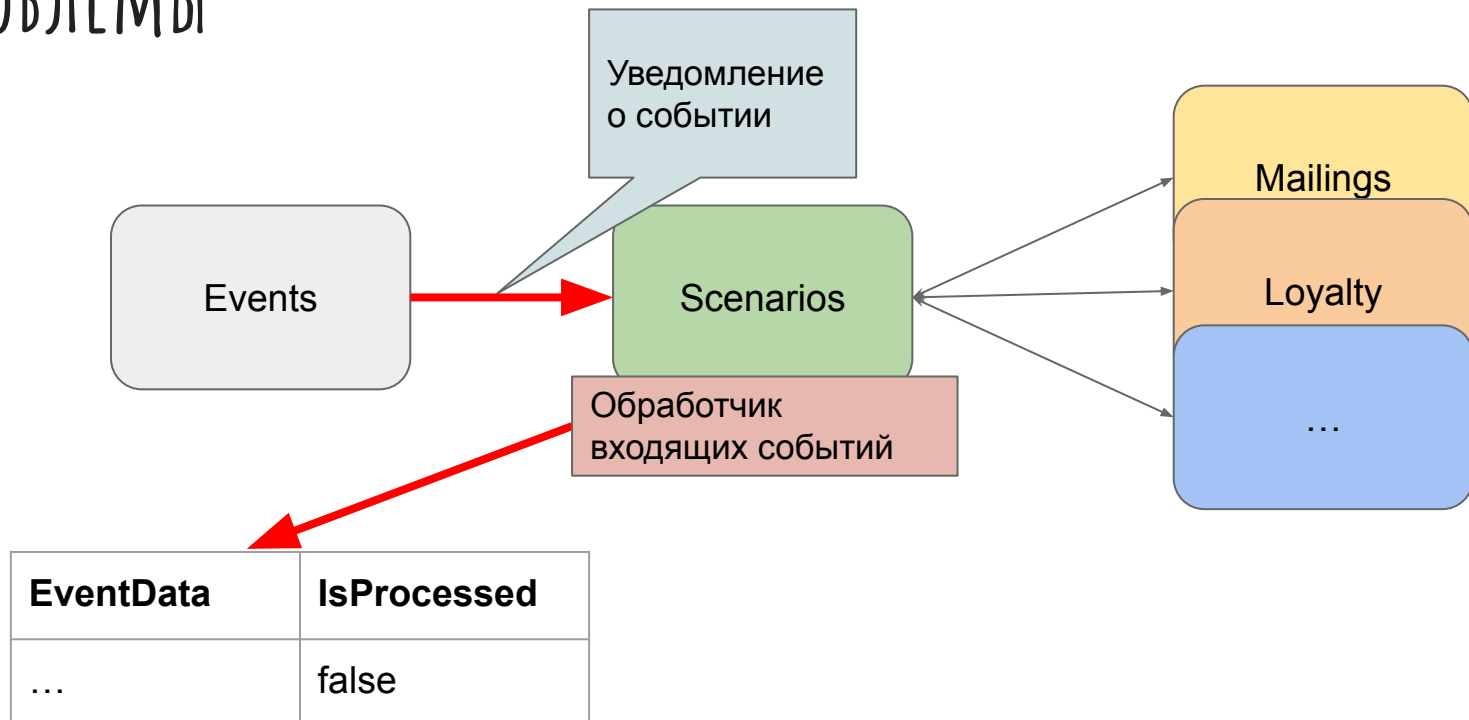
# НАИВНАЯ РЕАЛИЗАЦИЯ



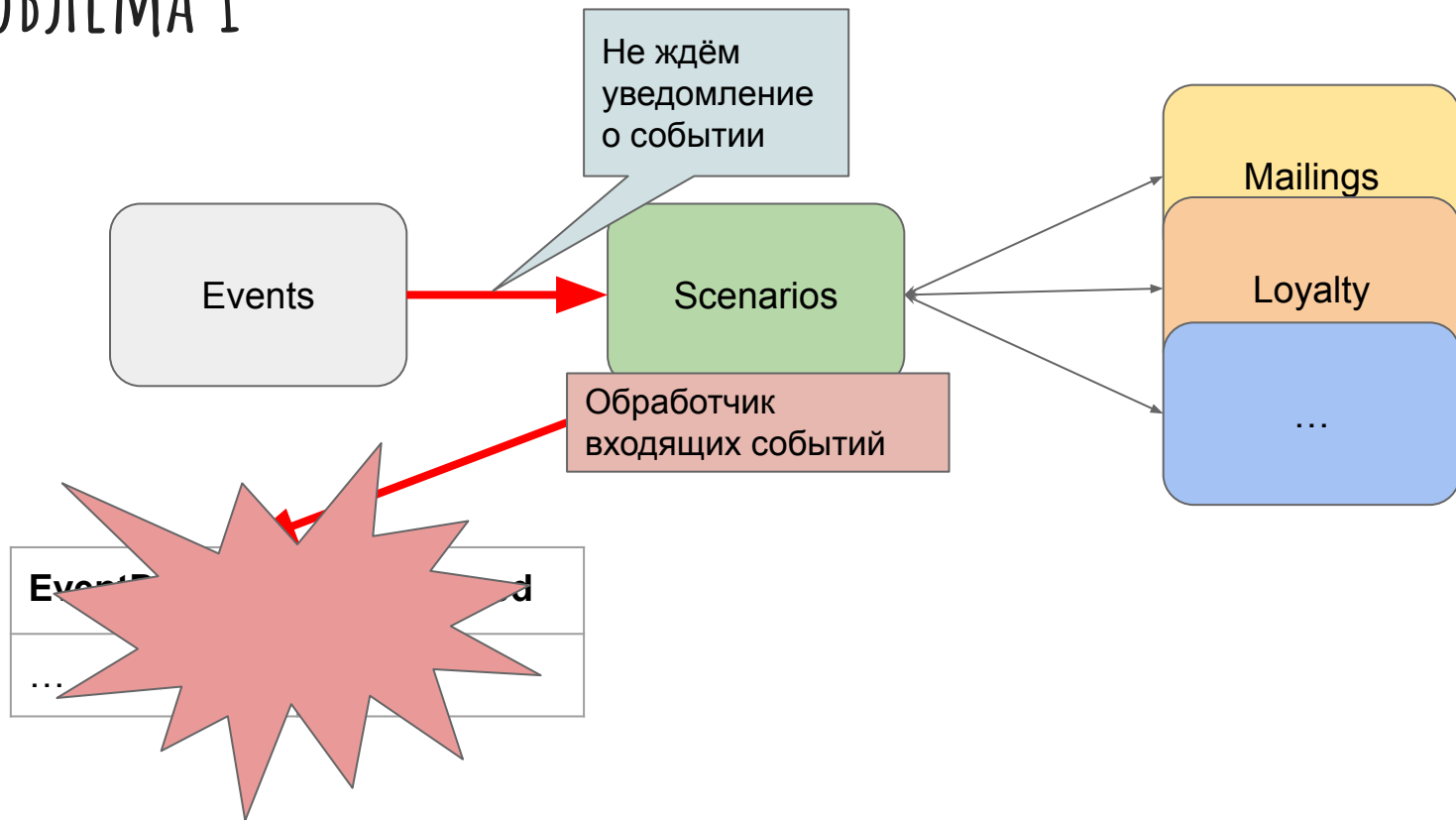
EventData	IsProcessed
...	true

COMMIT

# ПРОБЛЕМЫ

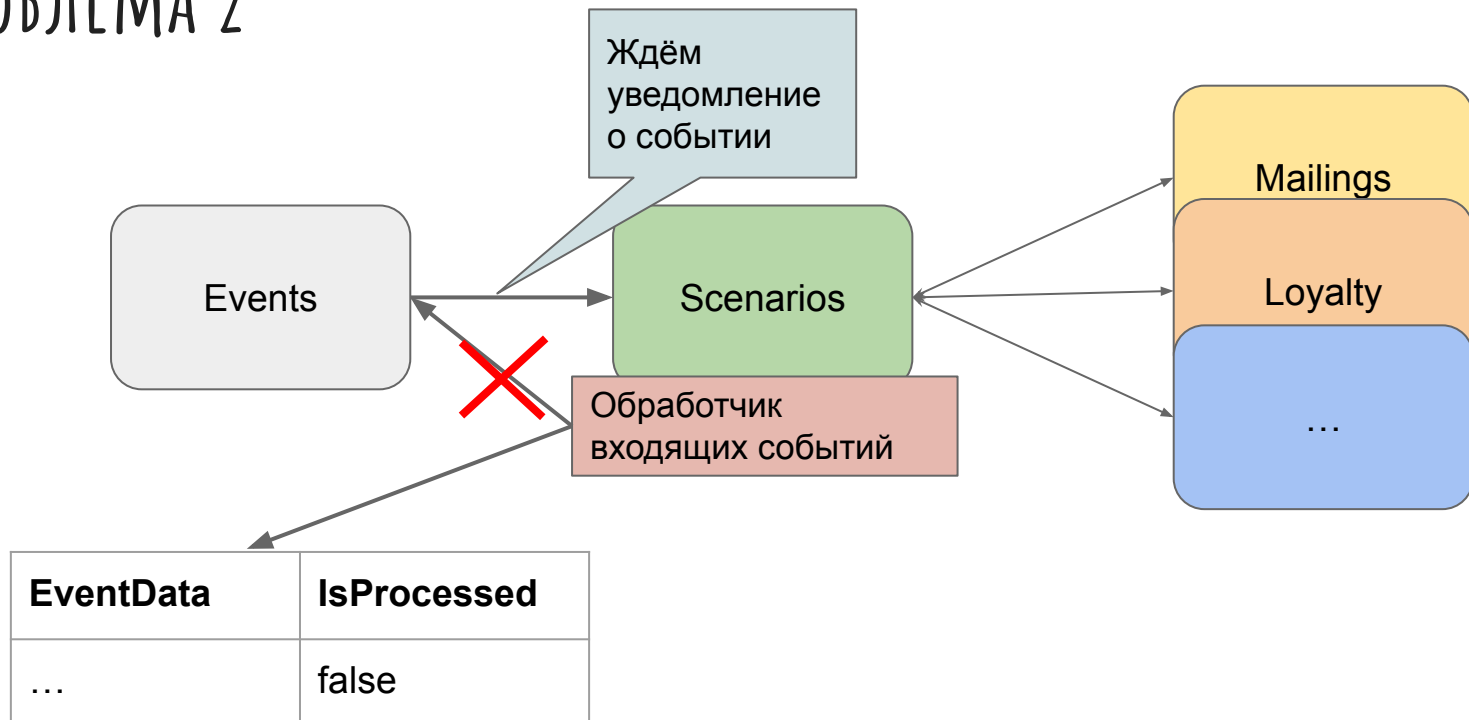


# ПРОБЛЕМА 1

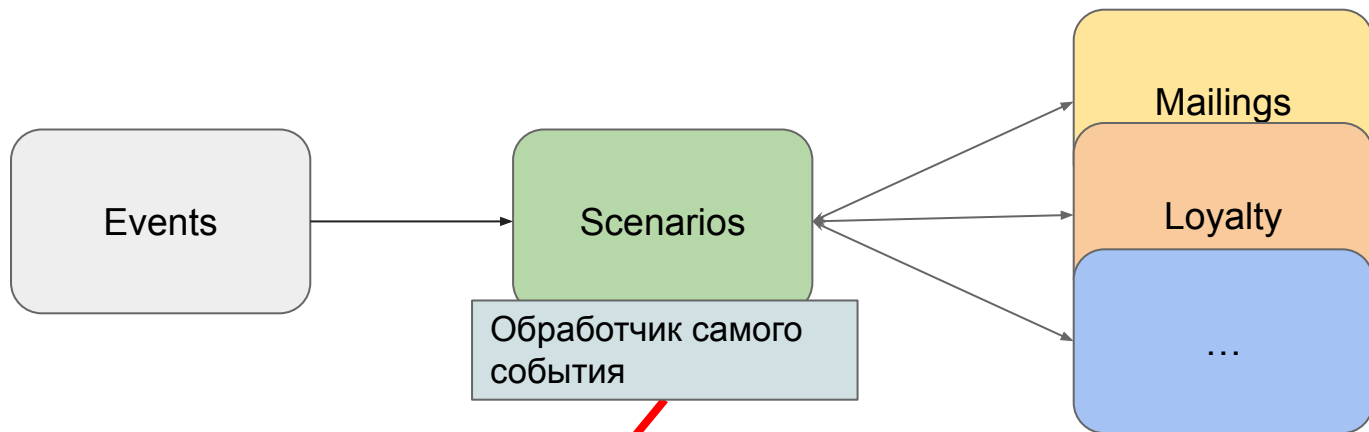




# ПРОБЛЕМА 2



# ПРОБЛЕМА 3



EventData	IsProcessed
...	true

COMMIT



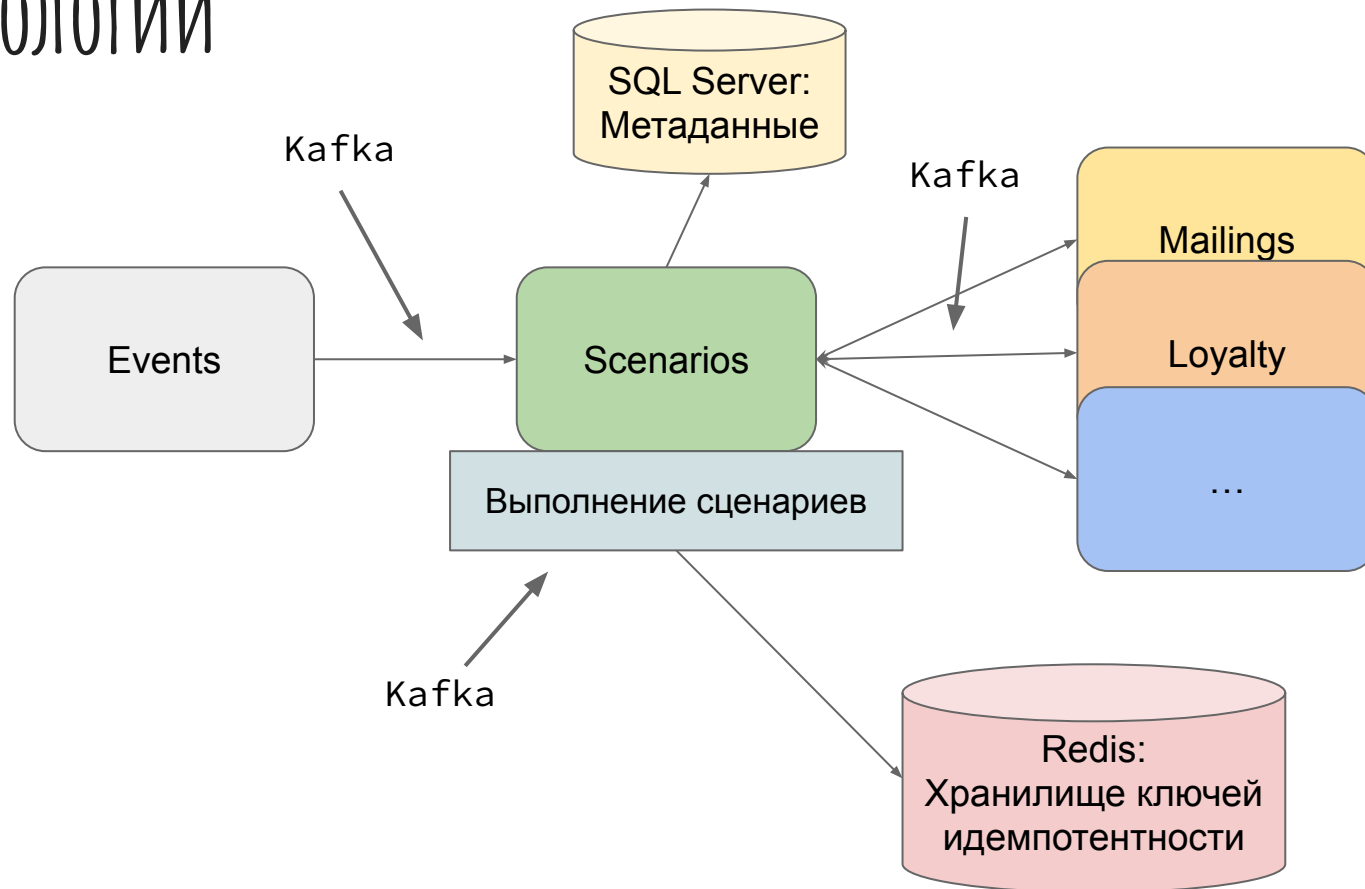
# ПРОБЛЕМА 4

Не соответствует требованиям по производительности, так как RDBMS не масштабируются горизонтально

# Выводы

- Даже в наивной реализации приходится думать о exactly-once гарантиях
