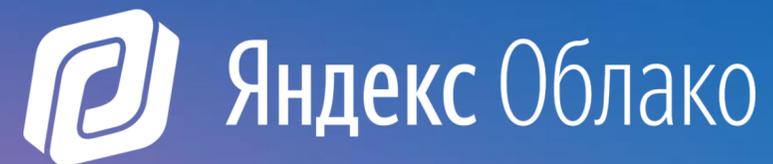


Яндекс



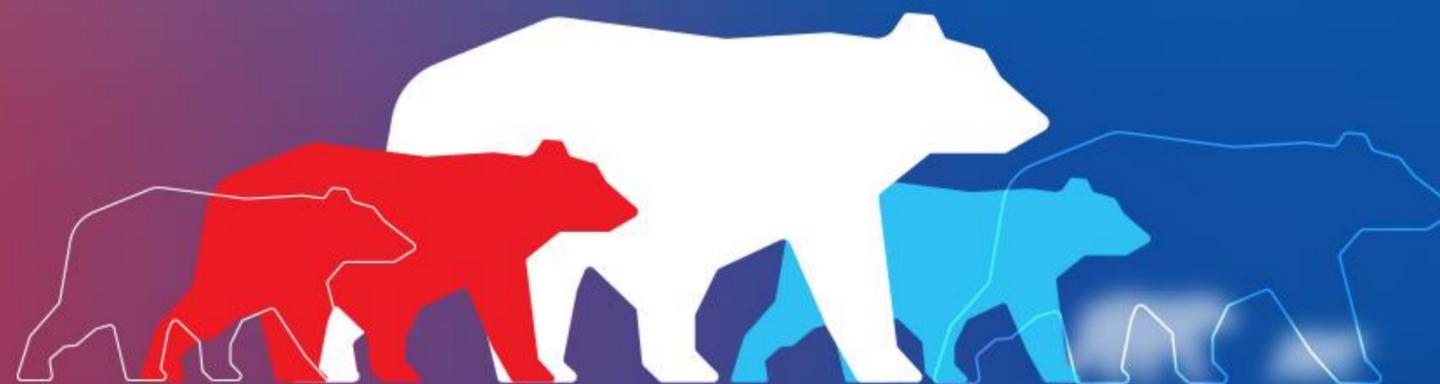
# Реализация геораспределённой персистентной очереди сообщений на примере Yandex Message Queue

Василий Богонатов, разработчик сервиса



## HighLoad<sup>++</sup> Siberia 2019

Профессиональная конференция  
для разработчиков высоконагруженных  
систем



# План

- 01 | Зачем нужны очереди сообщений?
- 02 | Гарантии доставки сообщений
- 03 | Сравнение популярных реализаций
- 04 | Свойства и API YMQ
- 05 | Внутреннее устройство YMQ

01

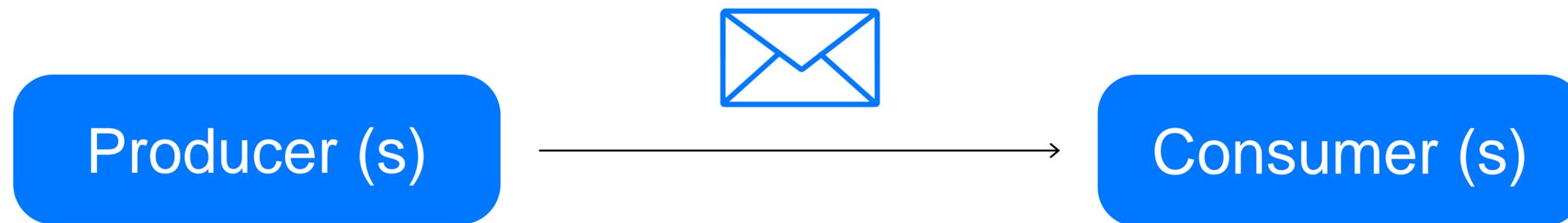
**Зачем нужны очереди  
сообщений?**

# Где мы используем очереди сообщений?

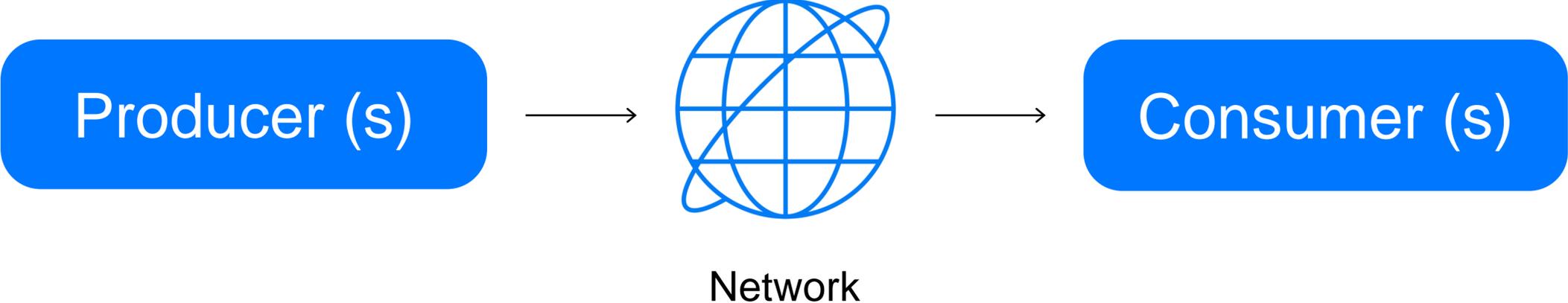


- › Доставка логов до хранилища
- › Вычислительные задачи
- › Обработка потока событий
- › Взаимодействие микросервисов

