

# Кafka и высоконагруженные системы

Роговский Денис



# План доклада

1. Что есть нагруженная система
2. Системы очередей сообщений
3. Выбор нашей компании
4. Кеш и его реализации

# Как очереди вообще касаются дата инженера

- Необходимость обработки большого количества сообщений от разных производителей данных
- Разделение обработки между серверами
- Отправка результатов разным потребителям

# Как очереди вообще касаются дата инженера. Цифры

Ежедневно приходит более 2 миллиарда запросов с 10+ серверов

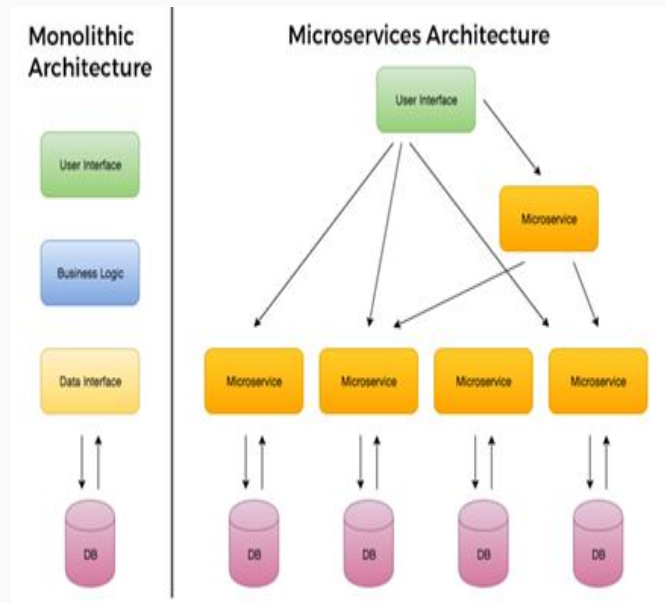
7 сервисов только от отдела машинного обучения

Несколько источников для каждого сервиса

Несколько потребителей

# Высоконагруженные системы

- Множество независимых сервисов
- Асинхронная обработка запросов (сообщений)
- Общение сервисов между собой

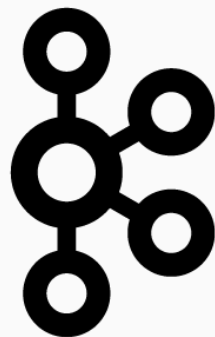


# Варианты Message Queue

- RabbitMQ
- Kafka
- RocketMQ
- Artemis
- ....



**RocketMQ**



**kafka**

# Kafka vs RabbitMQ

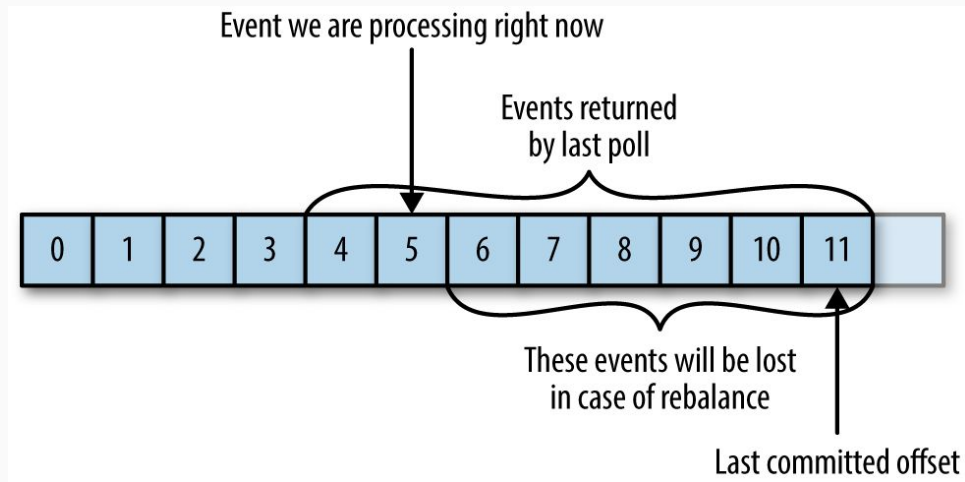
RabbitMQ - следит за сообщениями, при обработке сообщения, удаляет его, если не обработали, возвращает в очередь. Тратятся ресурсы сервера на отслеживание состояния сообщения, повторное прочтение невозможно.