- Если требуется высокая производительность и масштабируемость, нужно использовать соответствующие технологии

# ТЕХНОЛОГИИ



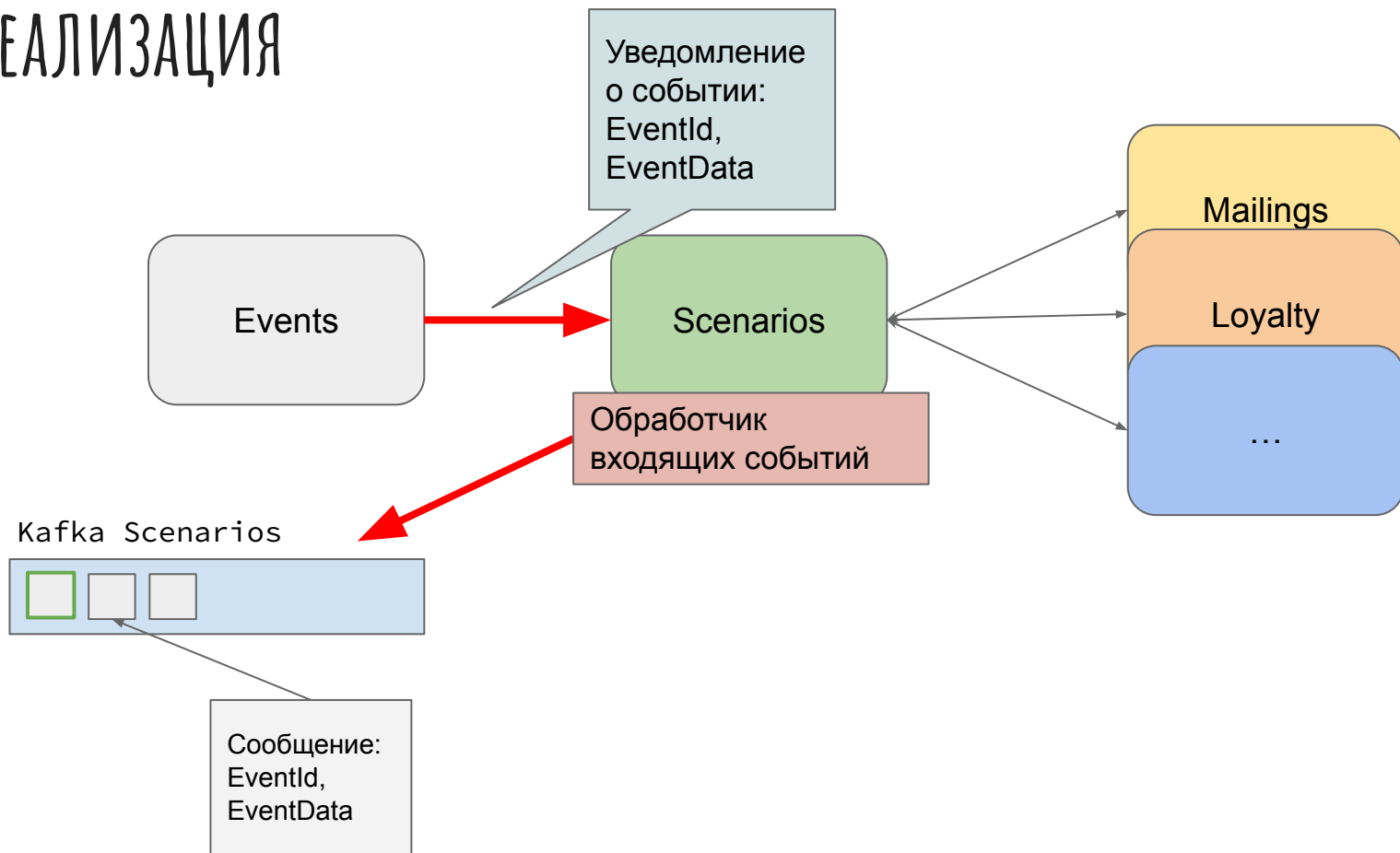
# ПОЧЕМУ КАФКА?

- Брокер сообщений, решает похожие с RabbitMQ задачи
- Persistent с репликацией: накопление миллионов необработанных сообщений в топике – не страшно
- Горизонтально масштабируется

# ЧТО ТАКОЕ ЭТА ВАША ИДЕМПОТЕНТНОСТЬ?

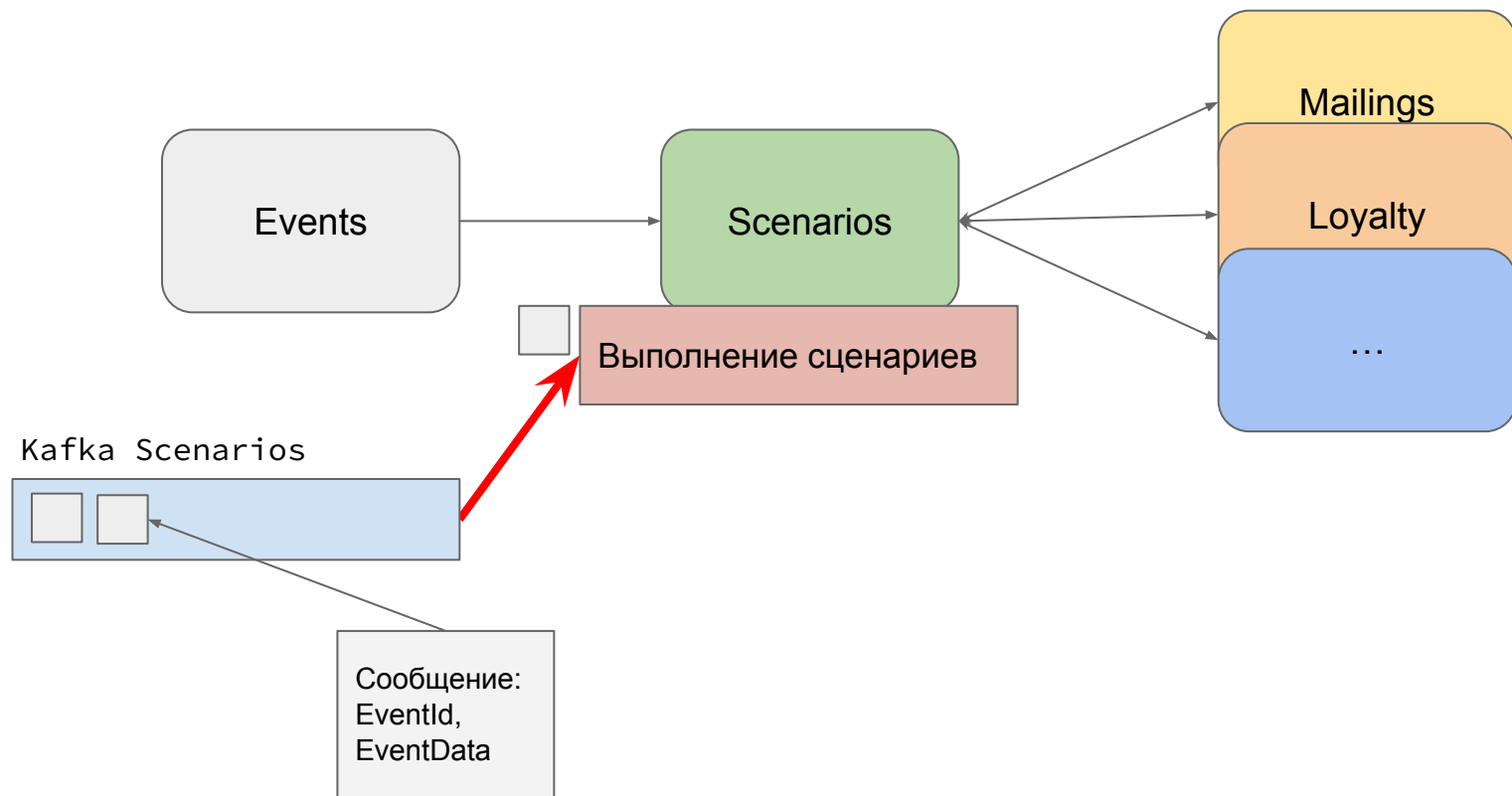
Свойство операции, при повторении которой не будет происходить новых сайд-эффектов