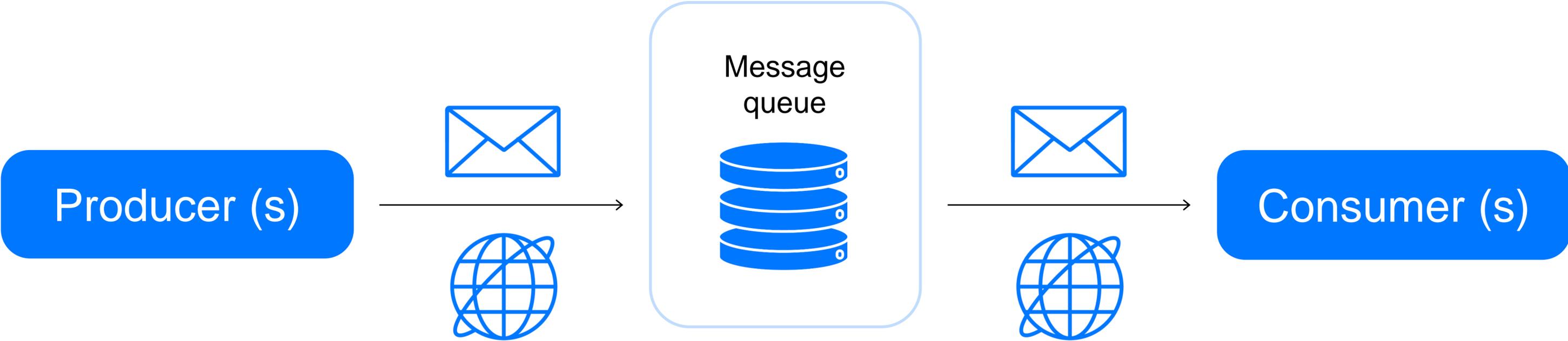
# Задача надёжной передачи сообщений



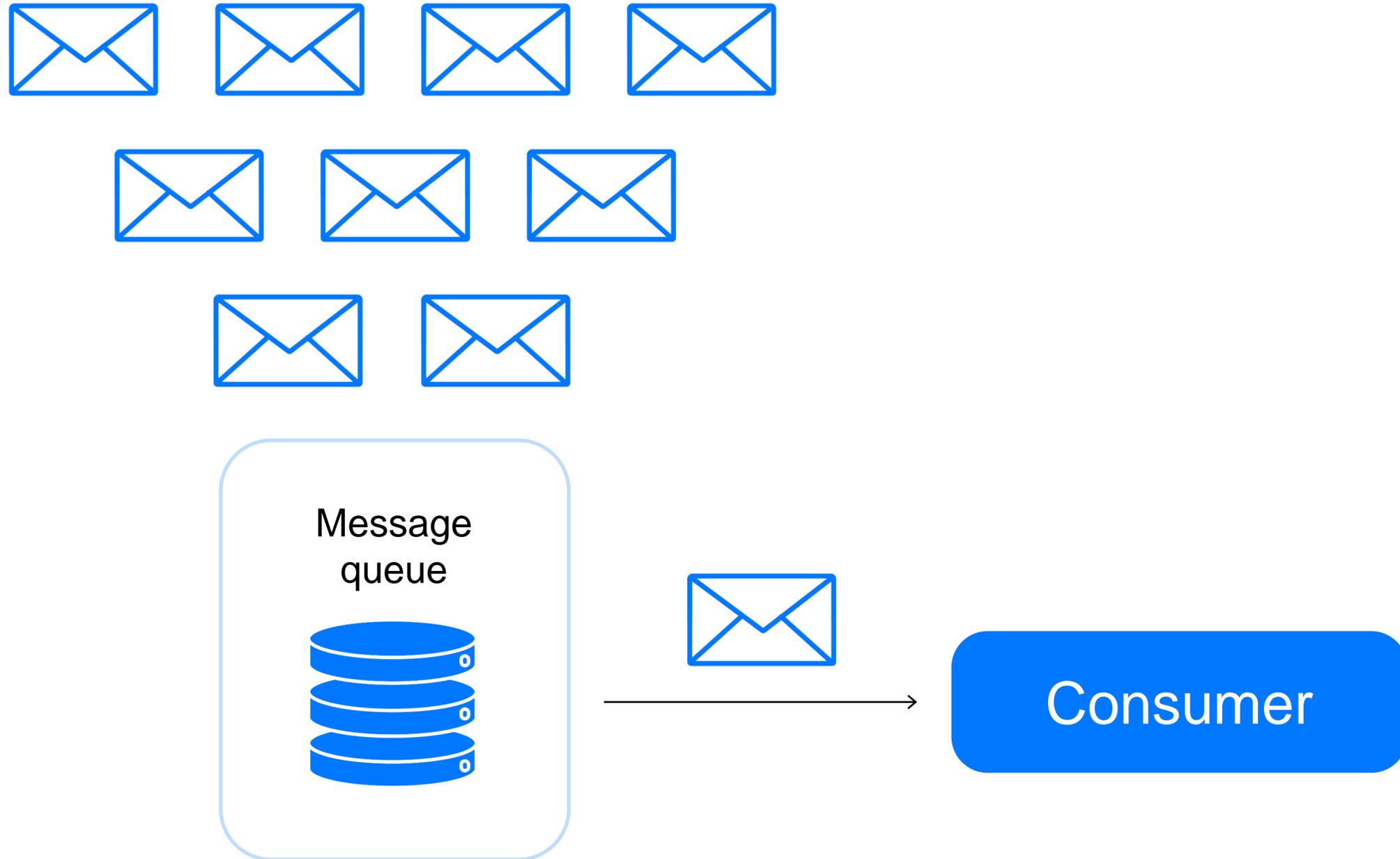
# Наивное решение



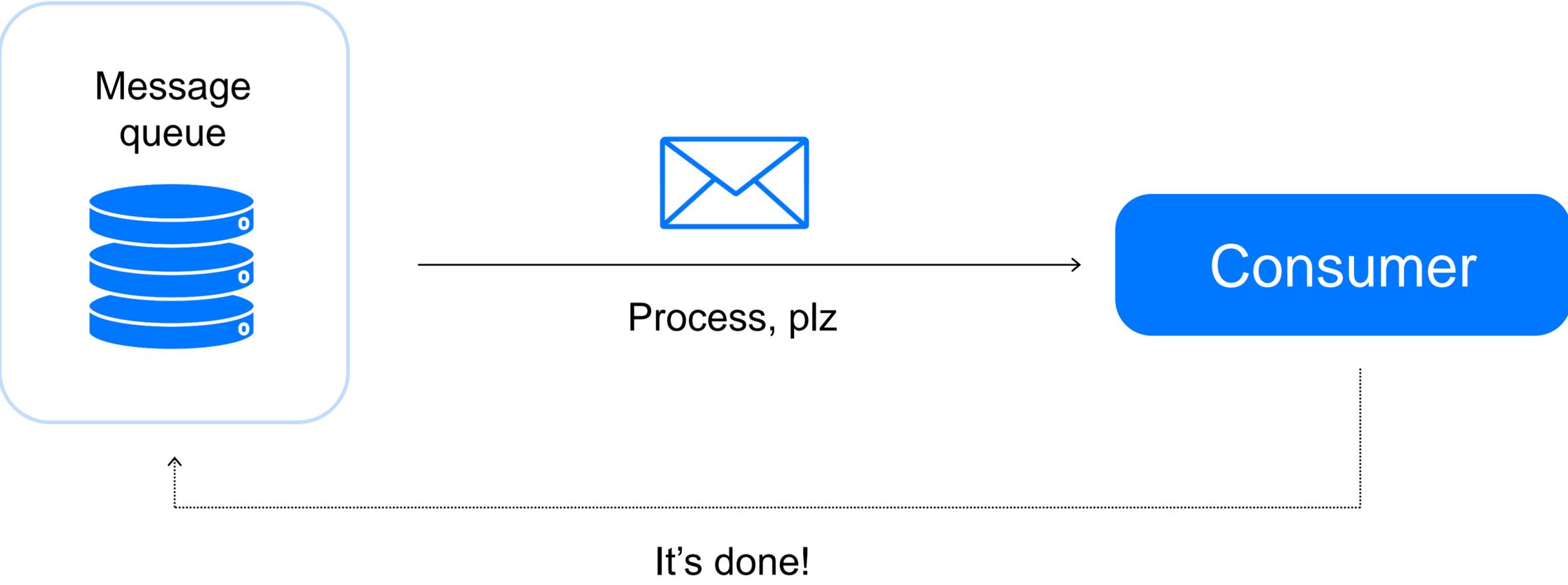
# Решение с очередью сообщений



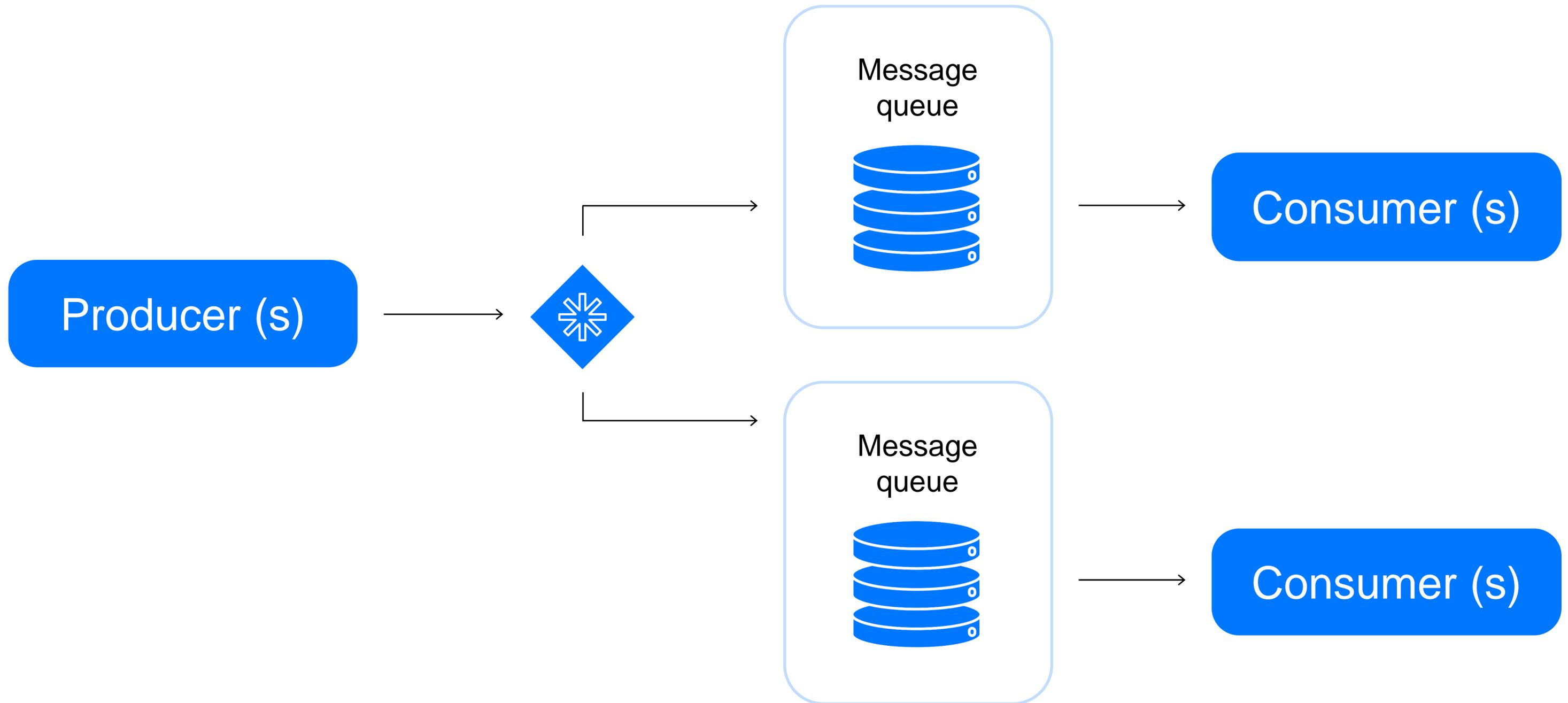
# Плюшки очередей: хранение сообщений



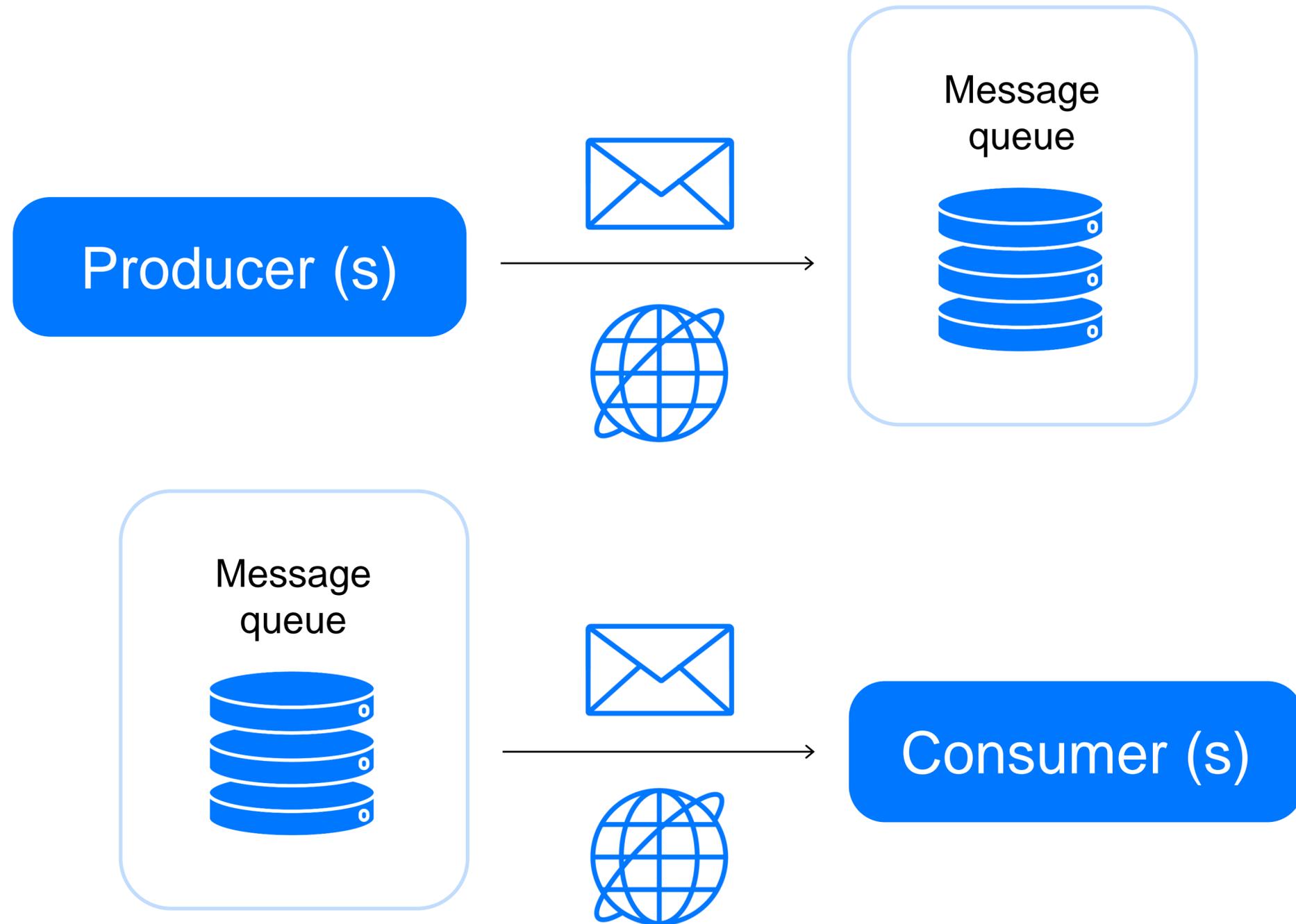
# Плюшки очередей: контроль обработки



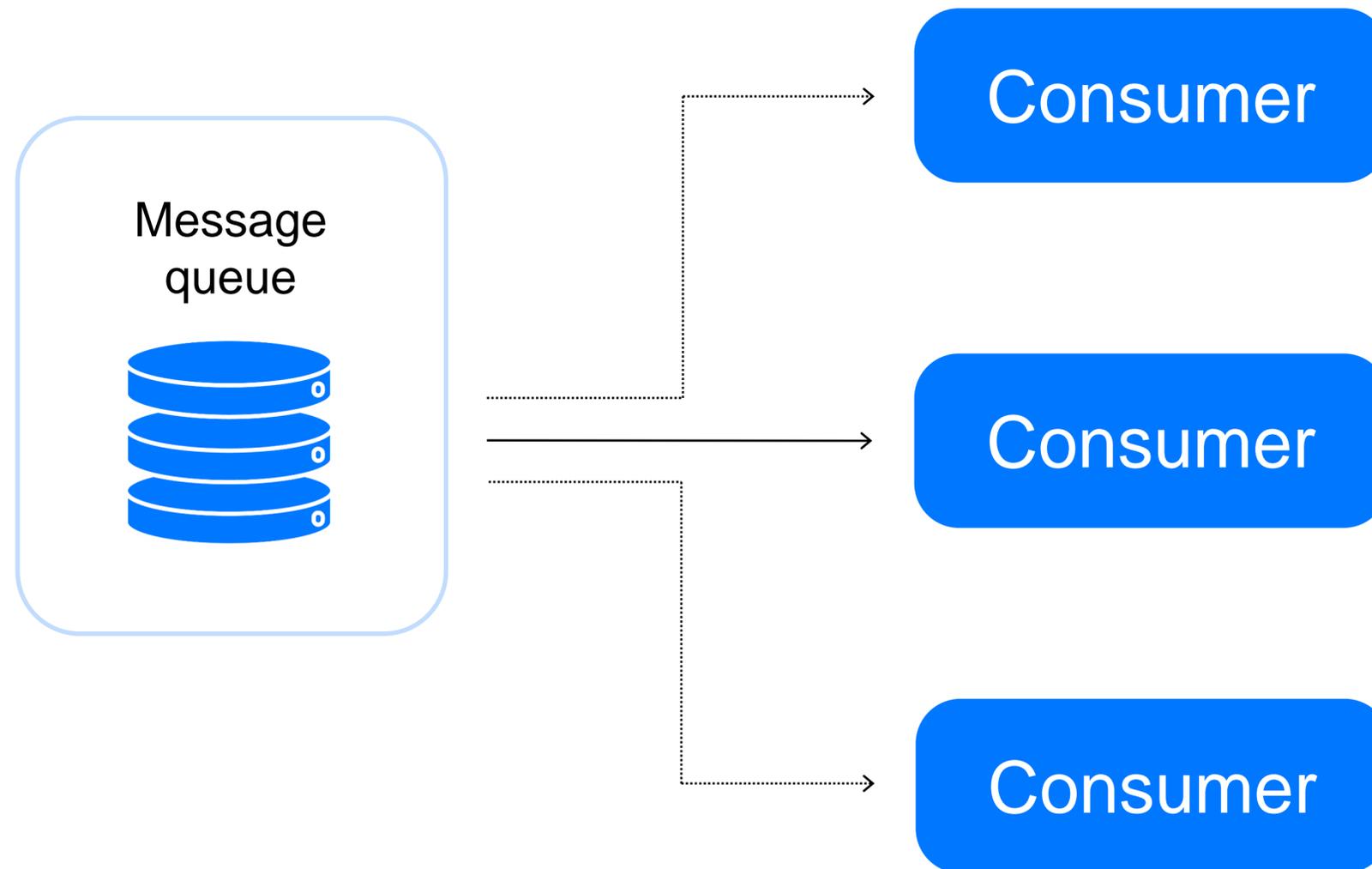
# Плюшки очередей: маршрутизация



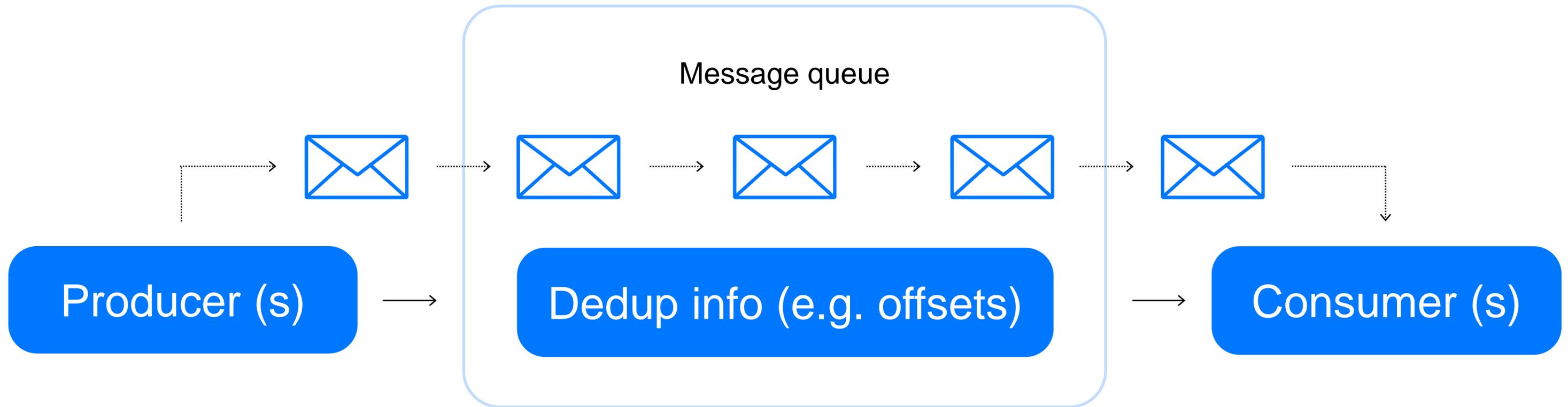
# Плюшки очередей: развязка компонентов



# Плюшки очередей: лёгкое масштабирование



# Плюшки очередей: порядок и дедупликация



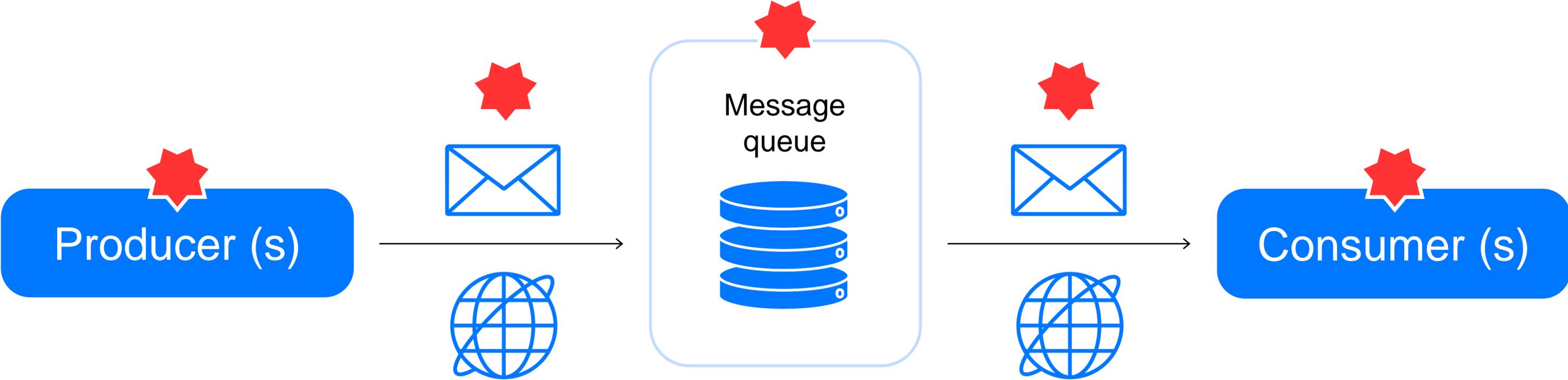