# Kafka vs RabbitMQ

Кafka - записывает сообщения на диск в определенной последовательности, сообщение удаляется либо по времени, либо при достижении объема памяти.

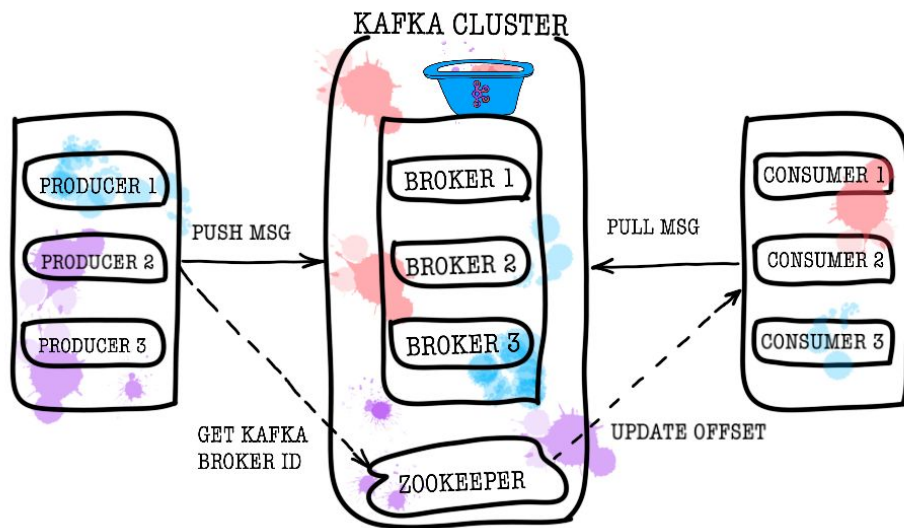
-





# Особенности кафки

- Обработка разной группой потребителей одного сообщения
- Упорядоченность очереди в рамках части топика
- При добавлении потребителя ребалансит нагрузку
- Могут появиться дубли!!



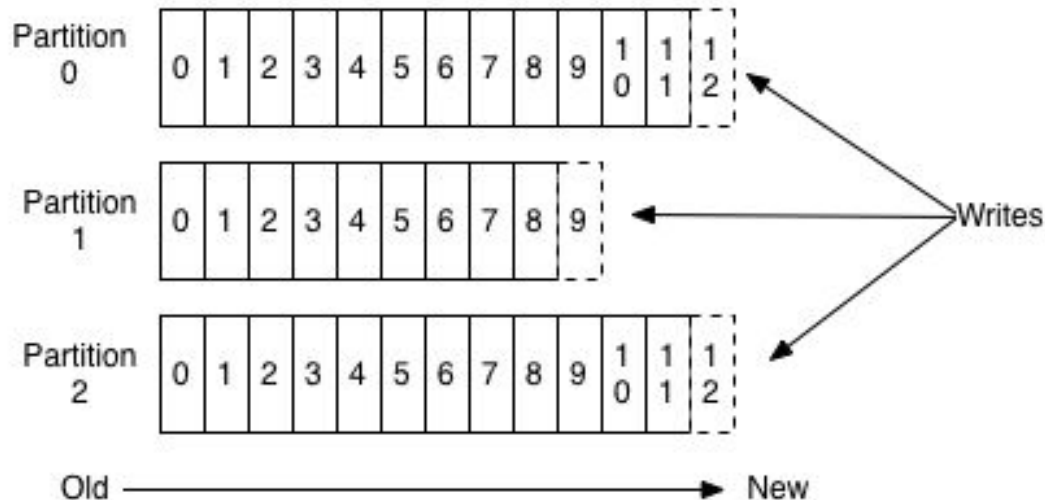
# Реализация топика кафки

В одном топике может существовать множество партиций.