# РЕАЛИЗАЦИЯ

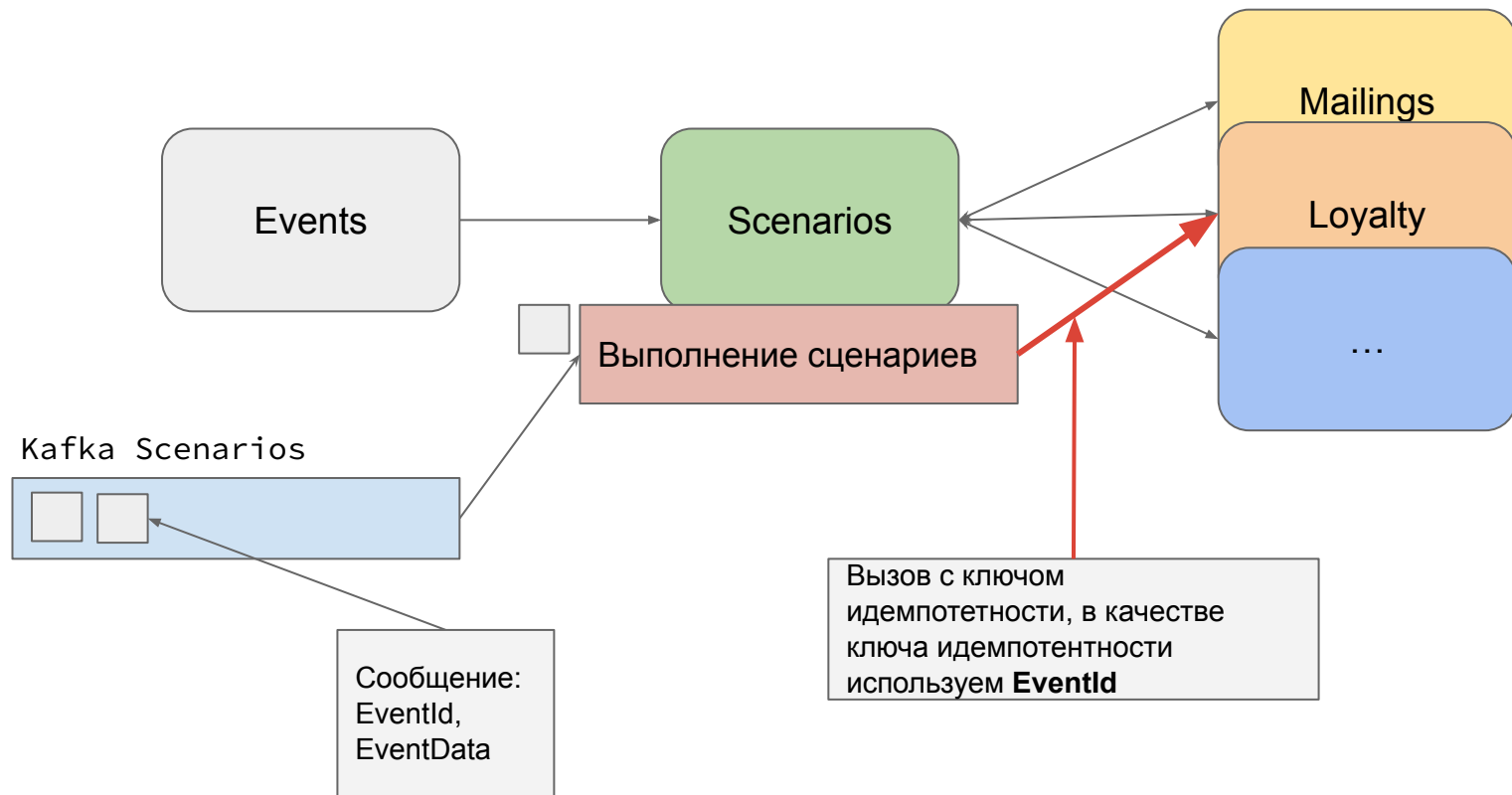




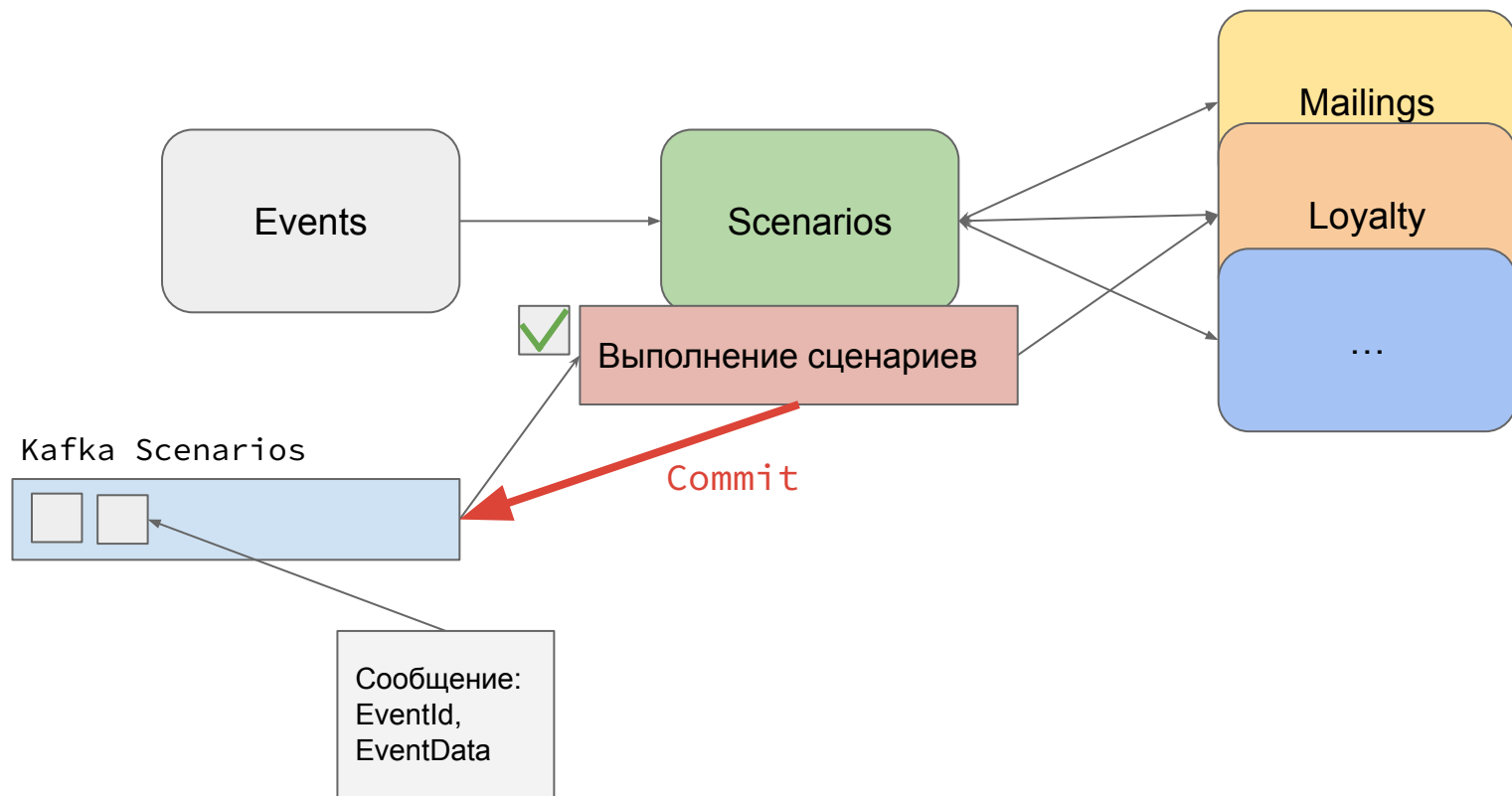
# РЕАЛИЗАЦИЯ



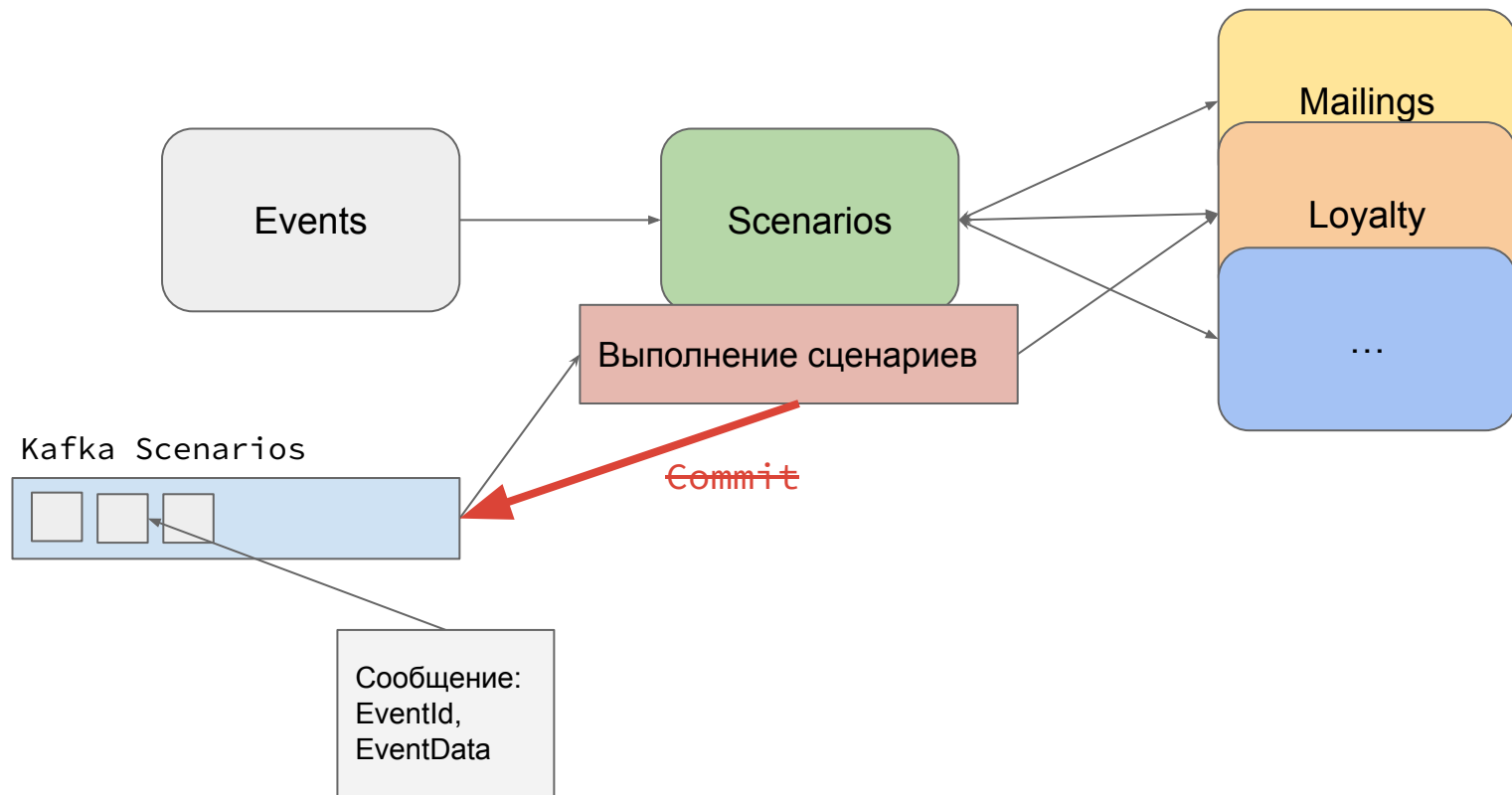
# РЕАЛИЗАЦИЯ



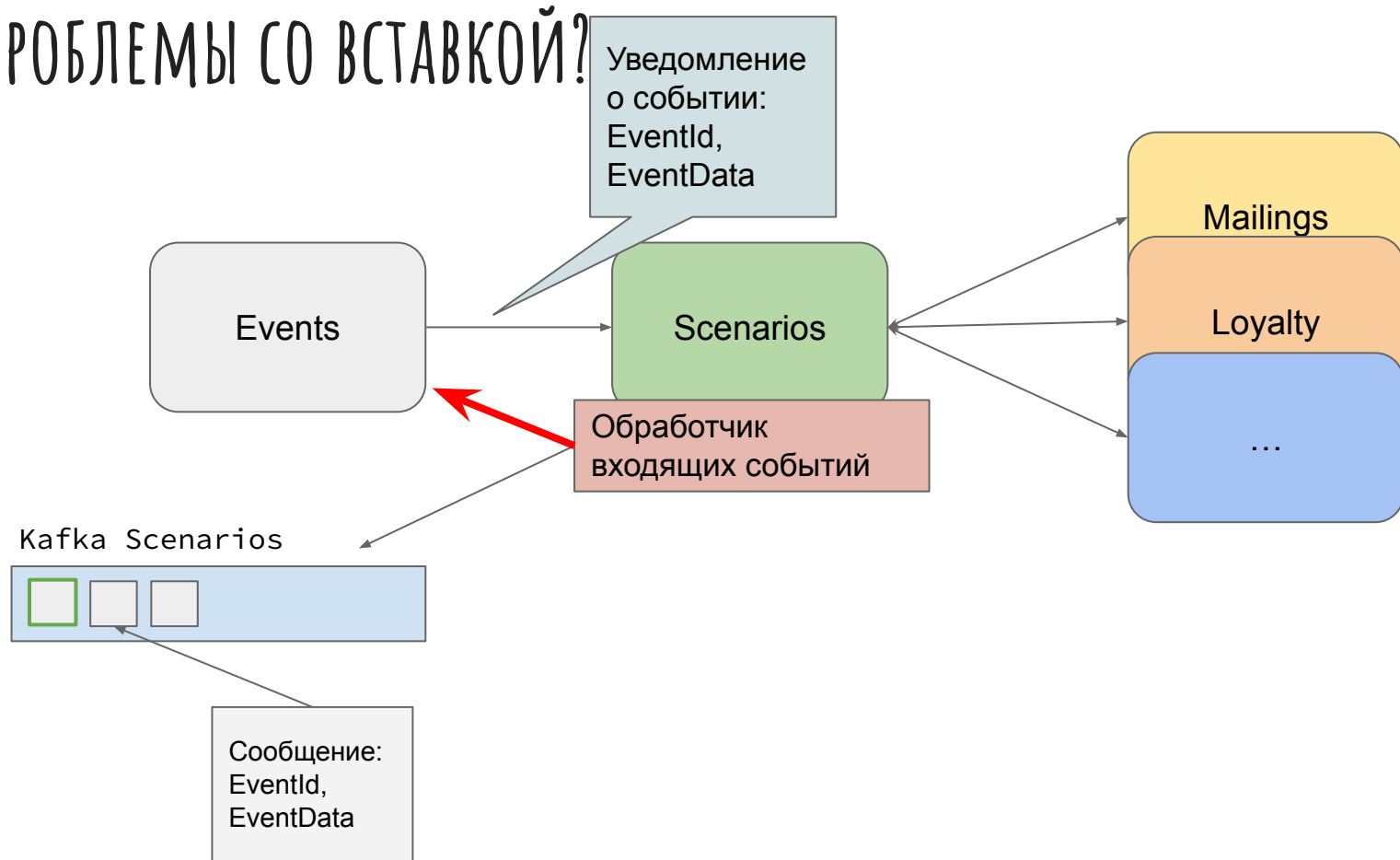
# РЕАЛИЗАЦИЯ



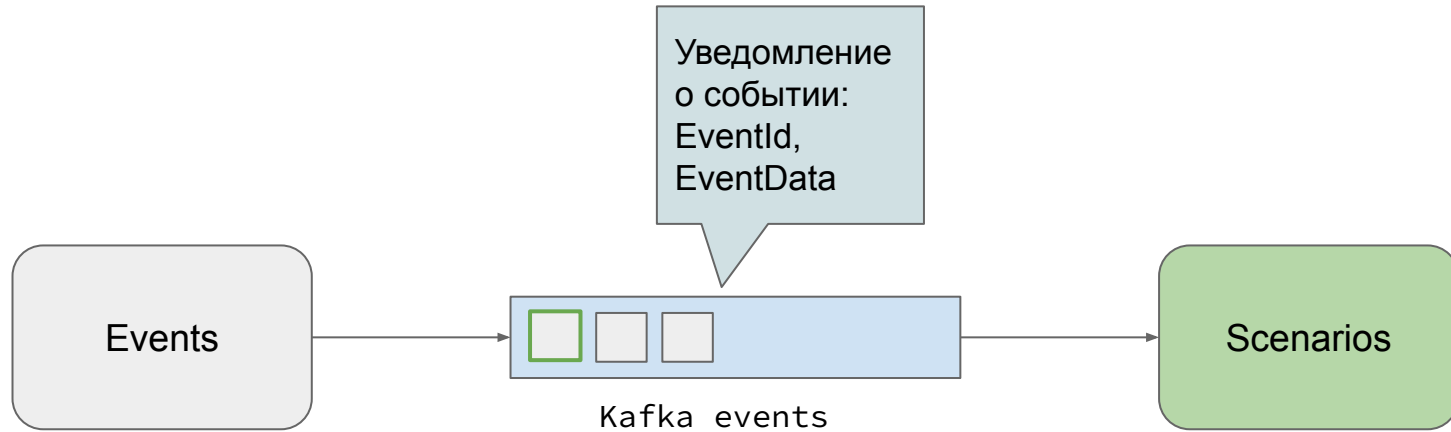
# ПРОБЛЕМЫ С КОММИТОМ?