# Потому что можем ~~НИ~~ нужно!

- › RabbitMQ
- › Apache Kafka
- › Apache ActiveMQ
- › Amazon SQS
- › Amazon Kinesis
- › Azure Service Bus
- › Azure Storage Queue
- › Google PubSub
- › IBM MQ
- › Redis Pub/Sub
- › Apache Apollo
- › Beanstalkd
- › HornetQ
- › IronMQ
- › NATS
- › ZeroMQ
- › И многие другие...

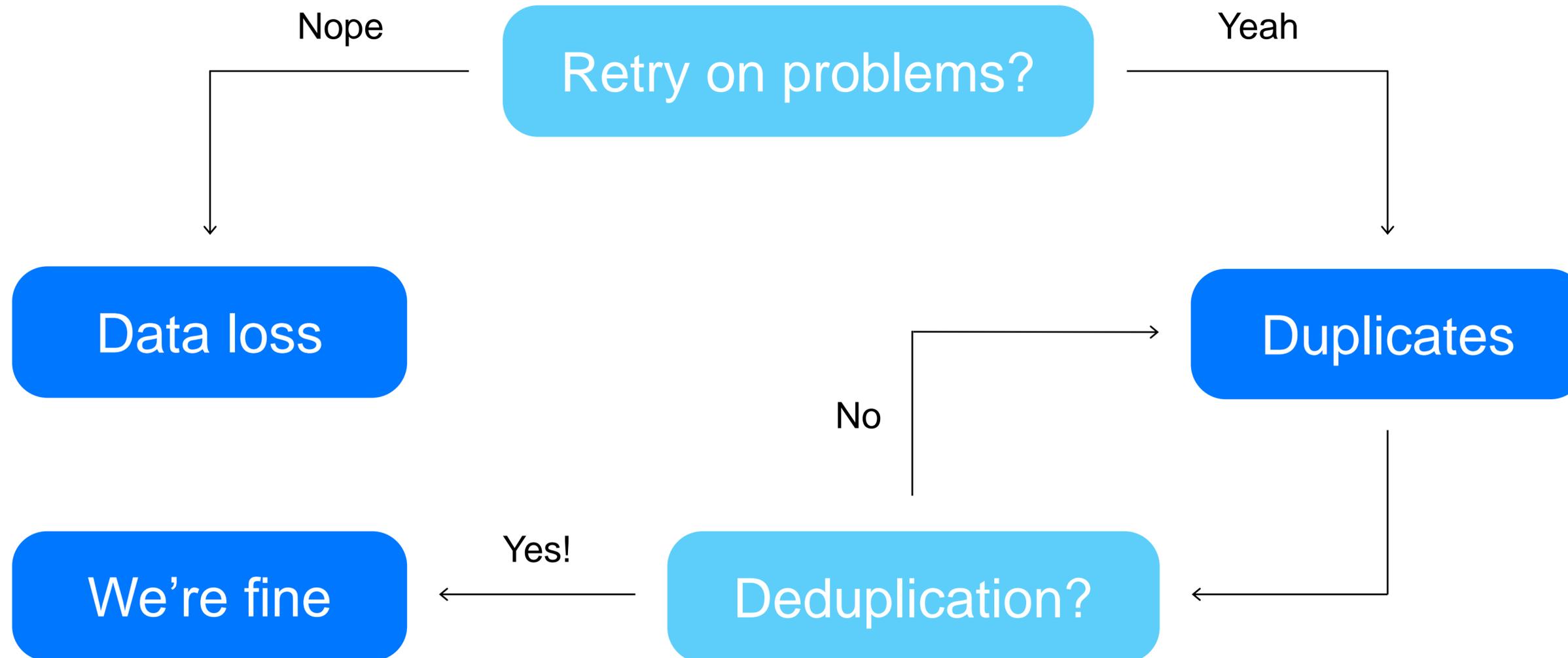
02

# Гарантии доставки сообщений

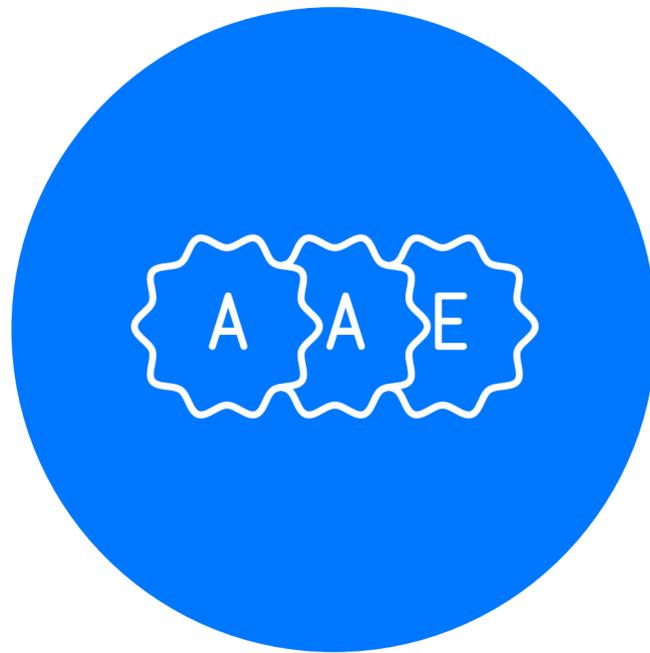
# Откуда они берутся?



# Что делать при отказах (помимо паники)?



# Виды гарантий



- › At-most-once — максимум раз, теряем данные
- › At-least-once — минимум раз, получаем дубли
- › Exactly-once — ровно раз, идеально!