# ПРОБЛЕМЫ СО ВСТАВКОЙ?

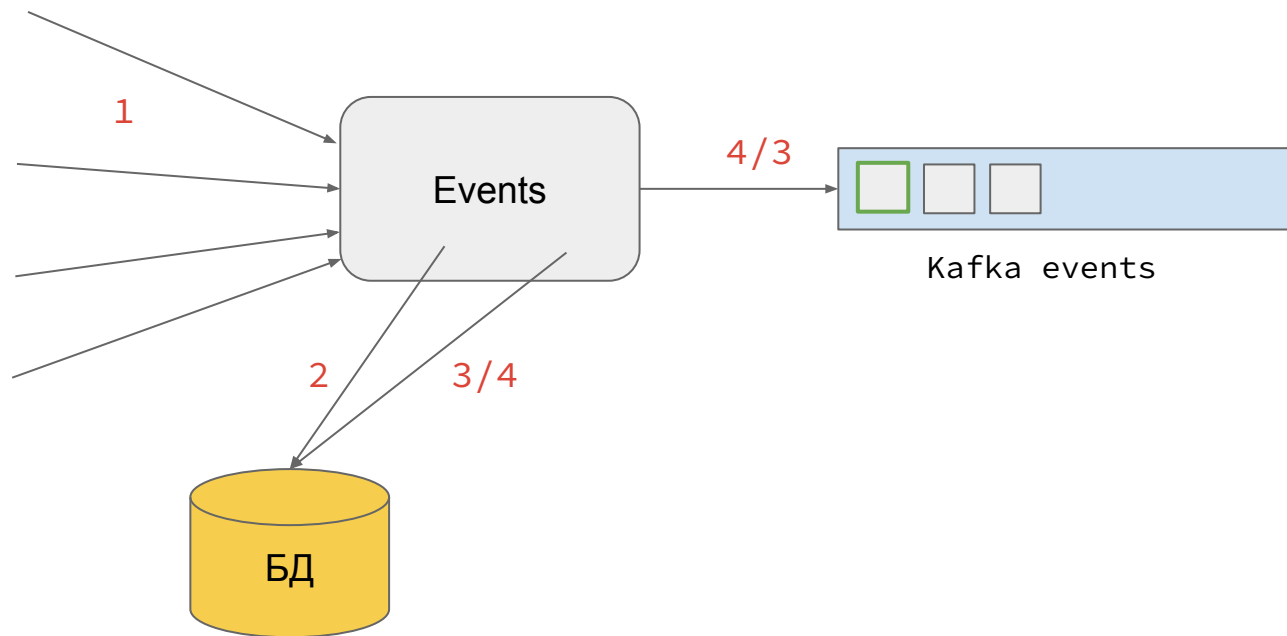


# ОТКУДА ЕЩЁ МОГУТ БЫТЬ ДУБЛИ СООБЩЕНИЙ?

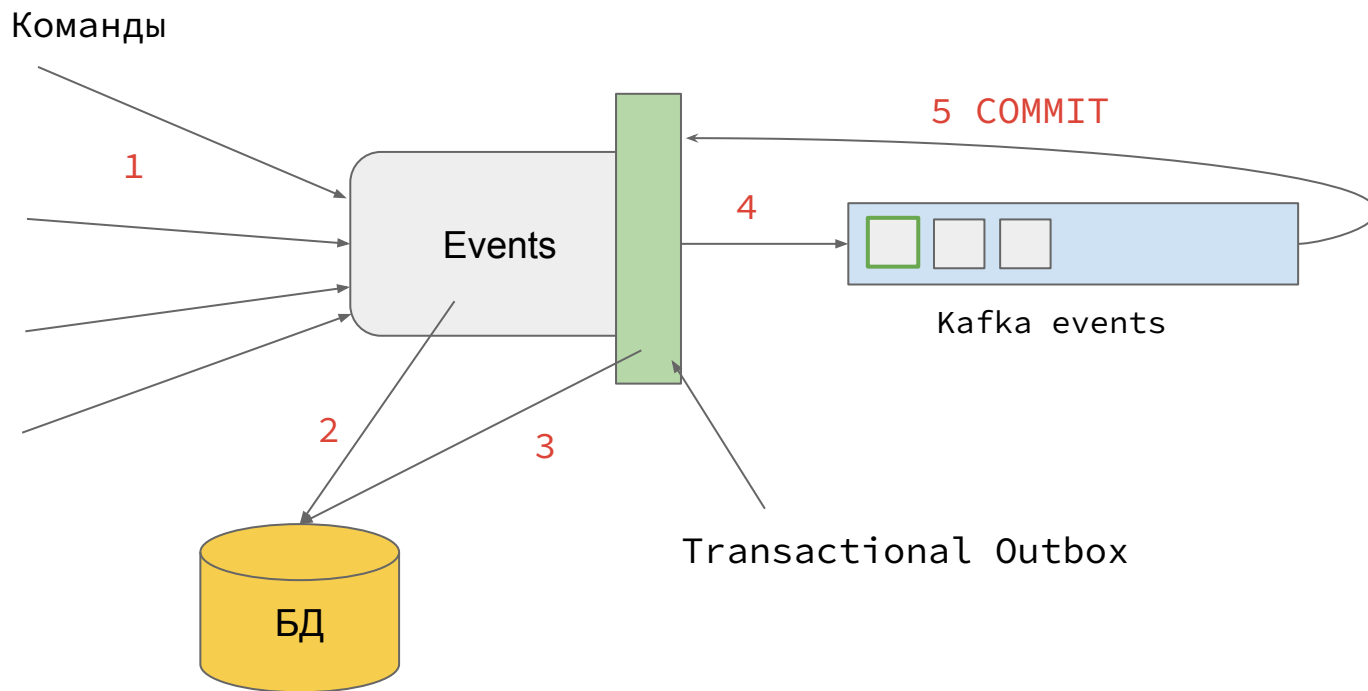


# КАК УВЕДОМЛЯЕМ О СОБЫТИЯХ?

Команды

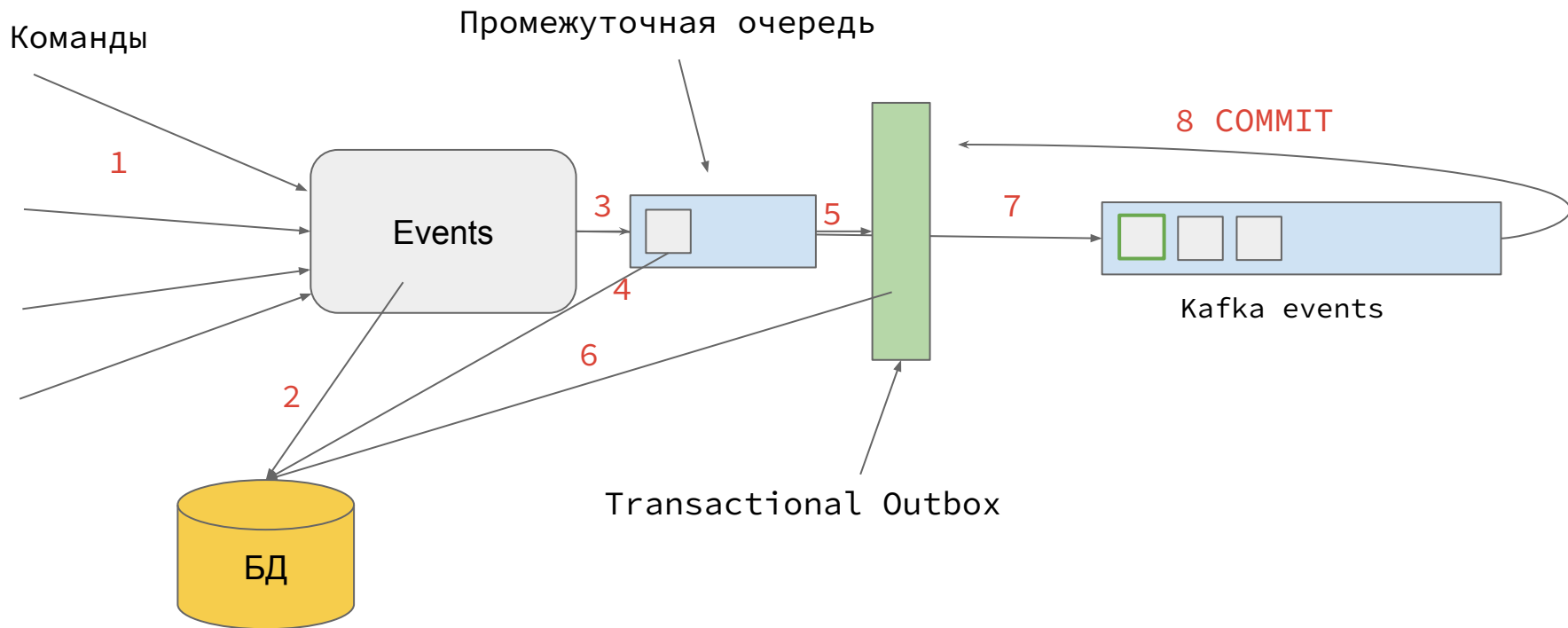


# КАК УВЕДОМЛЯЕМ О СОБЫТИЯХ?

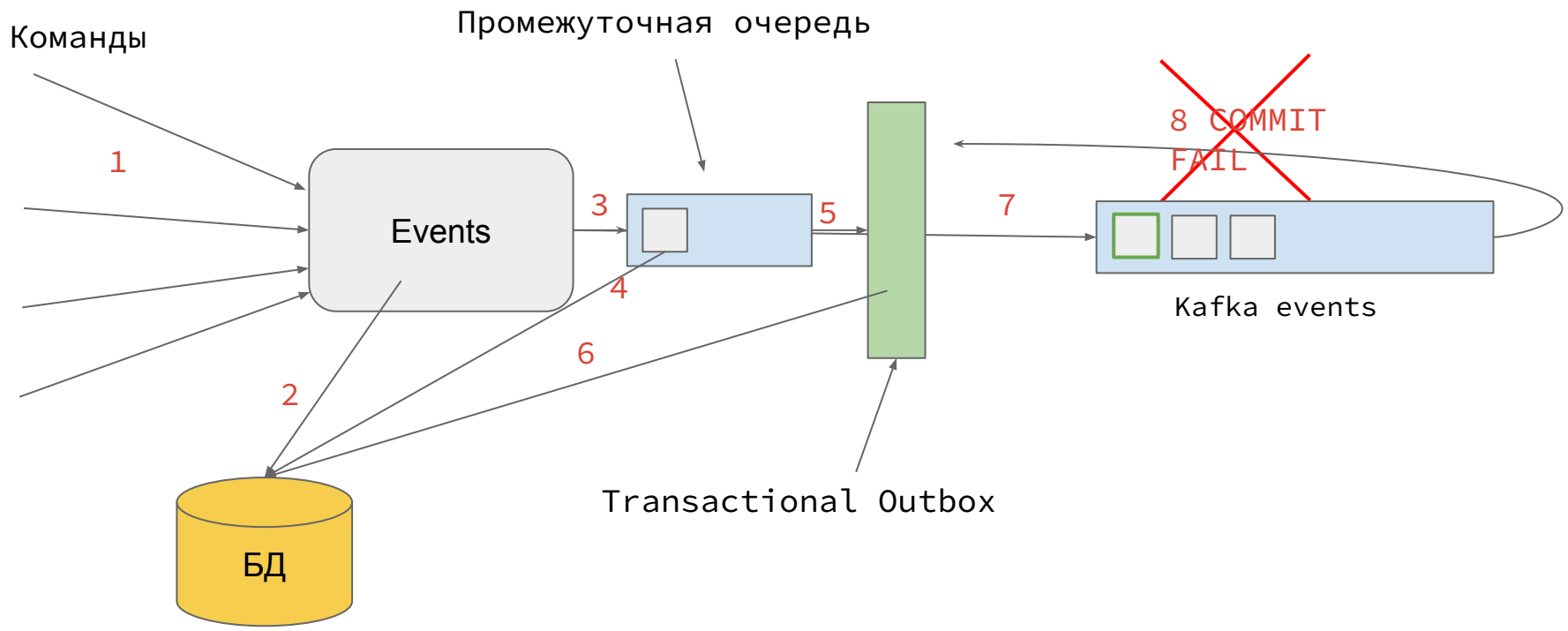




# КАК УВЕДОМЛЯЕМ О СОБЫТИЯХ?



# ГДЕ ЖЕ ТУТ ДУБЛИ?



# А КАК ЖЕ КАФКА IDEMPOTENT PRODUCER?

Он может только уменьшить количество дублей

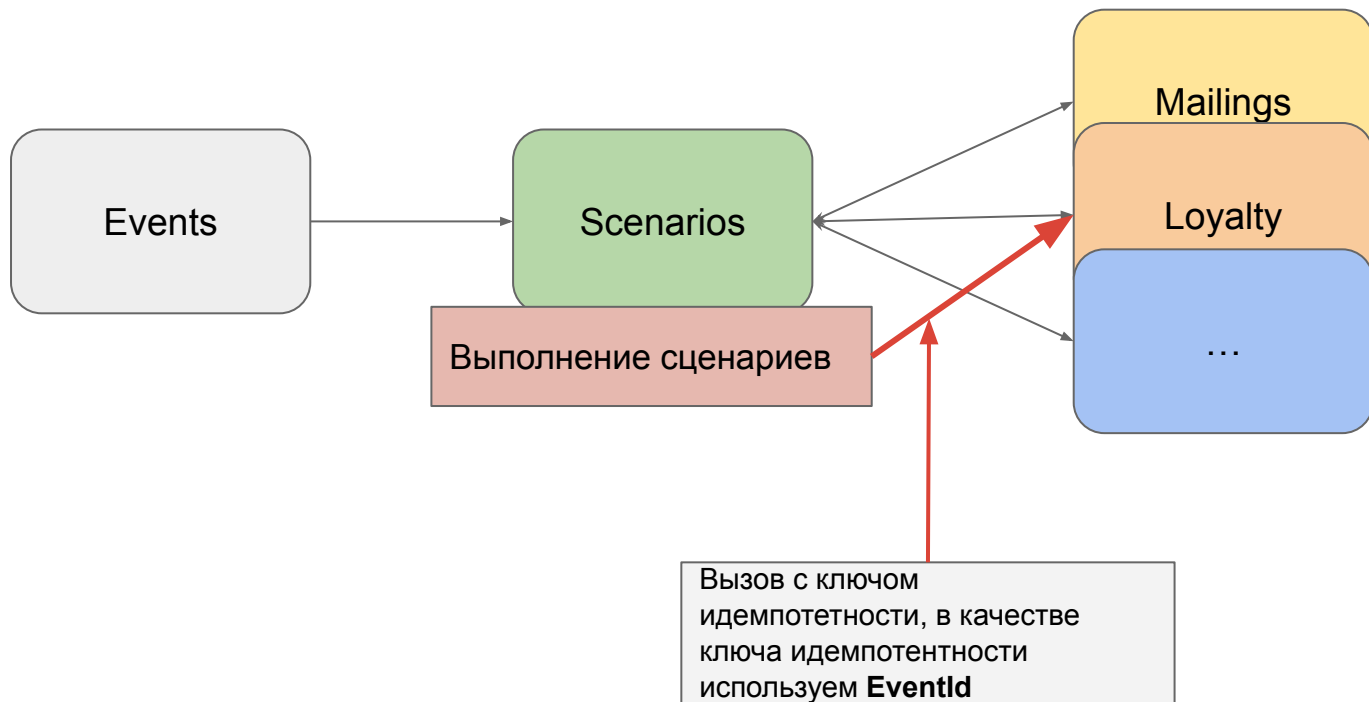
- Падения подов/ноды/сети
- Запрет на самостоятельные ретраи

When `enable.idempotence` is set to `true`, no manual retries are required, in fact performing retries in your application code will still cause duplicates.