# Примеры гарантий на практике



- › At-most-once — реалтаймовые данные, почта
- › At-least-once — почти везде
- › Exactly-once — (вы верите, что у вас так)

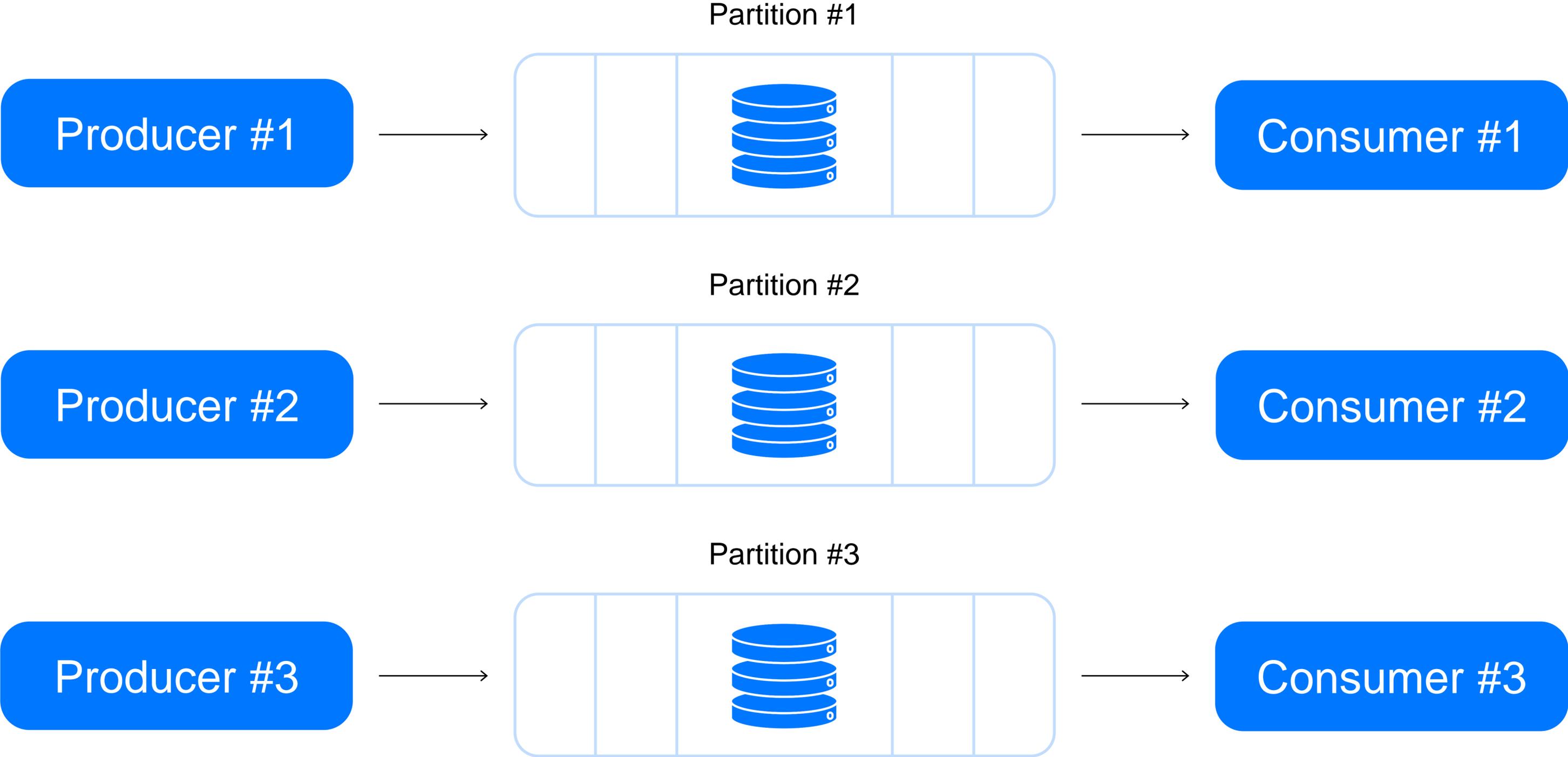
# Комбинации гарантий



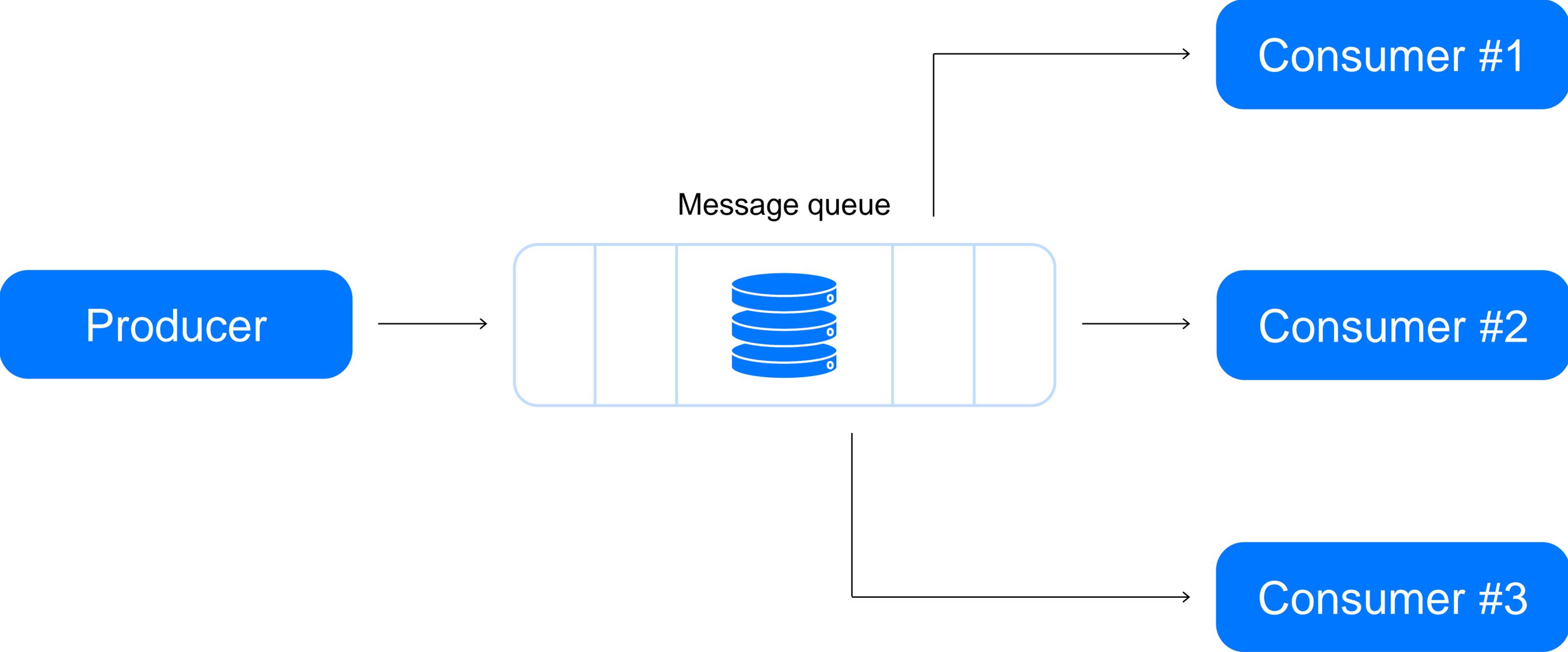
- | **At-least-once + at-most-once  $\neq$  Exactly-once**
- | **At-least-once + at-most-once = дубли + потери**
- | **Exactly-once + ???-once = ???-once**



# Apache Kafka



# RabbitMQ

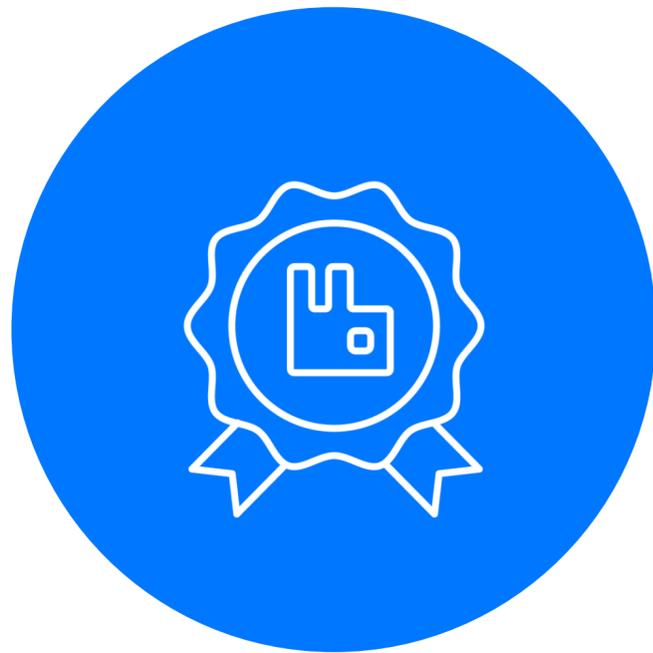


# Apache Kafka: гарантии доставки



- › At-most-once — коммитим offset до обработки
- › At-least-once — коммитим offset после неё
- › Exactly-once\*\*\* — с 0.11.0.0

# RabbitMQ: гарантии доставки



- › At-most-once — no ack
- › At-least-once — manual ack

# int fsync (int fd)



- › fsync синхронизирует состояние файла в памяти с состоянием на диске

# Apache Kafka: надёжность хранения



- › Отвечает без явного вызова `fsync`

# RabbitMQ: надёжность хранения



- › Отвечает после `fsync()`, но может переключиться на несинхронизированное зеркало в качестве мастера

**Bonus level:**  
**зачем свою запилили-то?**

# И всё-таки?



- › Целый велопарк у внутренних клиентов
- › Максимально простое API и готовые библиотеки
- › Стабильная работа в условиях –1 ДЦ
- › Действительно надёжное хранение
- › Наличие экспертизы и платформы
- › Очереди для экосистемы Яндекс.Облака

04

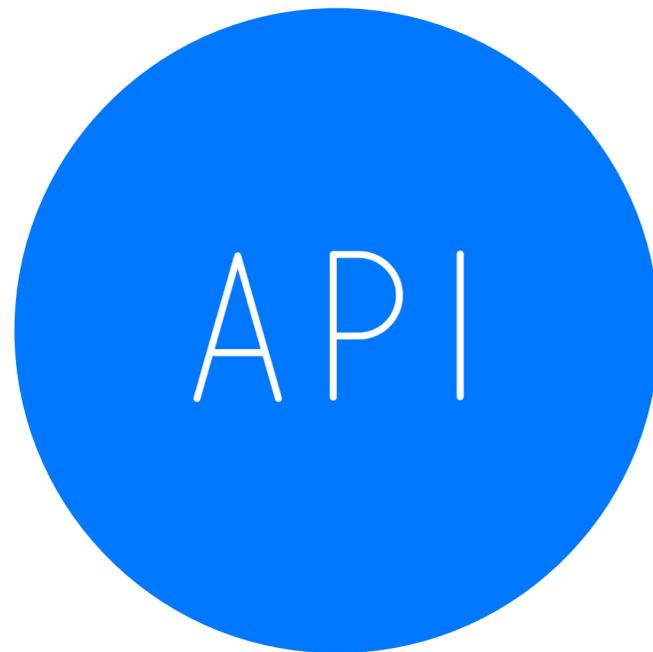
# Свойства и API YMQ

# Основные свойства очередей YMQ



- › Совместимы с API Amazon SQS (http-api)
- › Персистентные
- › Хранят данные в 3 репликах в разных ДЦ
- › Реализуют паттерн “competing consumers”
- › Бывают двух типов: standard и FIFO
- › Обеспечивают доставку сообщений at-least-once