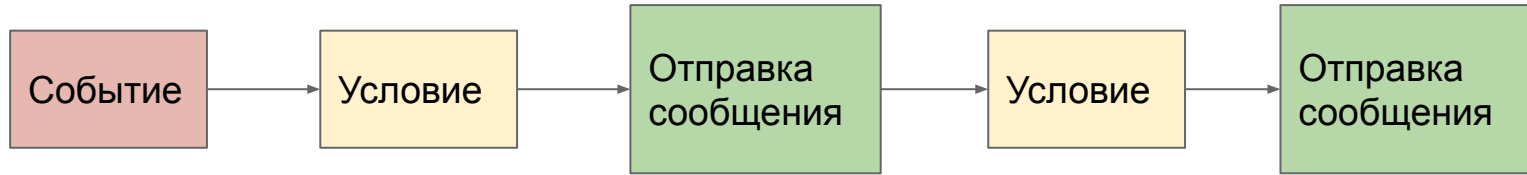
# ЕЩЁ ПАРА СЛОВ О TRANSACTIONAL OUTBOX

- Дорогая штука
- При больших нагрузках и идемпотентной переобработке мы не используем
- Вместо него можно попробовать [transaction log tailing](#)

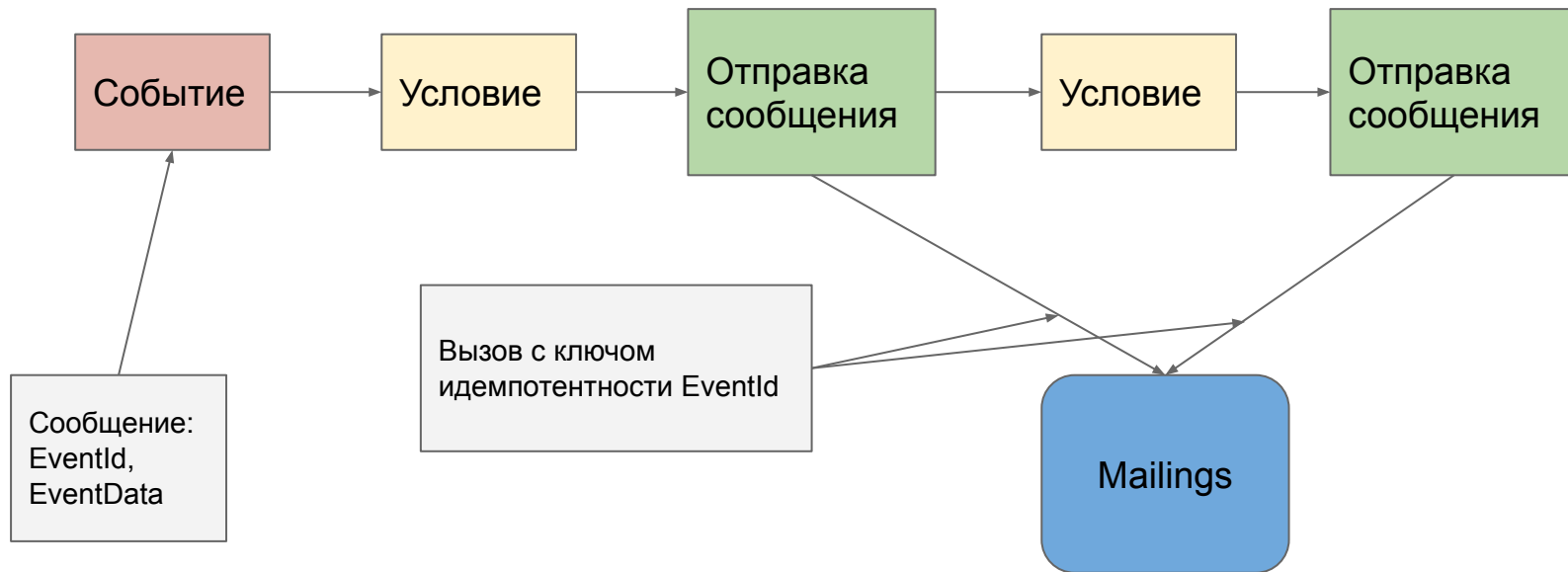
# ВЕРНЁМСЯ К САМИМ СЦЕНАРИЯМ



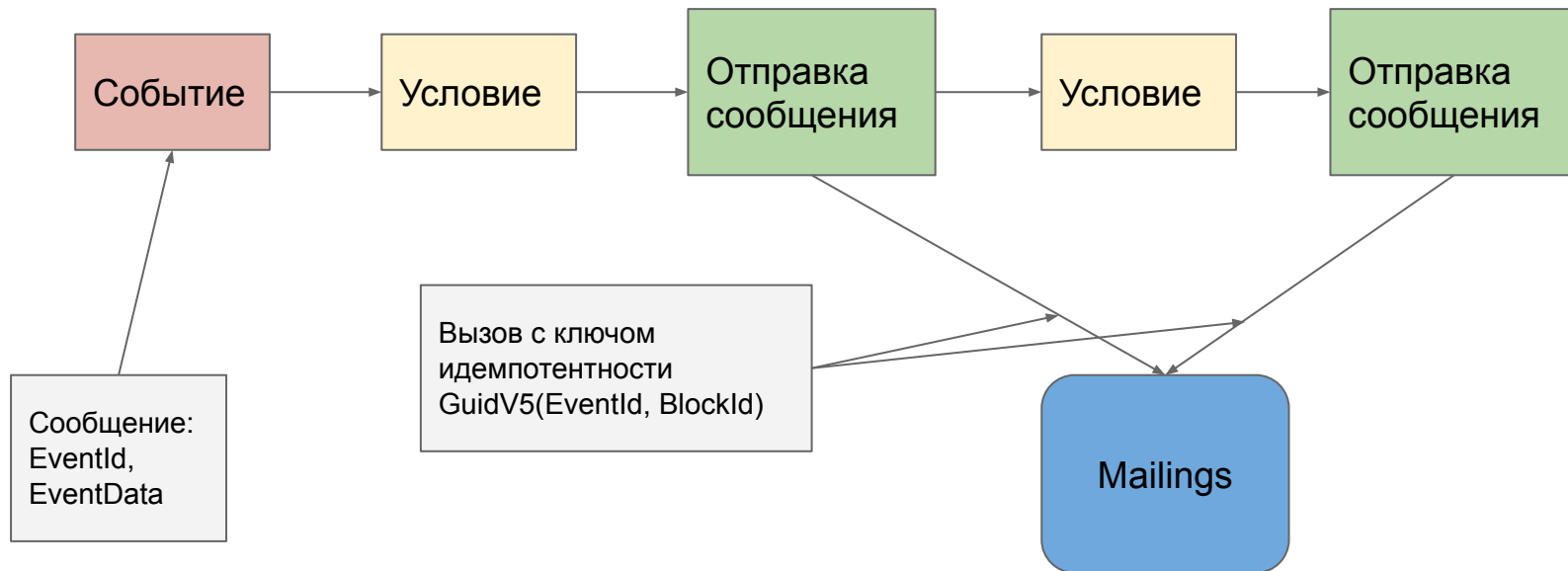
# НЕСКОЛЬКО ВЫЗОВОВ СЕРВИСА



# НЕСКОЛЬКО ВЫЗОВОВ СЕРВИСА

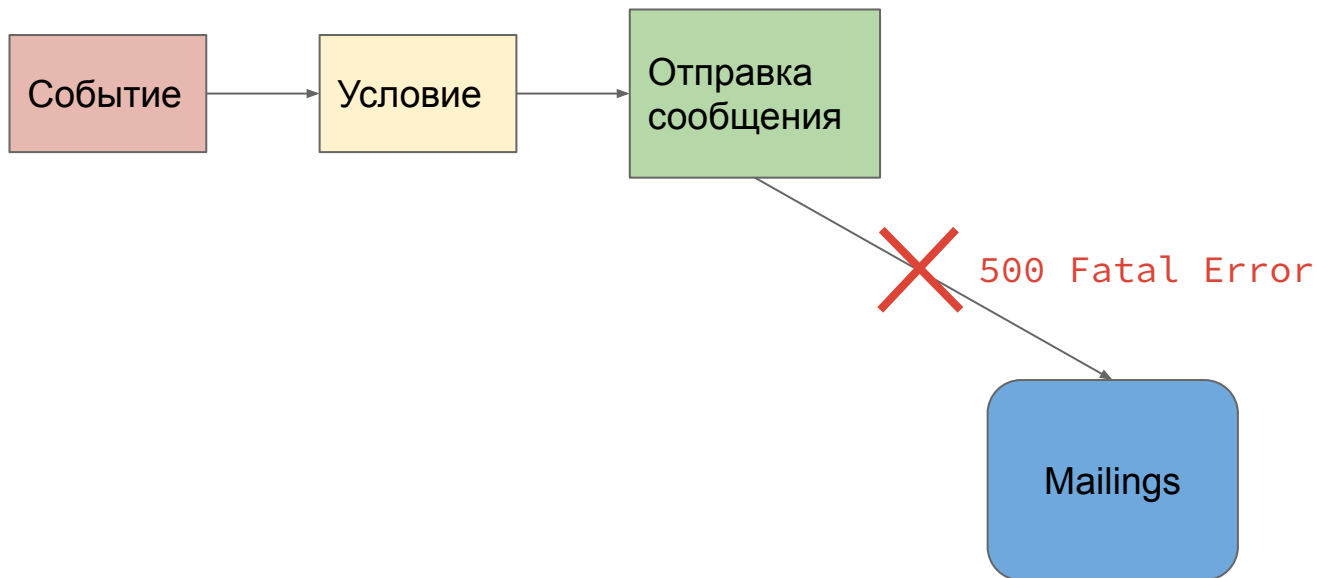


# НЕСКОЛЬКО ВЫЗОВОВ СЕРВИСА

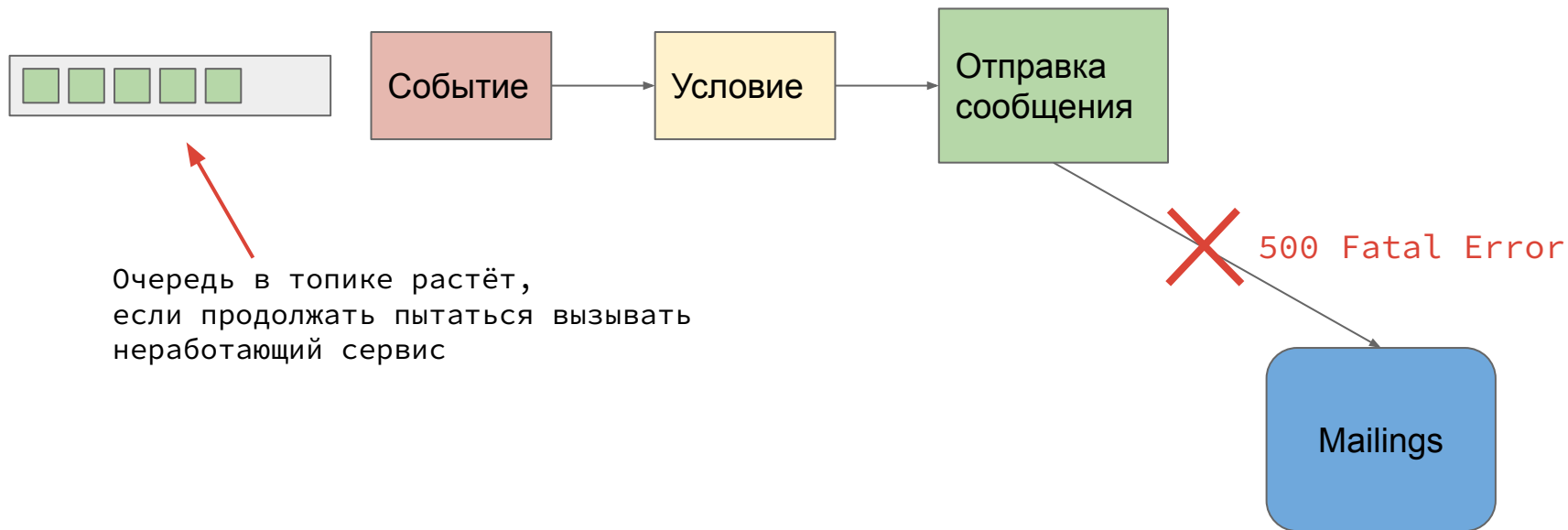




# ОБРАБОТКА ОШИБОК



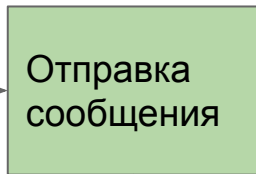
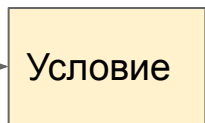
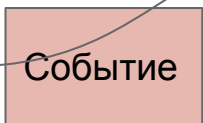
# ОБРАБОТКА ОШИБОК



# ОБРАБОТКА ОШИБОК

Исходный топик

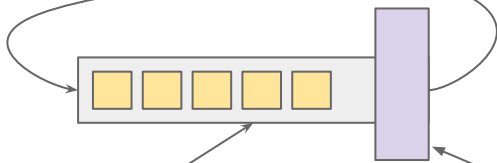
2: COMMIT



1



500 Fatal Error



Топик для задержки

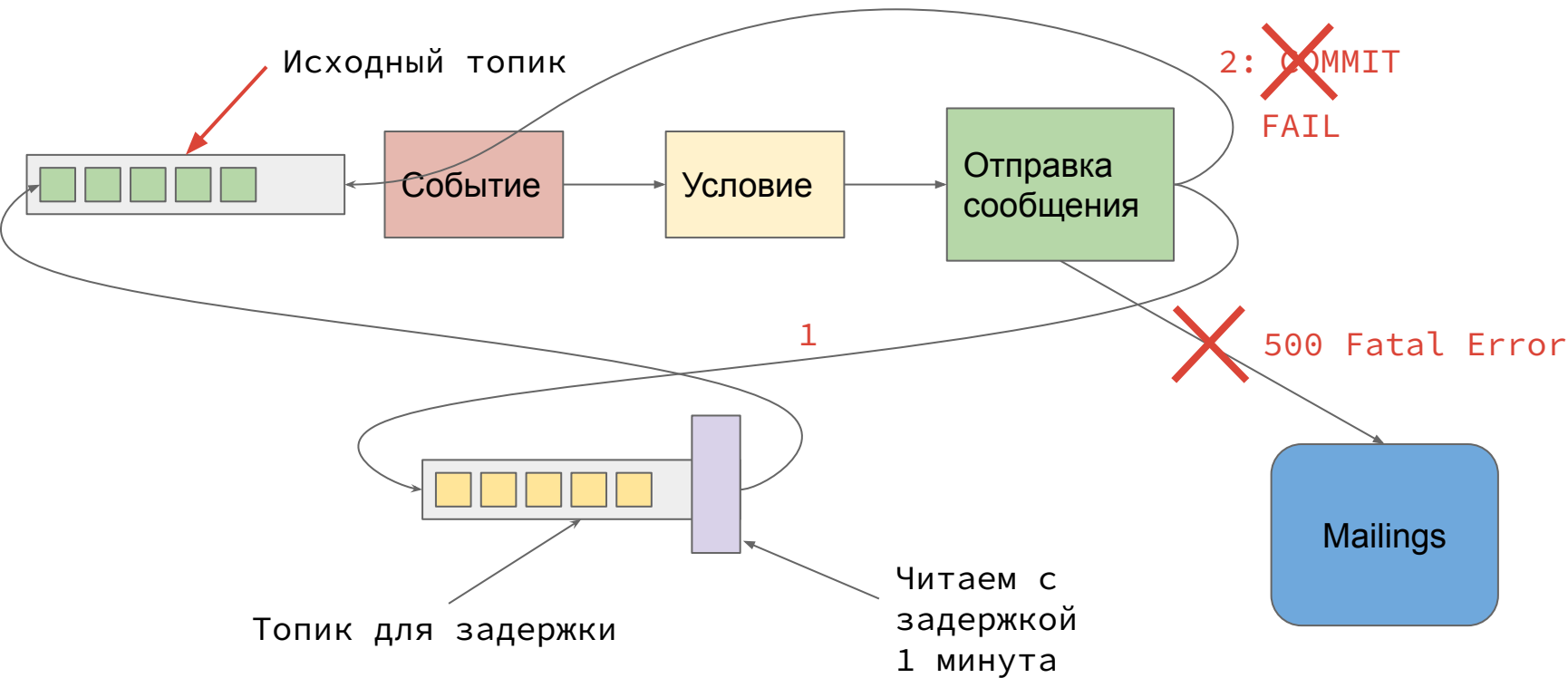
Читаем с задержкой 1 минута



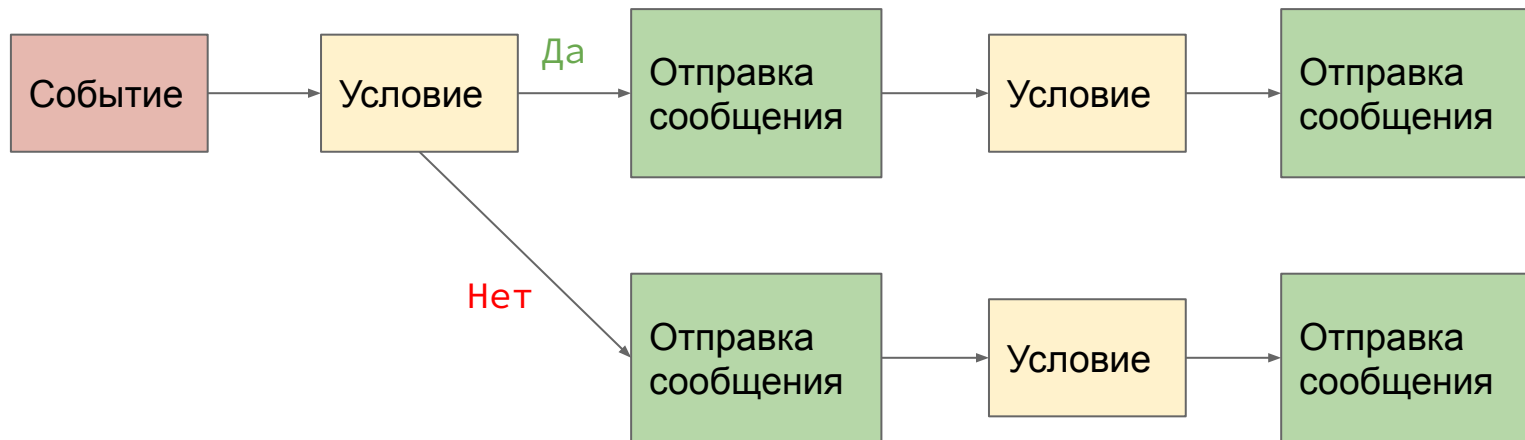
# НЮАНСЫ ПОВТОРОВ ДЛЯ ОБРАБОТКИ ОШИБОК

- Порядок будет нарушен
- Константное время повтора – плохо

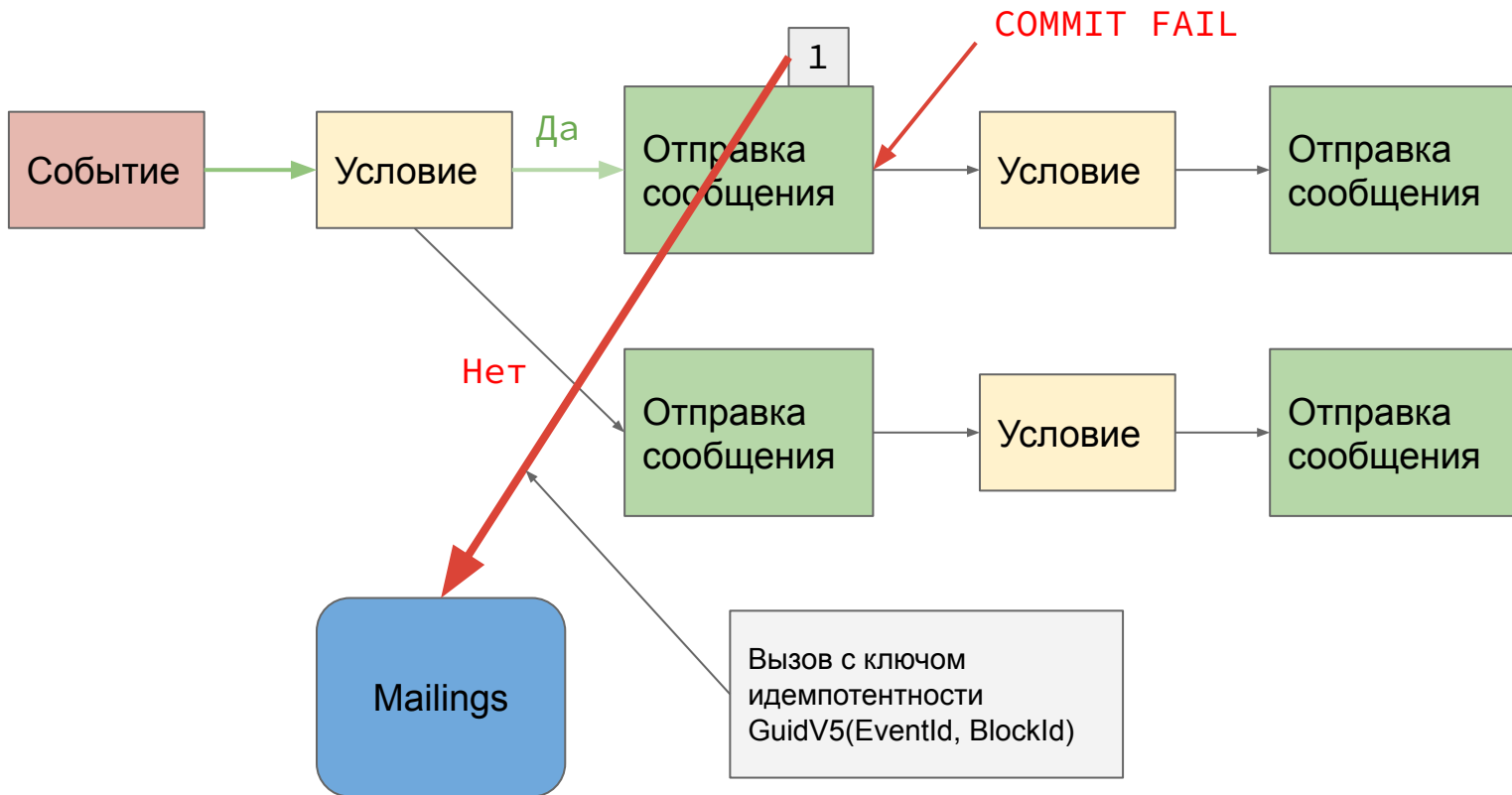
# ЕЩЁ ДУБЛИ



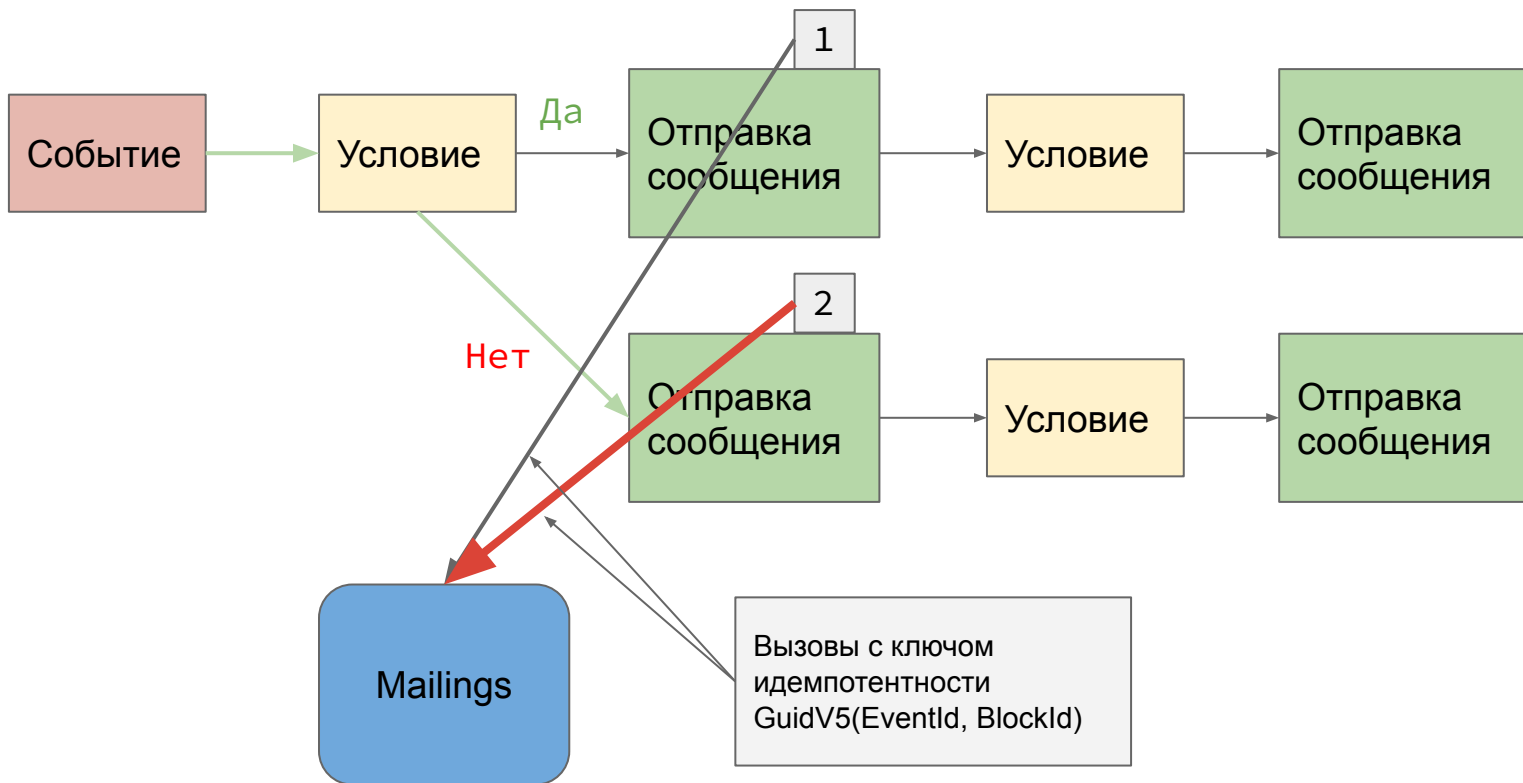
# ВЕТВЛЕНИЕ



# ВЕТВЛЕНИЕ



# ВЕТВЛЕНИЕ





# ИДЕЯ

Идемпотентное применение блоков: сохраняем результат применения блока, при повторе будем проходить по той же цепочке

# РЕЗУЛЬТАТ ПРИМЕНЕНИЯ БЛОКА

Каждый блок может:

- Остановить выполнение сценария
- Отправить в одно из исходящих соединений

# RACES

- Множество примеров дублей сообщений и переобработок
- Могут быть параллельным
- Значит необходима синхронизация