# Основные методы API



- › CreateQueue
- › SendMessage
- › ReceiveMessage
- › DeleteMessage

# Создание очереди

Producer

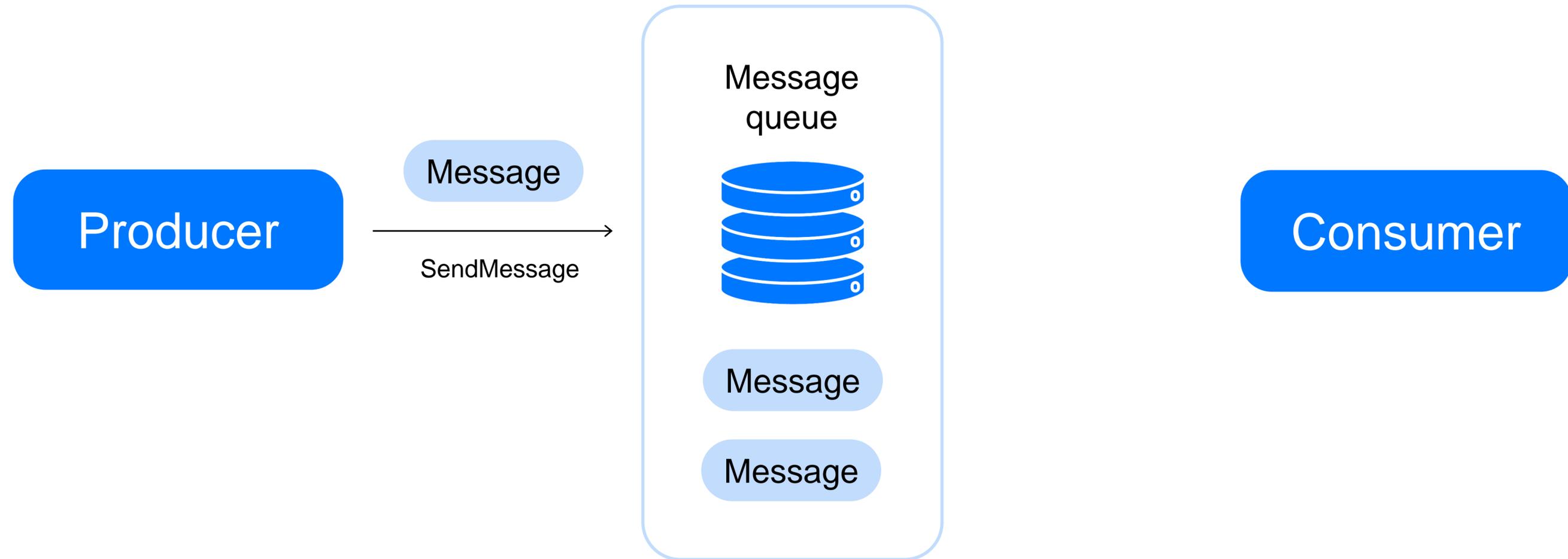


Consumer

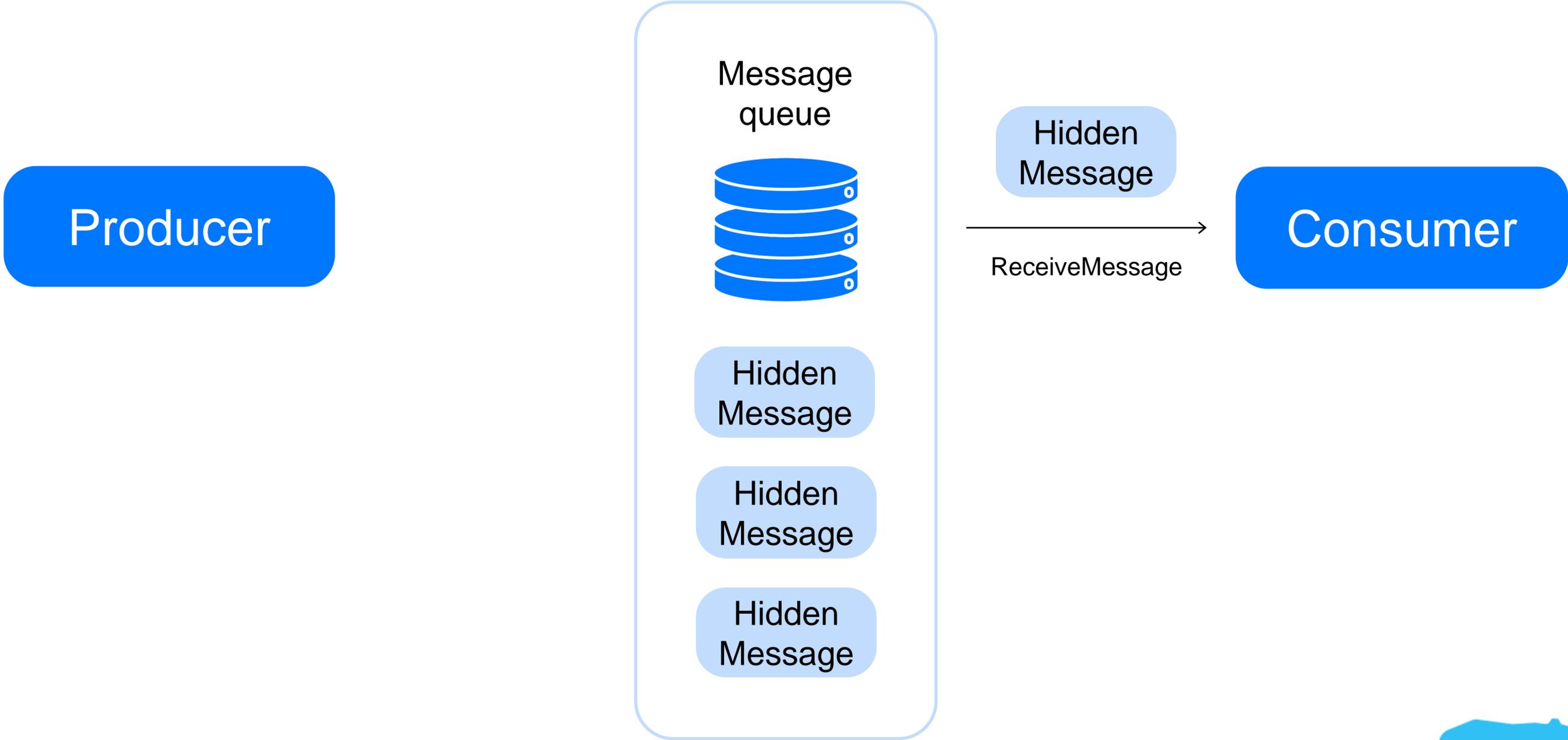
Client

CreateQueue

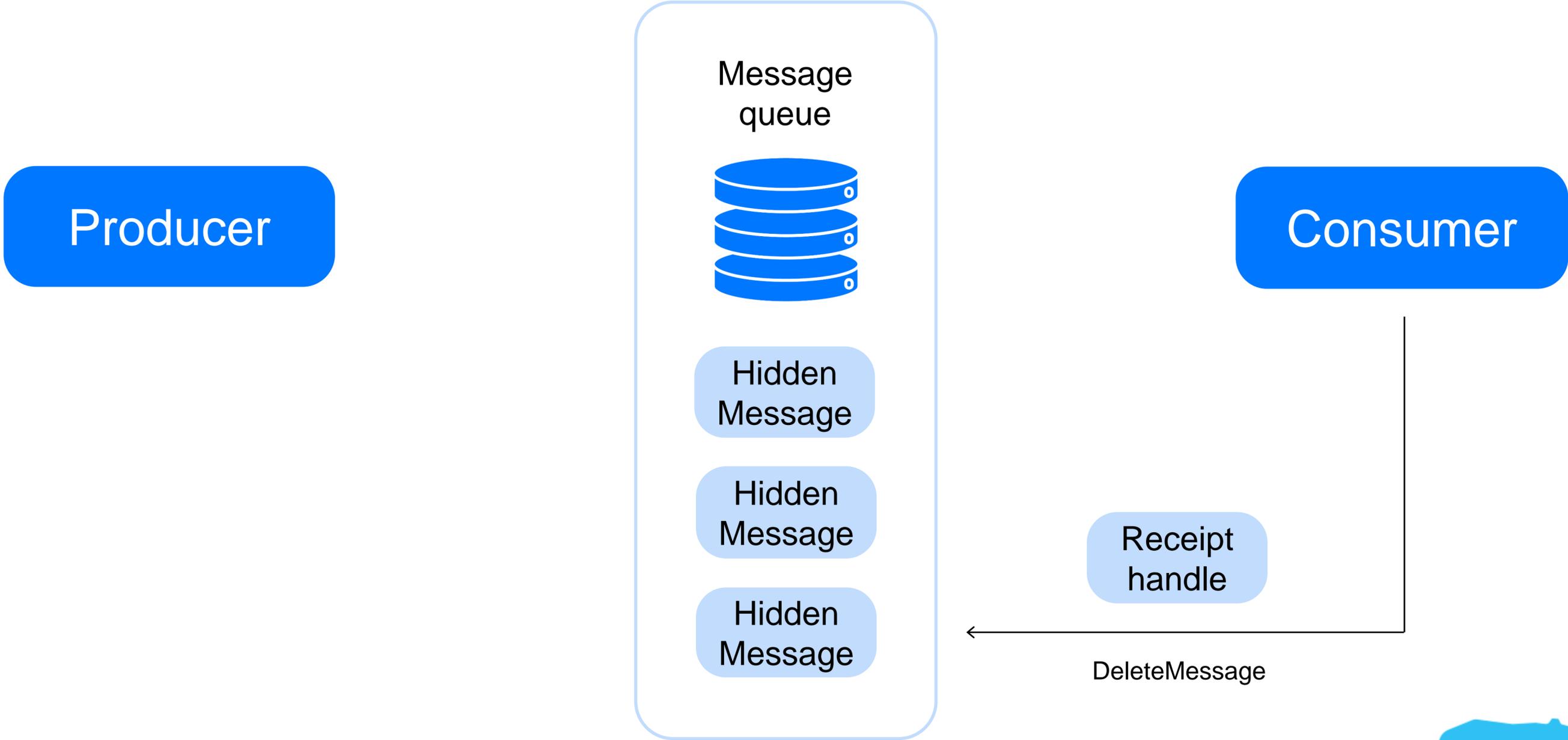
# Запись сообщения



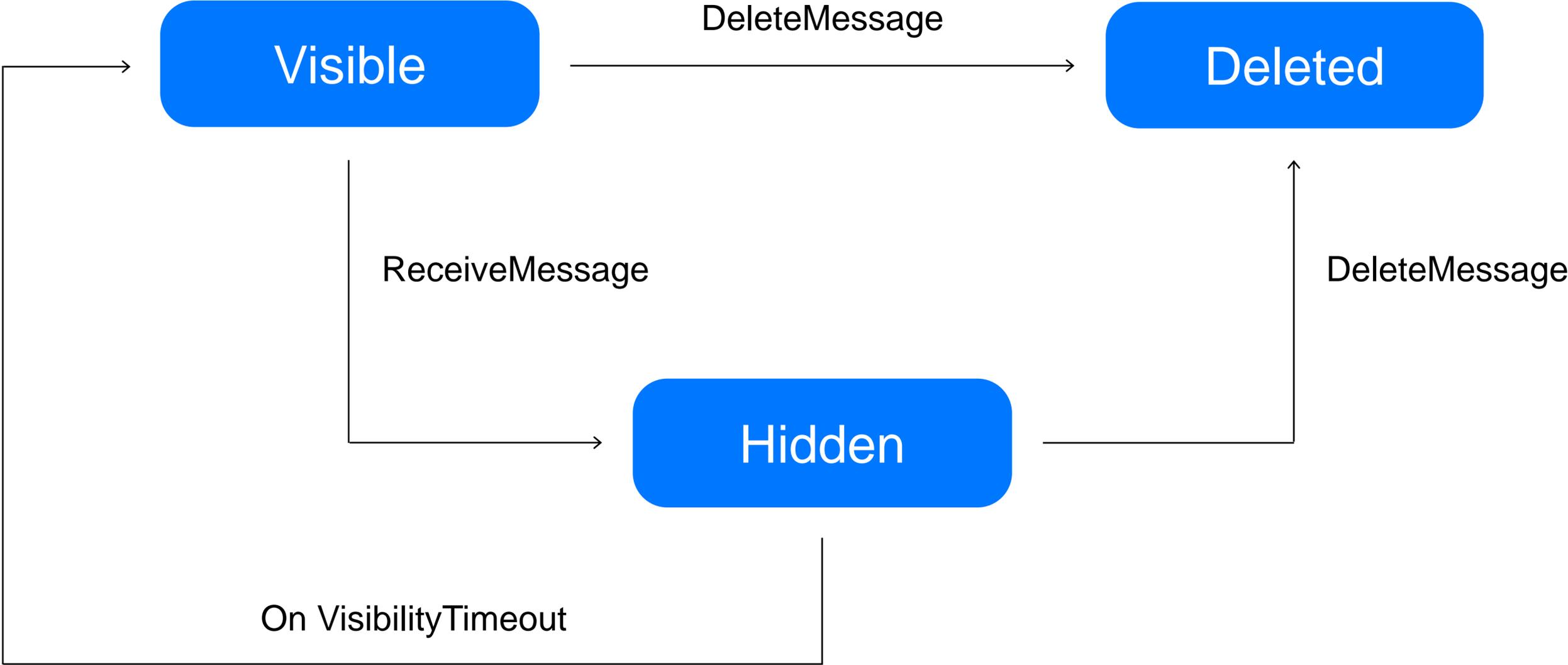
# Чтение сообщения



# Удаление сообщения



# Visibility timeout



05

**Внутреннее устройство YMQ**

# Платформа для Yandex Message Queue



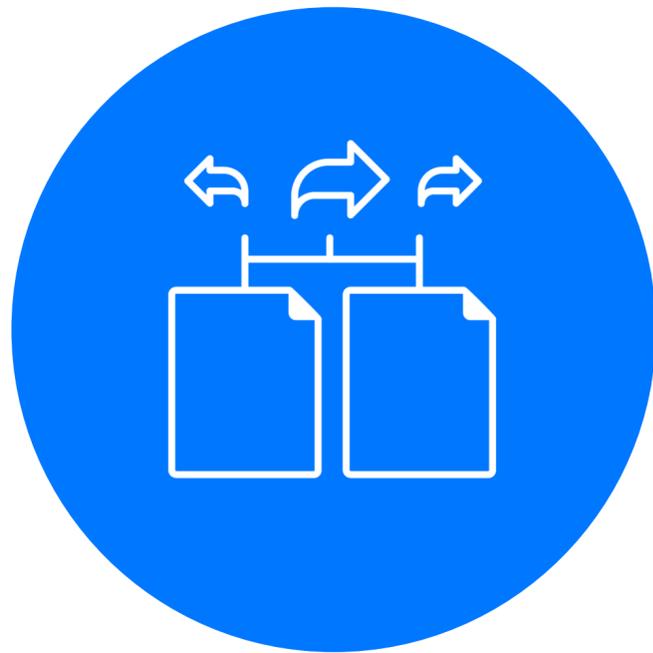
Основой YMQ является  
Yandex Database (YDB)

# Важные свойства Yandex Database



- › Распределённая
- › Горизонтально масштабируемая
- › Отказоустойчивая
- › Поддерживает ACID-транзакции
- › Поддерживает диалект SQL (YQL)
- › Строго консистентная
- › Очень крутая!