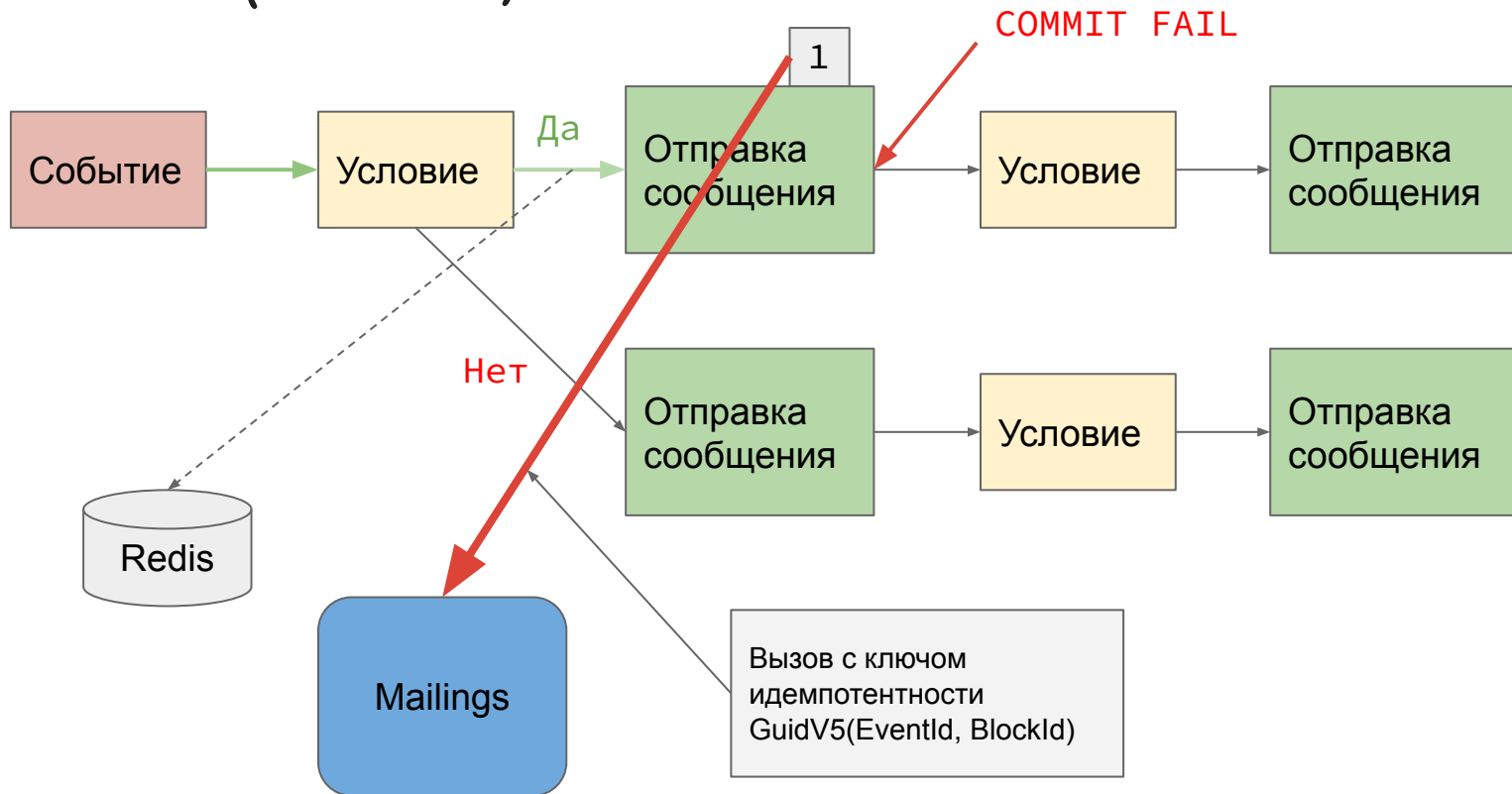
# REDIS

Быстрое key-value хранилище с горизонтальным масштабированием и (очень ограниченными) транзакциями

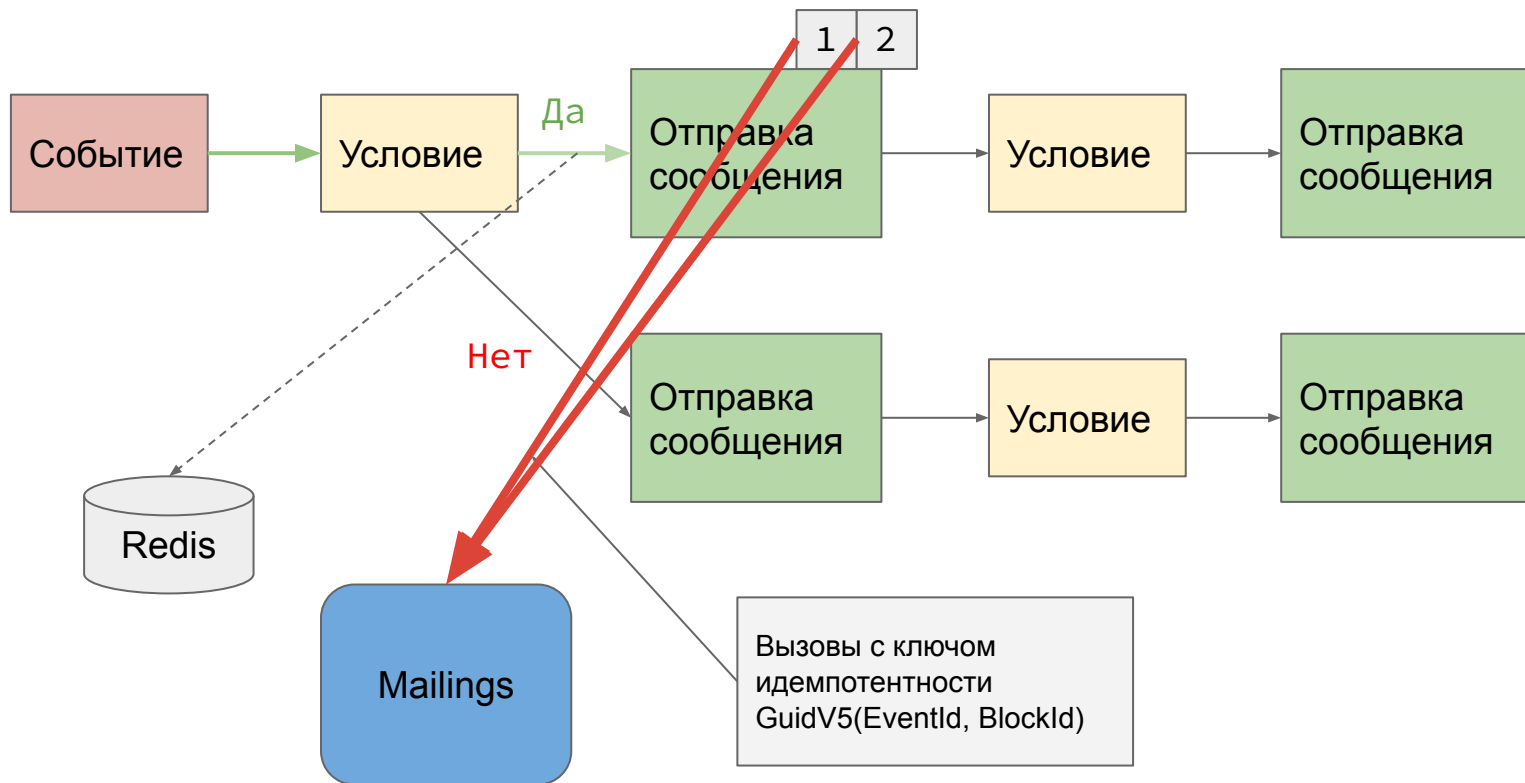
# АЛГОРИТМ

- Применить блок
- По ключу GuidV5(EventId, BlockId) установить результат через SETNX
- Если не удалось, прочитать по ключу результат и использовать его
- Оптимизация: до применения блока пробовать читать результат

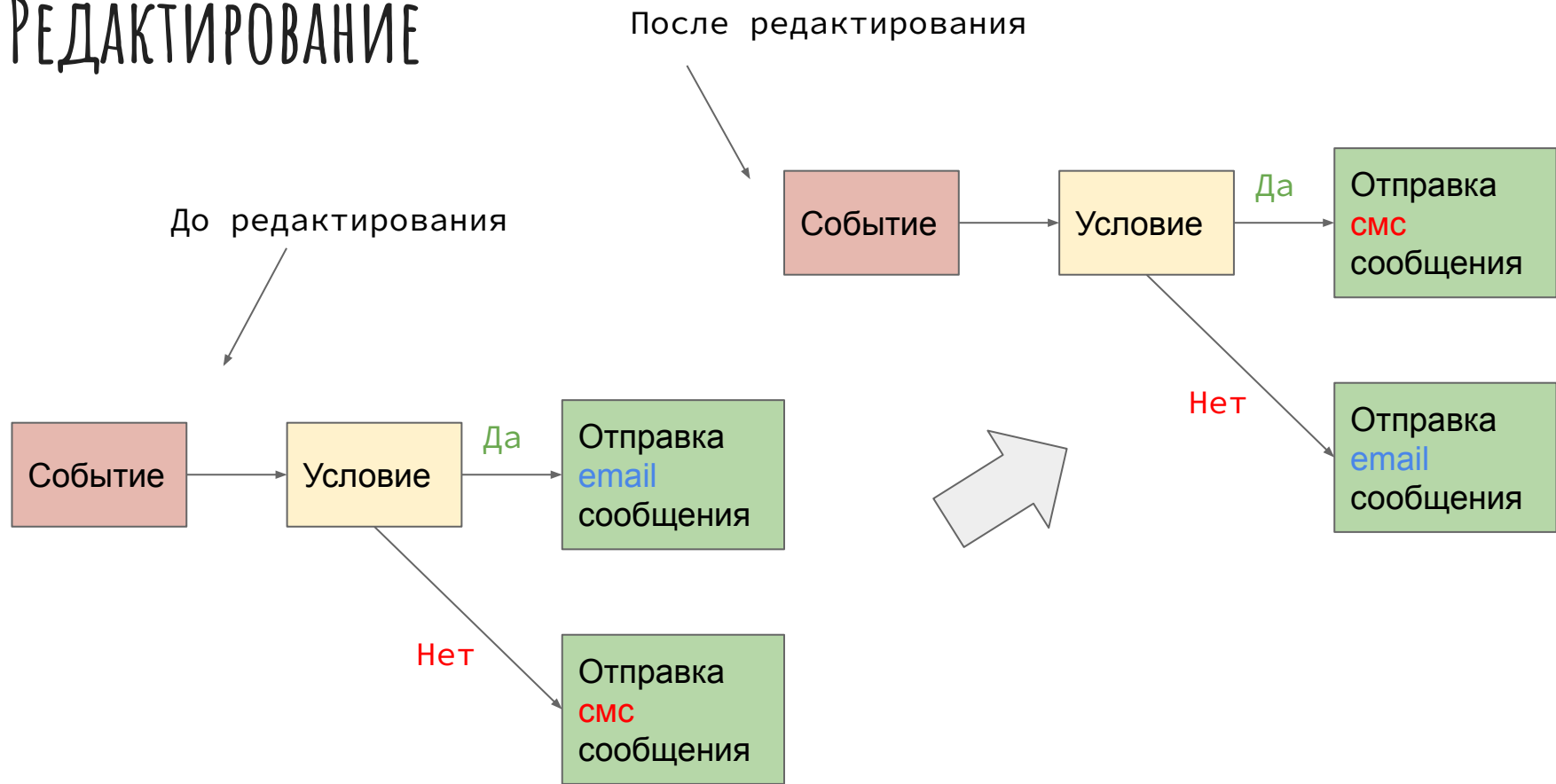
# ВЕТВЛЕНИЕ (РЕШЕНИЕ)



# ВЕТВЛЕНИЕ (РЕШЕНИЕ)

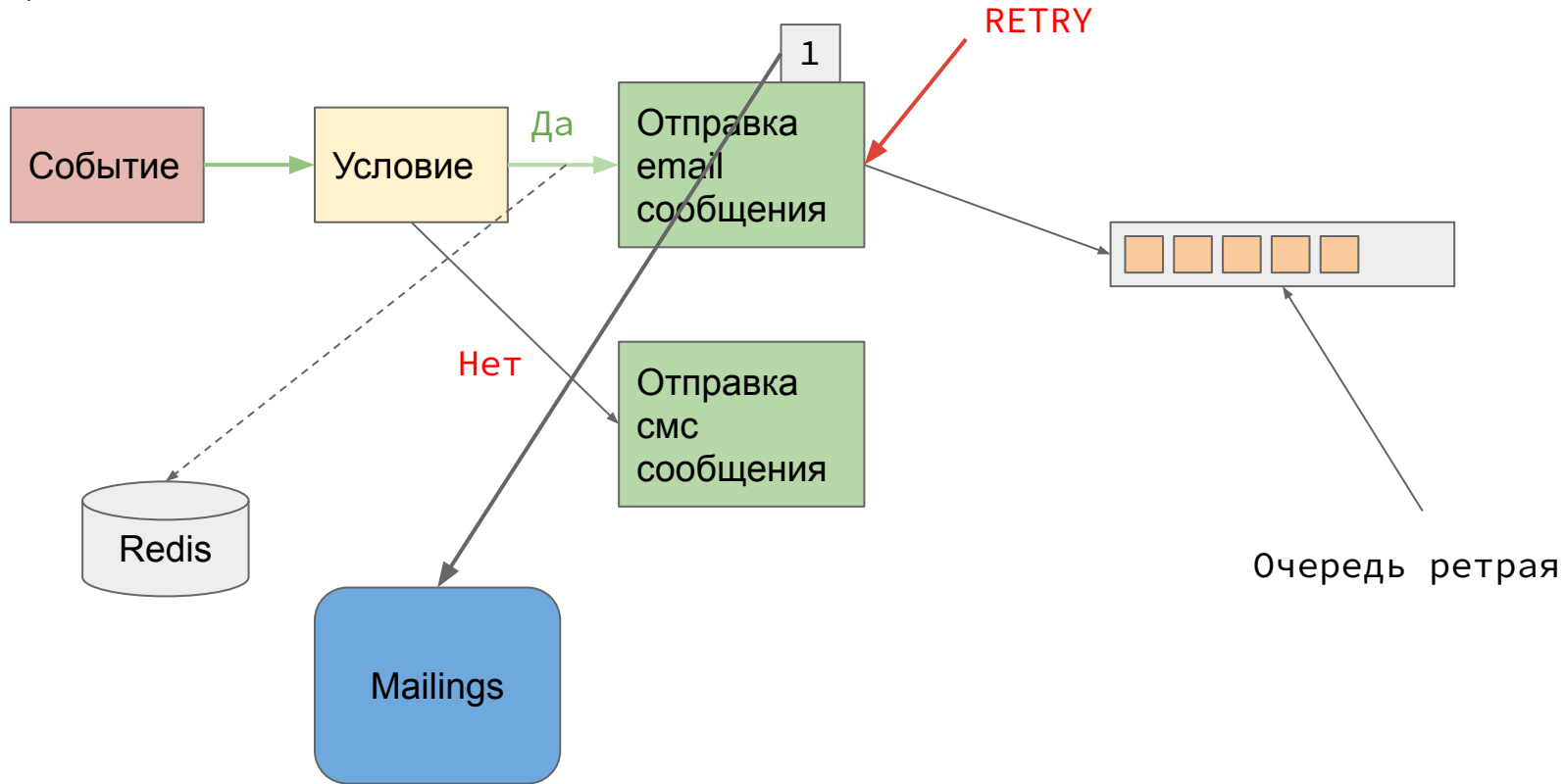


# РЕДАКТИРОВАНИЕ

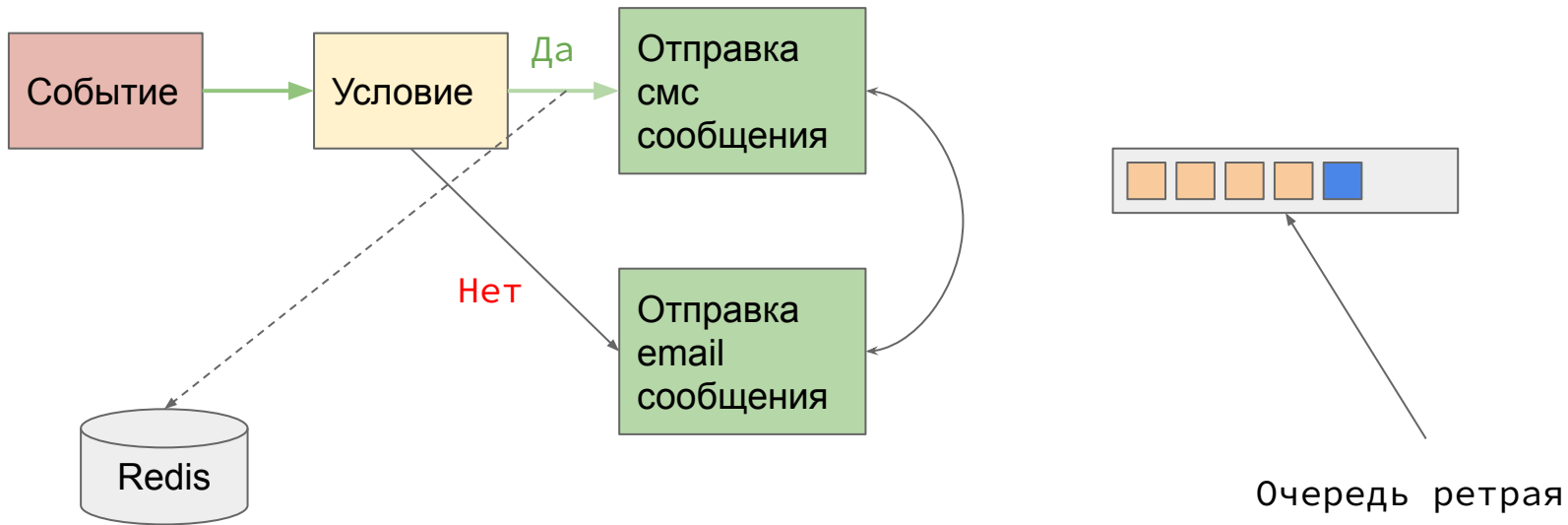




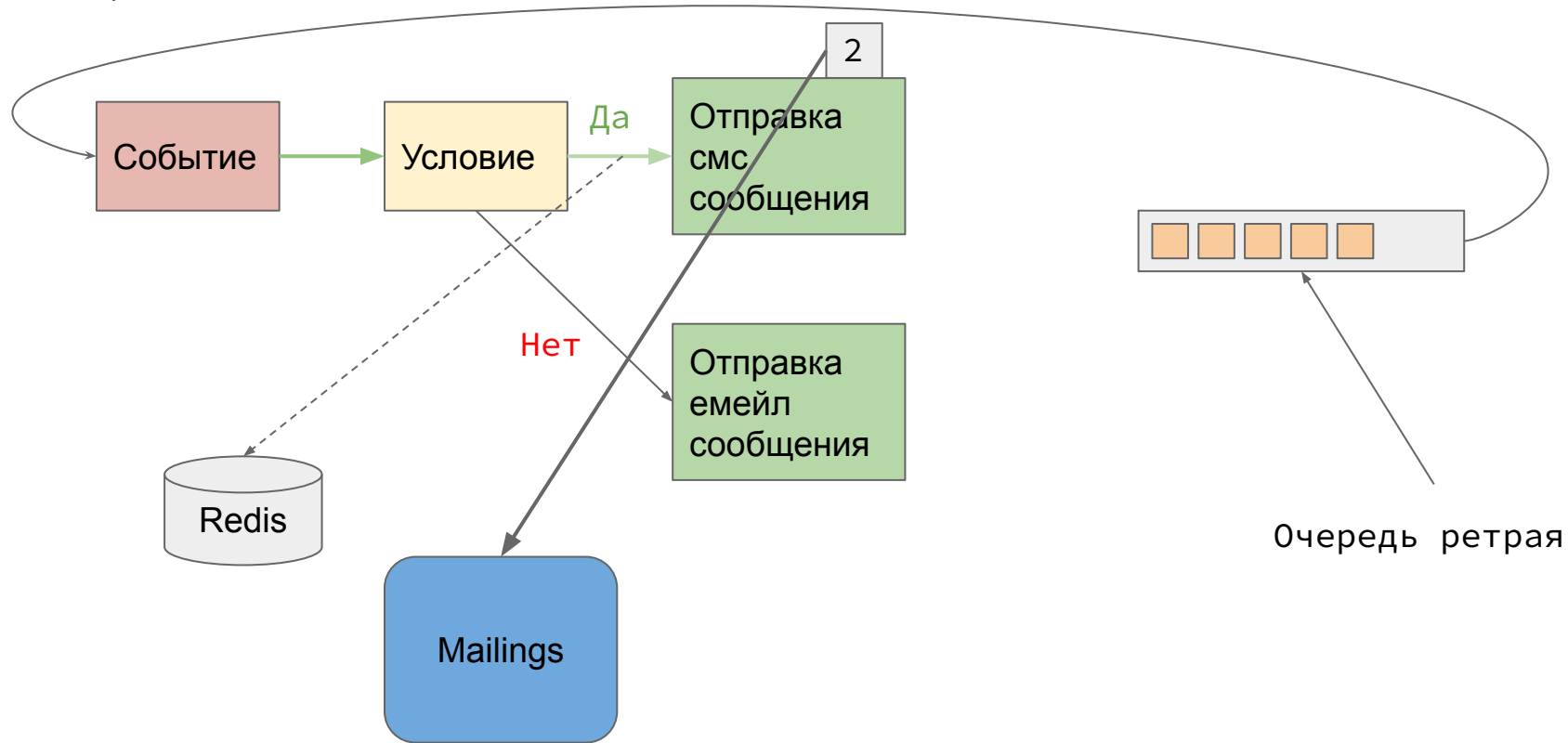
# РЕДАКТИРОВАНИЕ



# РЕДАКТИРОВАНИЕ



# РЕДАКТИРОВАНИЕ

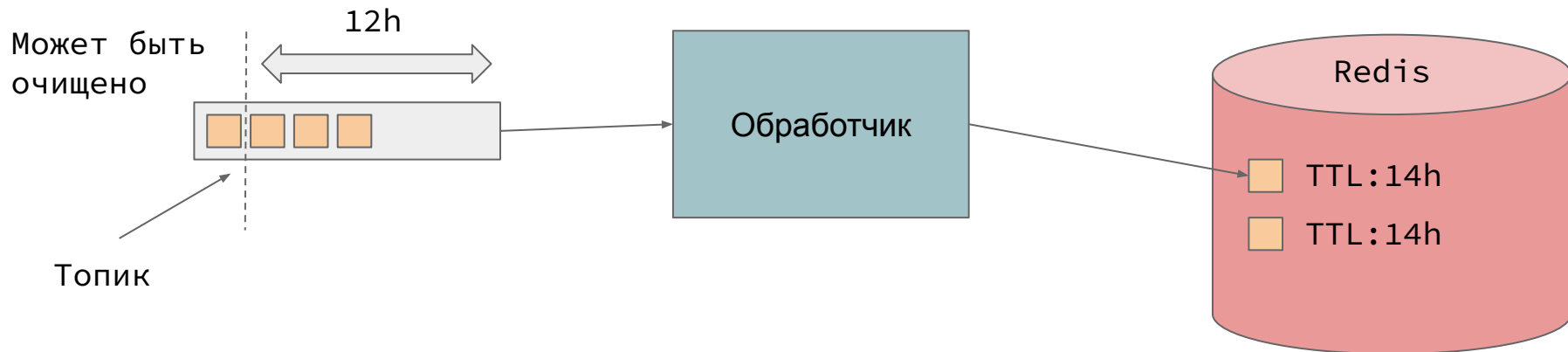


# ЧТО ДЕЛАТЬ ПРИ РЕДАКТИРОВАНИИ МЕТАДААННЫХ?

- Хранить историю метаданных
- Гарантировать использование нужной версии метаданных

# СКОЛЬКО ХРАНИТЬ КЛЮЧИ ИДЕМПОТЕНТНОСТИ?

- Retention в кафке исходя из скорости обработки
- Redis TTL не меньше Retention в кафке из-за ретраев
- Обновление TTL при успешных чтениях (GETEX + PXAT)



# ЧТО ДЕЛАТЬ НА СТОРОНЕ ВЫЗЫВАЕМОГО МИКРОСЕРВИСА?

- Readonly? Ничего (но не всегда).
- Mutable? Транзакция + хранение ключа.
- Внешний вызов? Опционально ТО + идемпотентность.

# ХРАНЕНИЕ КЛЮЧА ПРЯМО В ДАННЫХ

CustomerMailings

<b>IdempotentKey</b>	<b>CustomerId</b>	<b>MailingId</b>	<b>DateTimeUtc</b>	<b>...</b>
un42ndt881	1	100	01.01.2020	..
infjs9a3003	2	200	02.01.2020	..

# ОТДЕЛЬНАЯ ТАБЛИЧКА ДЛЯ ОБОБЩЕННЫХ ТРАНЗАКЦИЙ

IdempotentKeys

<b>IdempotentKey</b>	<b>DateTimeUtc</b>
un42ndt881	01.01.2020
infjs9a3003	02.01.2020

+ очистка по времени



# Что получилось

- Всё работает
- Стоит как каменный мост из-за редиса :)

## Итого:

Redis™ 152 709,81 ₽ ×

Redis. Intel Cascade Lake.  
100% vCPU 27 216,00 ₽

Redis. Intel Cascade Lake.  
RAM 116 121,60 ₽

Быстрое сетевое  
хранилище — Redis 9 372,21 ₽

Исходящий трафик в  
интернет 0,00 ₽

152 709,81 ₽

в месяц ▾

# ЕЩЁ РАЗ, А ЗАЧЕМ НАМ РЕДИС?

- Борьба с состоянием гонки
- Хранение номера версии сценария
- Хранение результатов применения блоков

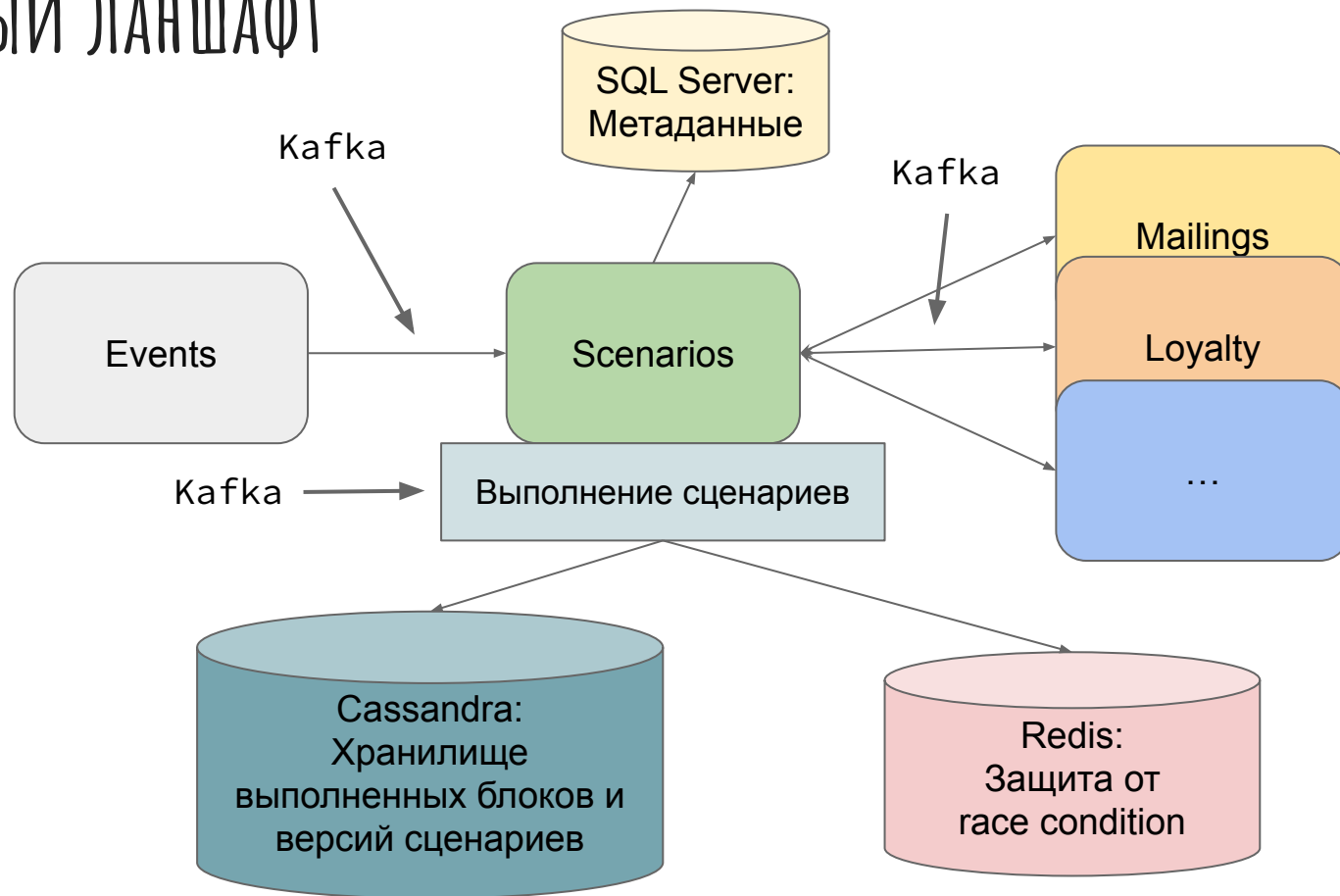
# ЕЩЁ РАЗ, А ЗАЧЕМ НАМ РЕДИС?

- Борьба с состоянием гонки
- ~~Хранение версии сценария~~
- ~~Хранение результатов применения блоков~~

# CASSANDRA: ЧТО-ТО ПОЛУЧШЕ

- Горизонтально масштабируется
- Надёжна
- Идеальна для кейса хранения ключ-значения
- Есть TTL
- Есть транзакции, но лучше не использовать

# НОВЫЙ ЛАНШАФТ



# САММАРИ

Для межсервисного взаимодействия:

- ТО
- Идемпотентность как зомби
- Правильный выбор технологий



# ЛИТЕРАТУРА

- Доклад по Kafka из первого дня :)
- [microservices.io](#) + [Microservices Patterns](#)
- [Kafka: The Definitive Guide](#)
- [Kafka VS Rabbit](#)

# КОНЕЦ

Вопросы приветствуются

Телеграмм **@timramone**