# Основные идеи



- › Очередь — набор шардов
- › Шард — набор таблиц в базе данных
- › Операции — запросы к таблицам шарда

# Анатомия очереди

## Queue inner table structure

Attributes

State

### Queue shard #1

Messages

MessagesData

Infly

Timestamps

### Queue shard #2

Messages

MessagesData

Infly

Timestamps

### Queue shard #3

Messages

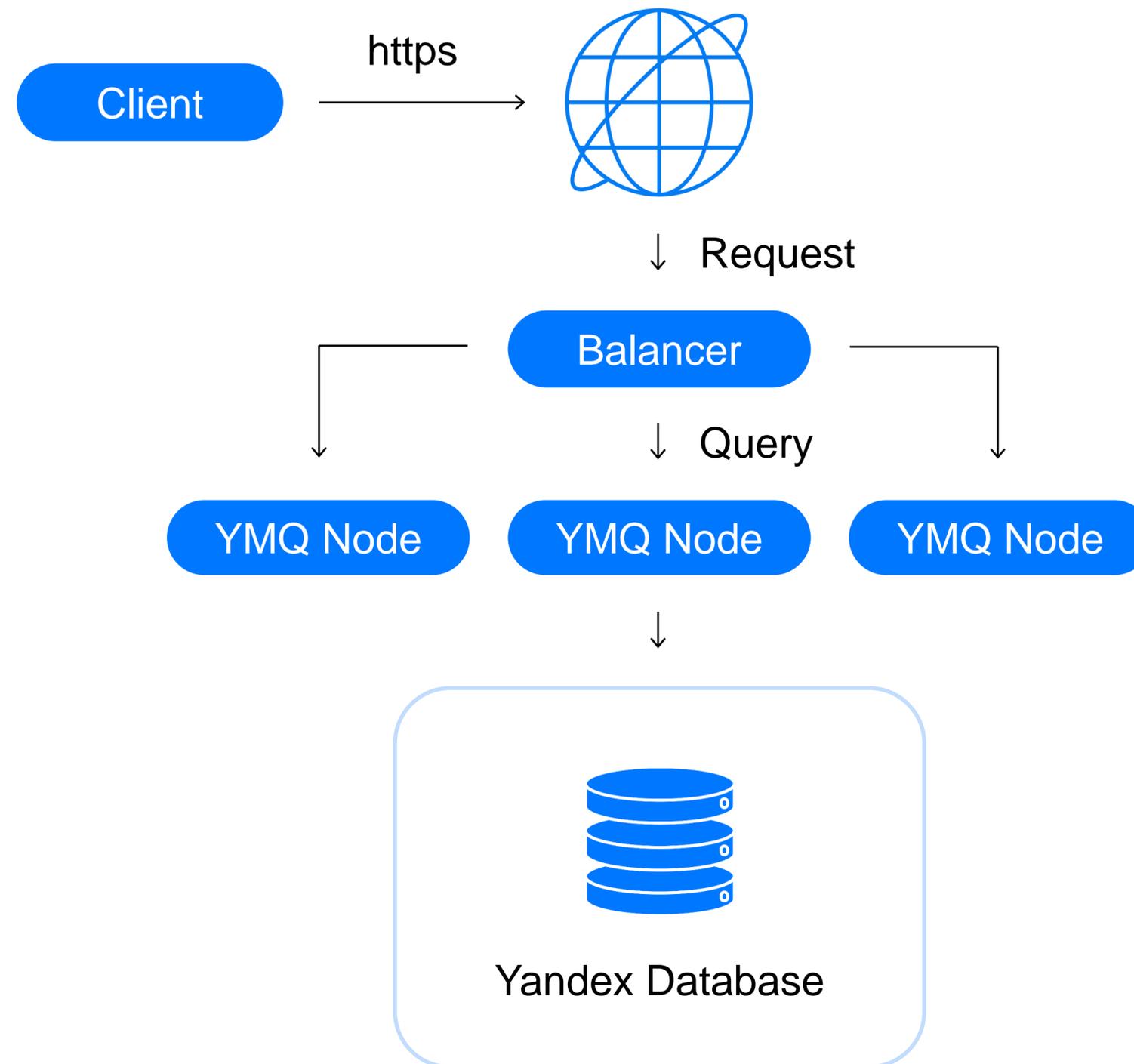
MessagesData

Infly

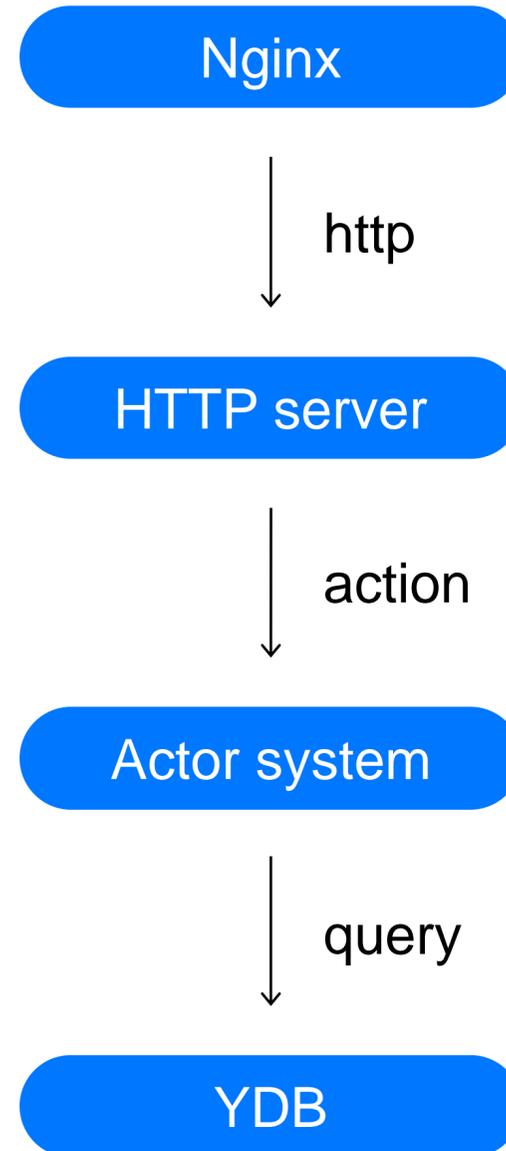
Timestamps



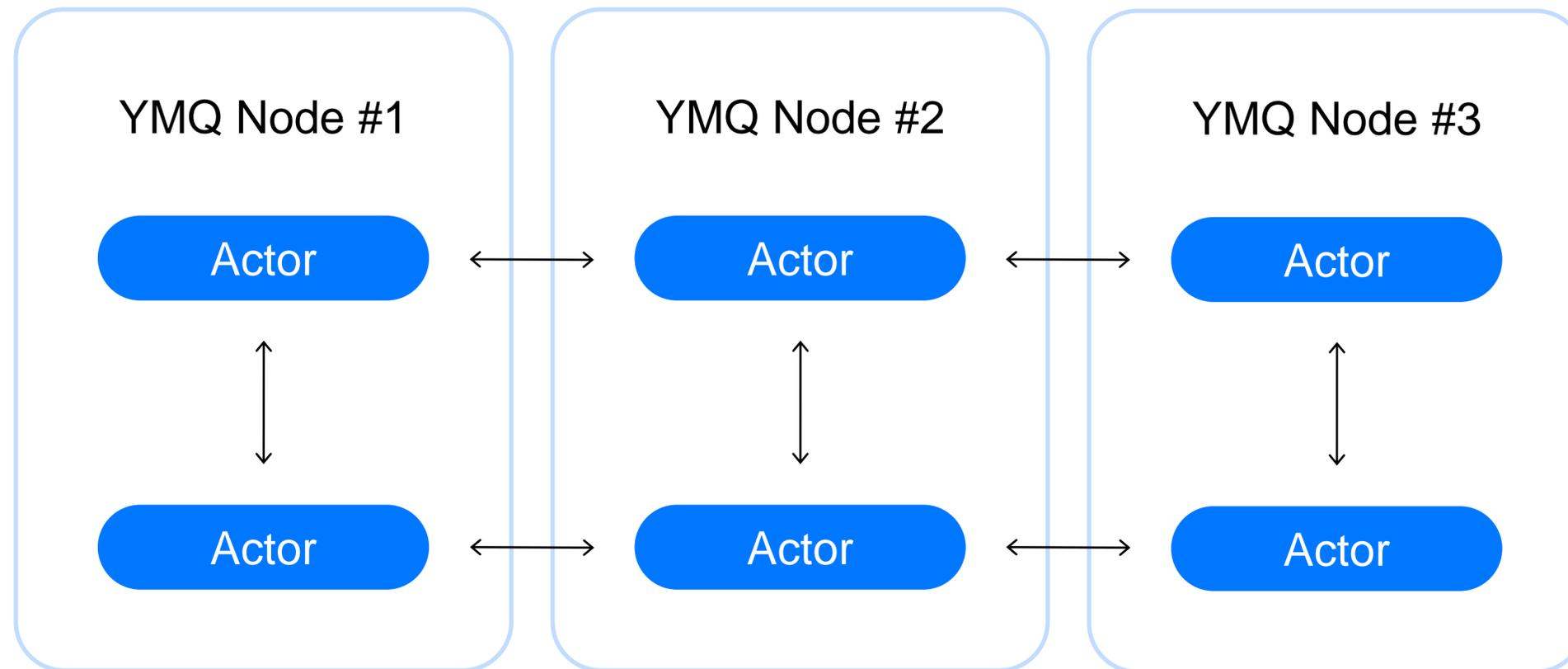
# Архитектура сервиса



# Внутри узла YMQ



# Акторная система

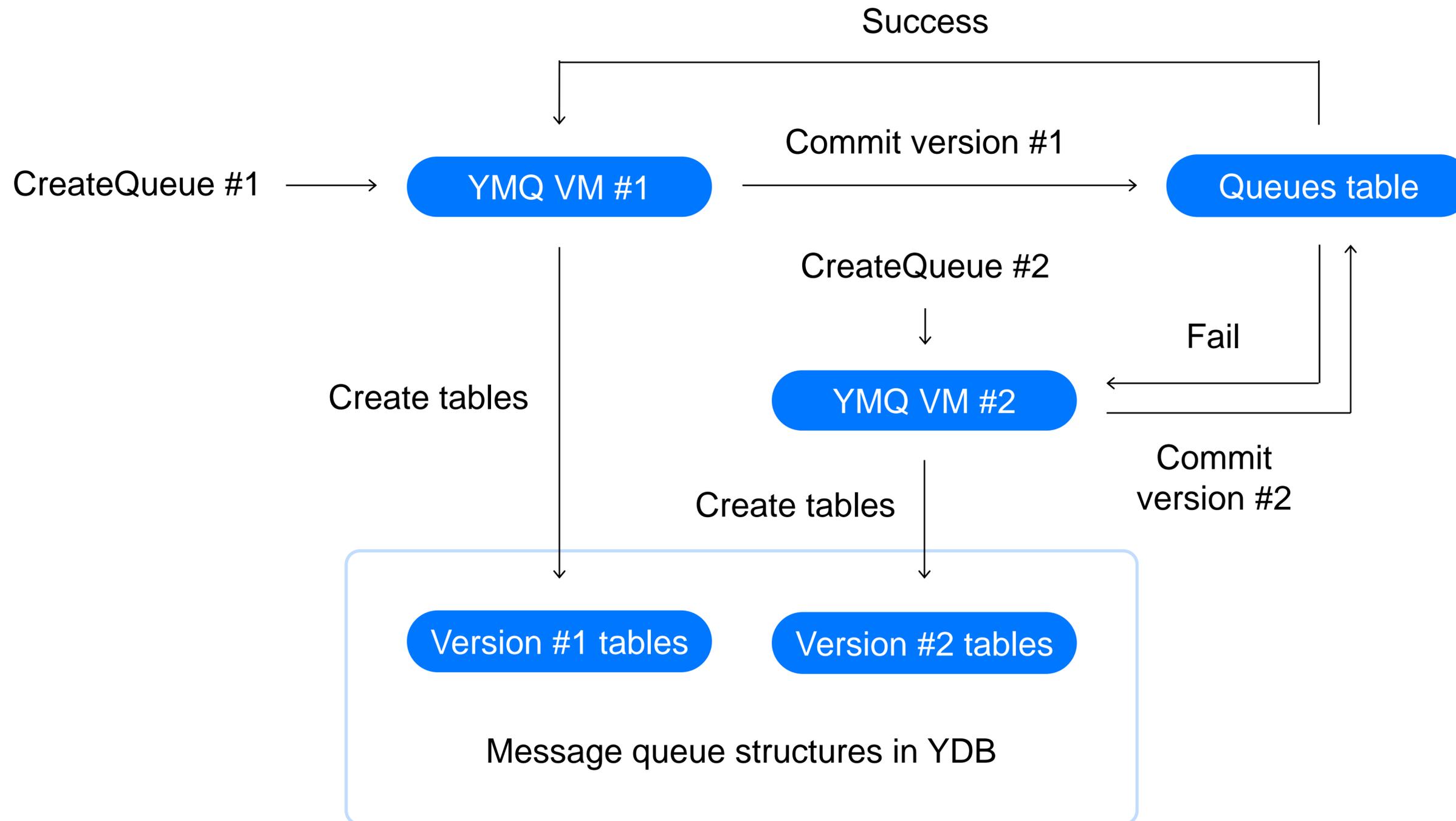


# Проблема создания очередей

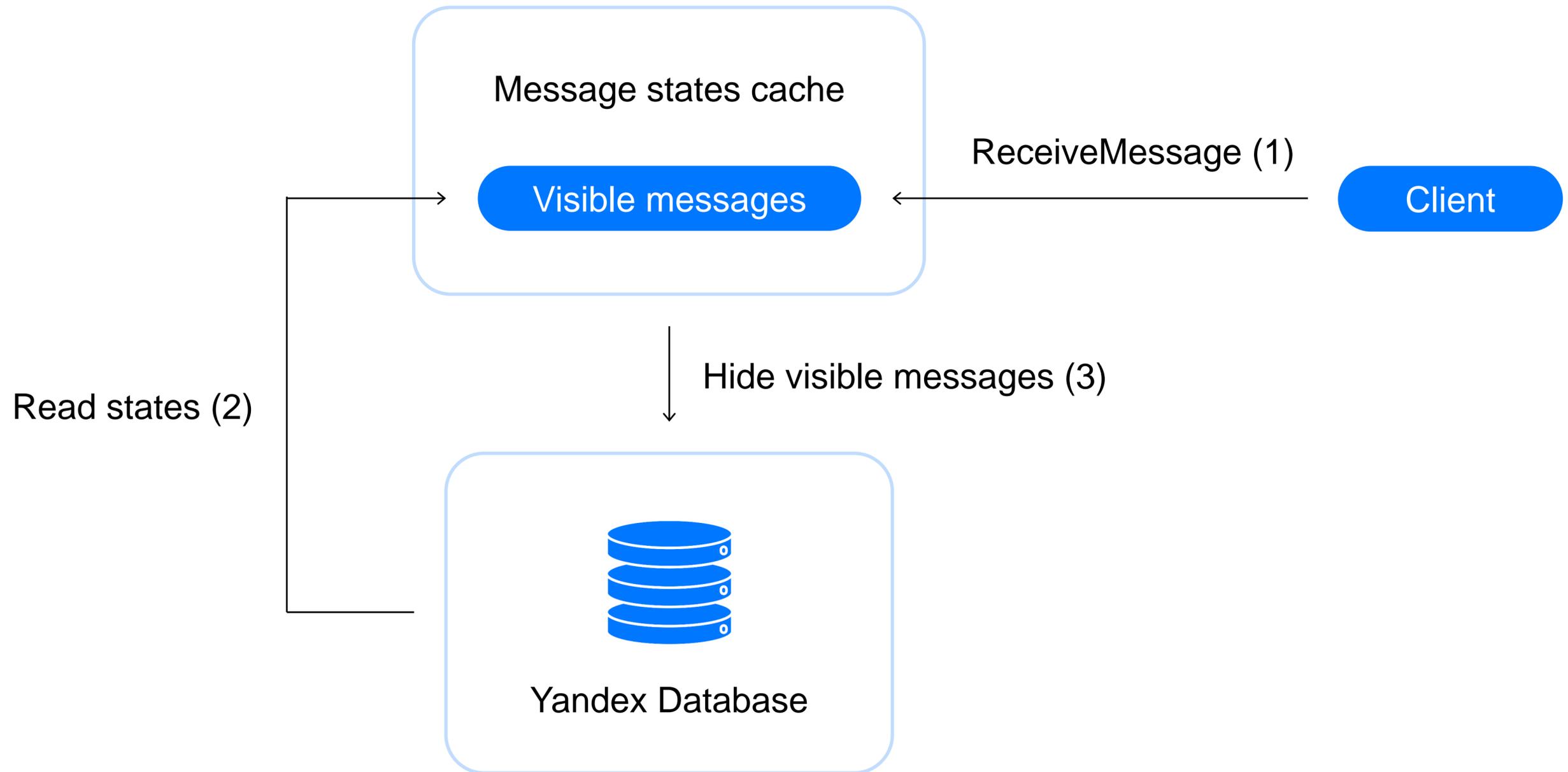


- › В YDB (пока что) нельзя транзакционно одновременно создать несколько таблиц и записать данные
- › Хост может упасть во время работы CreateQueue

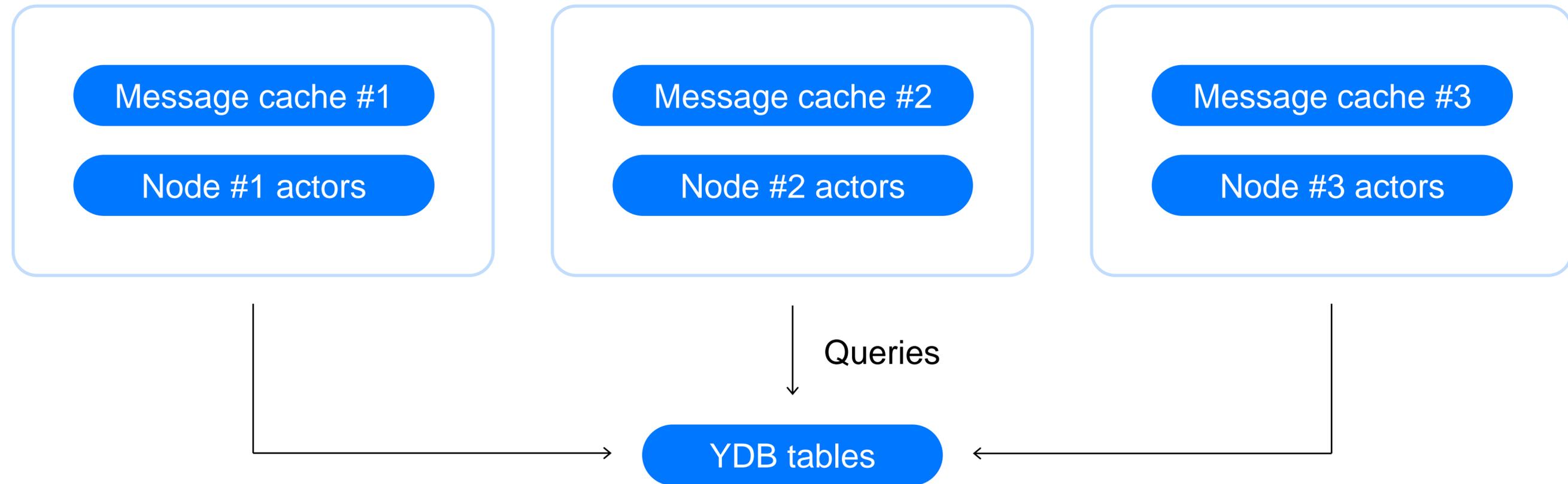
# Отказоустойчивое создание очередей



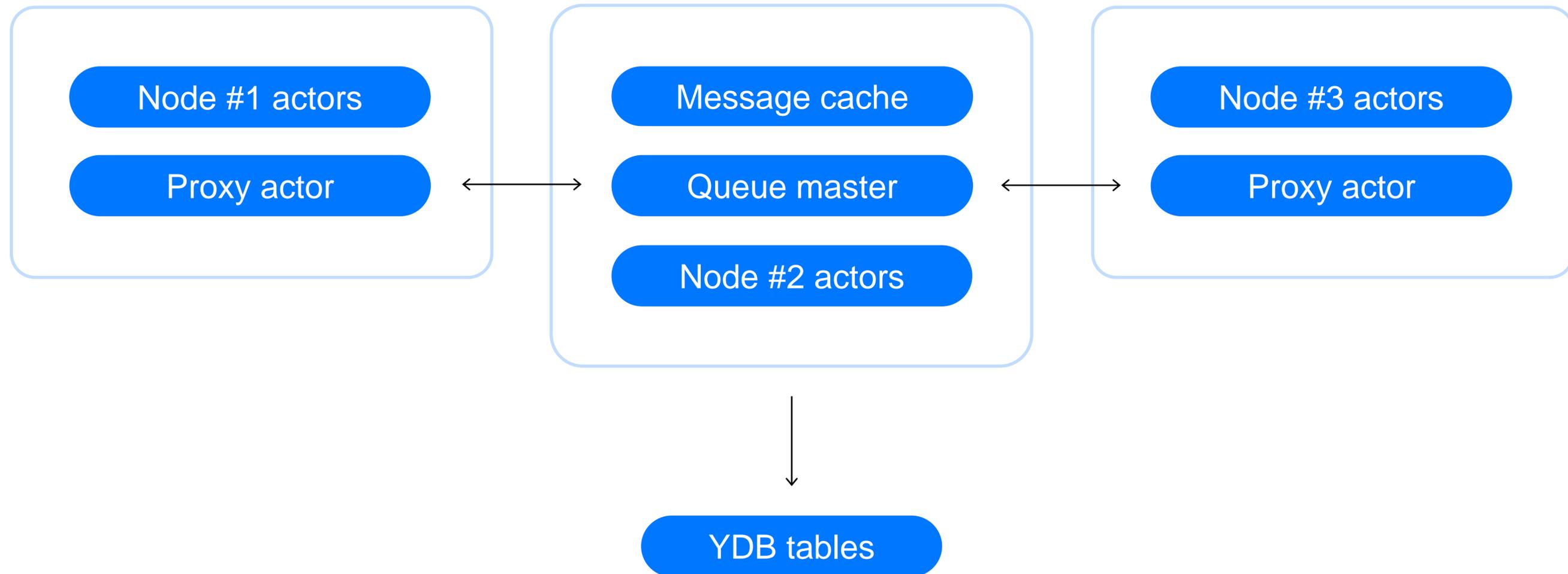
# Кэш сообщений



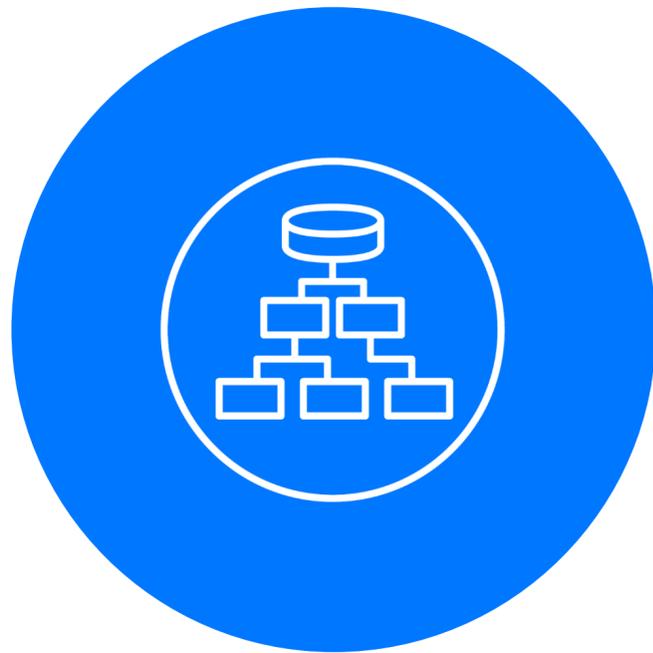
# Проблема согласованности кэшей сообщений



# Решение: мастера очередей



# Функции мастера очереди



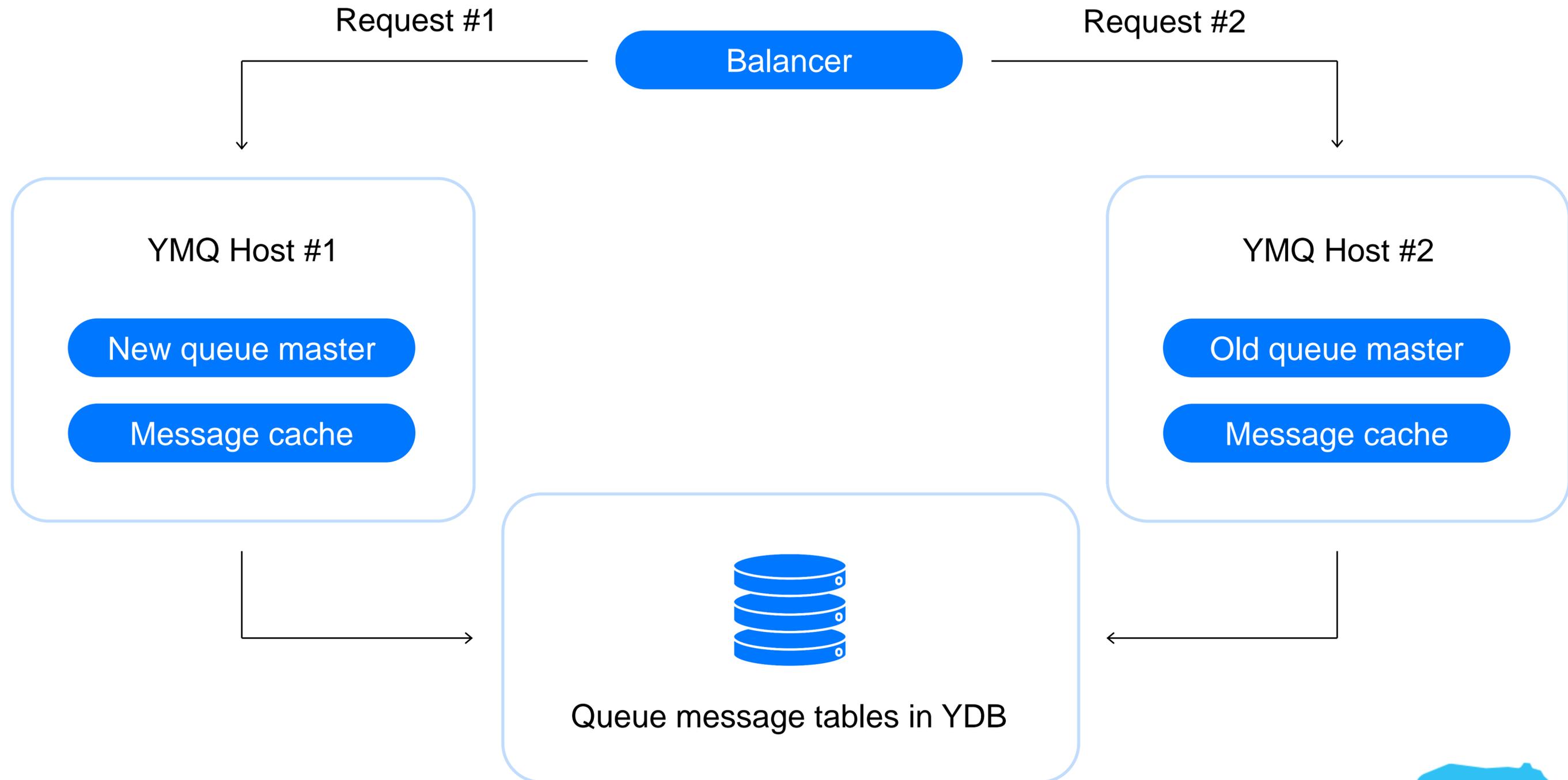
- › Кэширование информации
- › Единая точка управления запросами
- › Место сбора пользовательских метрик
- › Оптимизация потока запросов к базе

# Выбор мастера

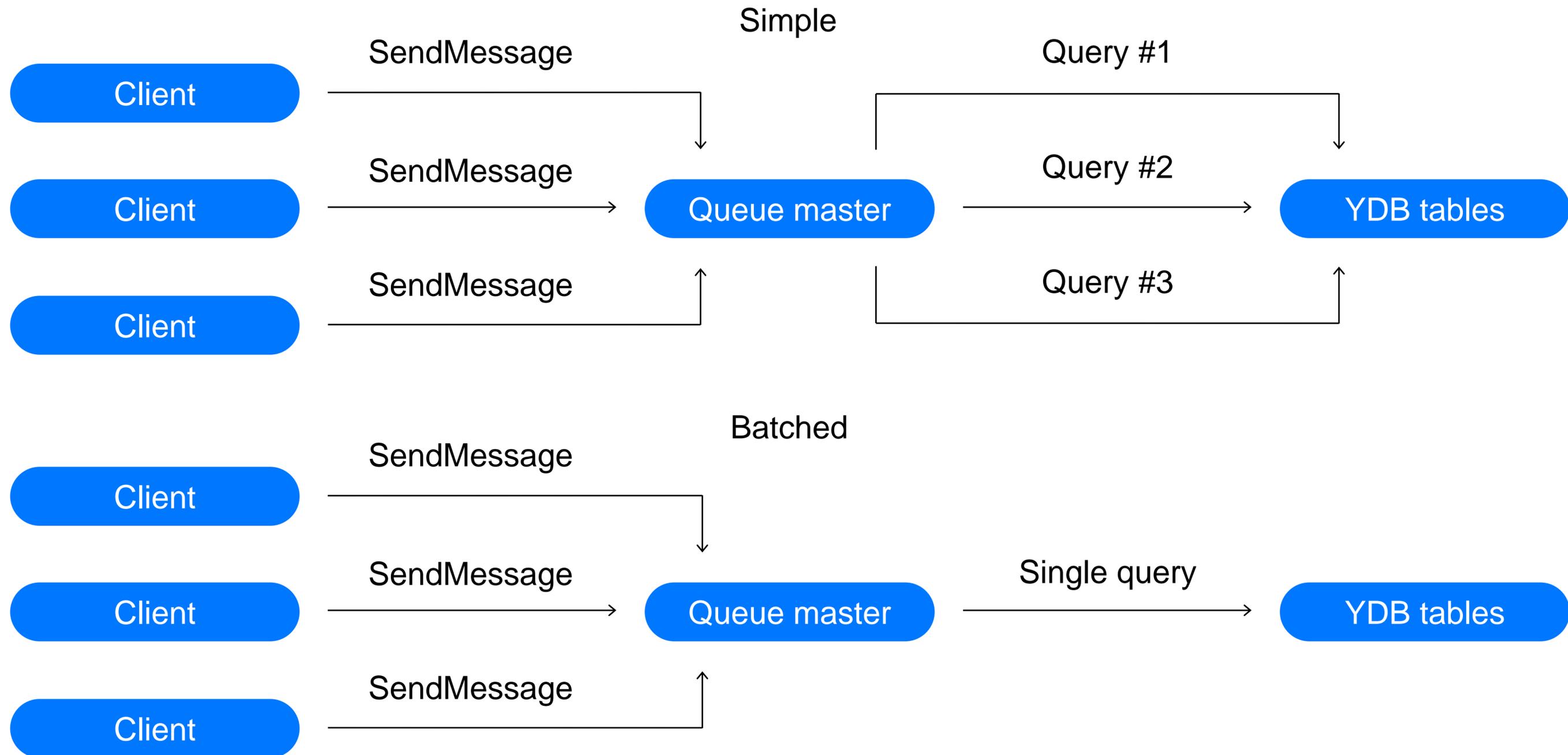


- › Выбирается за счет механизмов YDB
- › Мастера быстро переедут при отказе узла

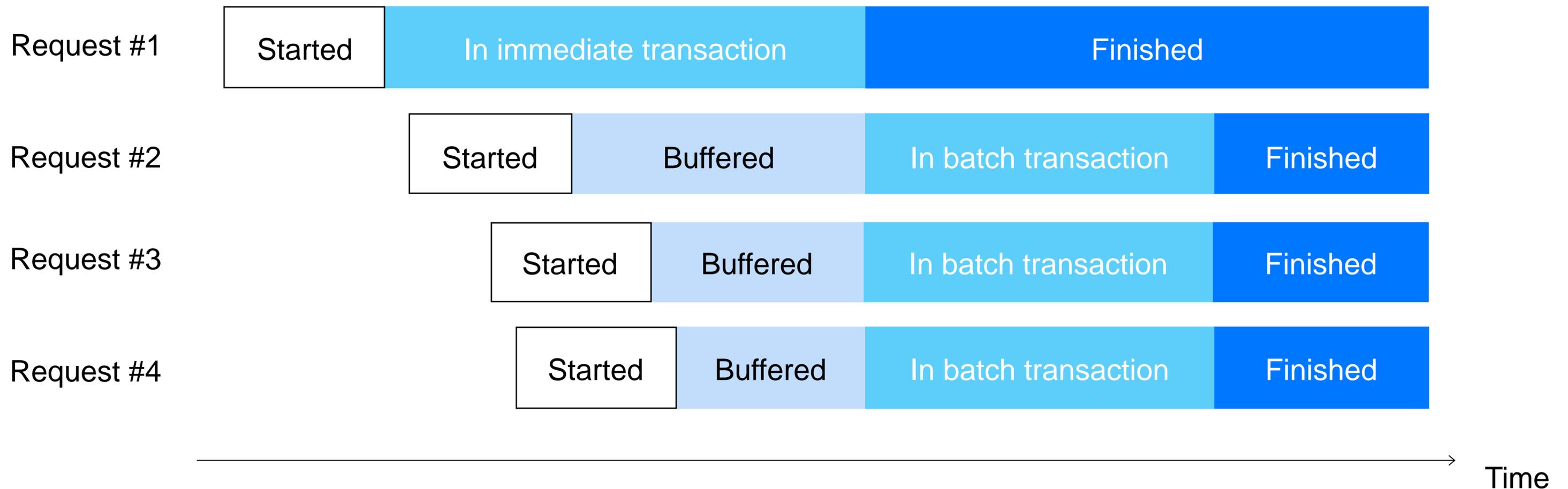
# Когда мастеров в кластере несколько



# Группировка однотипных запросов



# Временная диаграмма батчинга



# Что мы узнали сегодня



- › Очередь сообщений — полезный компонент
- › Exactly-once недостижима только за счет очереди
- › Apache Kafka и RabbitMQ — два различных подхода
- › Иногда делать своё решение выгоднее
- › YMQ построен на мощной платформе YDB
- › Распределённость привносит свои сложности



# Спасибо за внимание!

Василий Богонатов

Разработчик сервиса Yandex Message Queue

 [radix@yandex-team.ru](mailto:radix@yandex-team.ru)

 [@drween](https://t.me/drween)



**HighLoad<sup>++</sup>**  
Siberia 2019