При добавлении нового потребителя партиции ребалансируются.

Появляются дубли при чтении с тех же отступов от верха очереди.

## Anatomy of a Topic



# Кэш

Кэш - промежуточный буфер с быстрым доступом к нему, содержащий информацию, которая может быть запрошена с наибольшей вероятностью.  
(с) Джейсон Стедхем



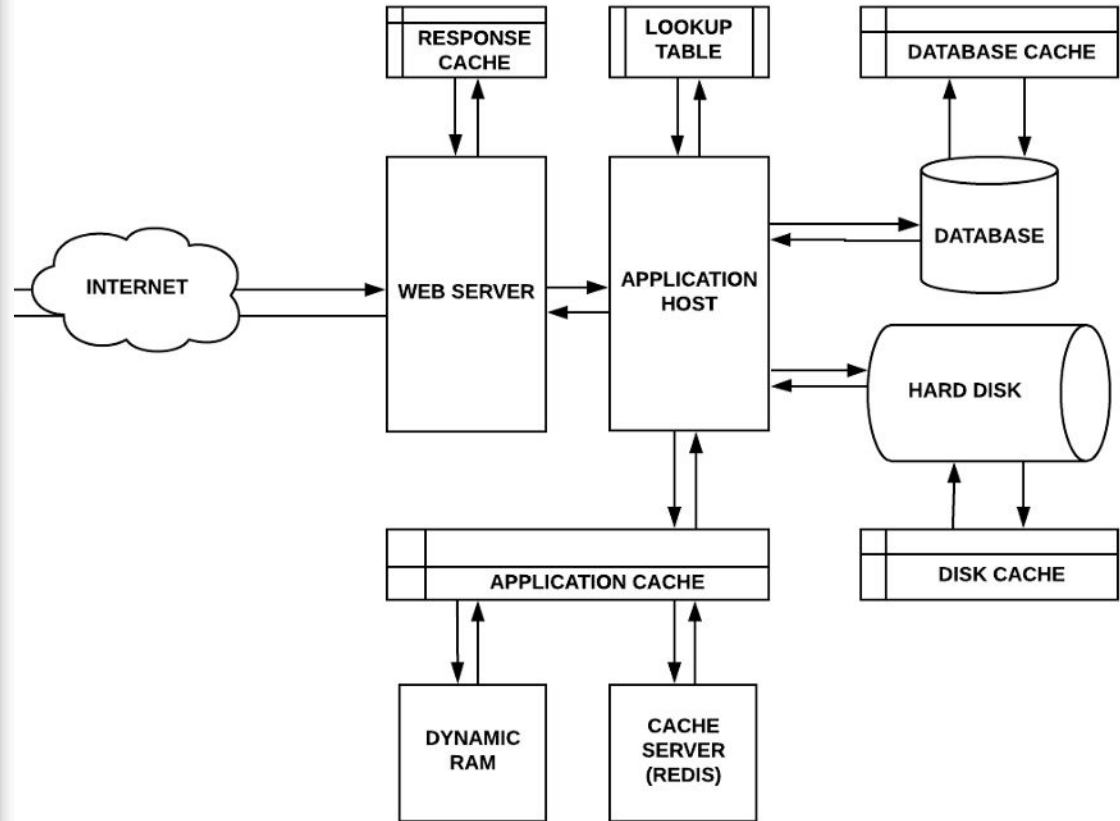
# Кэш

Под каждую задачу  
своя реализация кеша

Правильный объем  
данных в кэше

Кешировать как можно  
позже

Как реализовать



# Все это в разрезе работы отдела МЛ

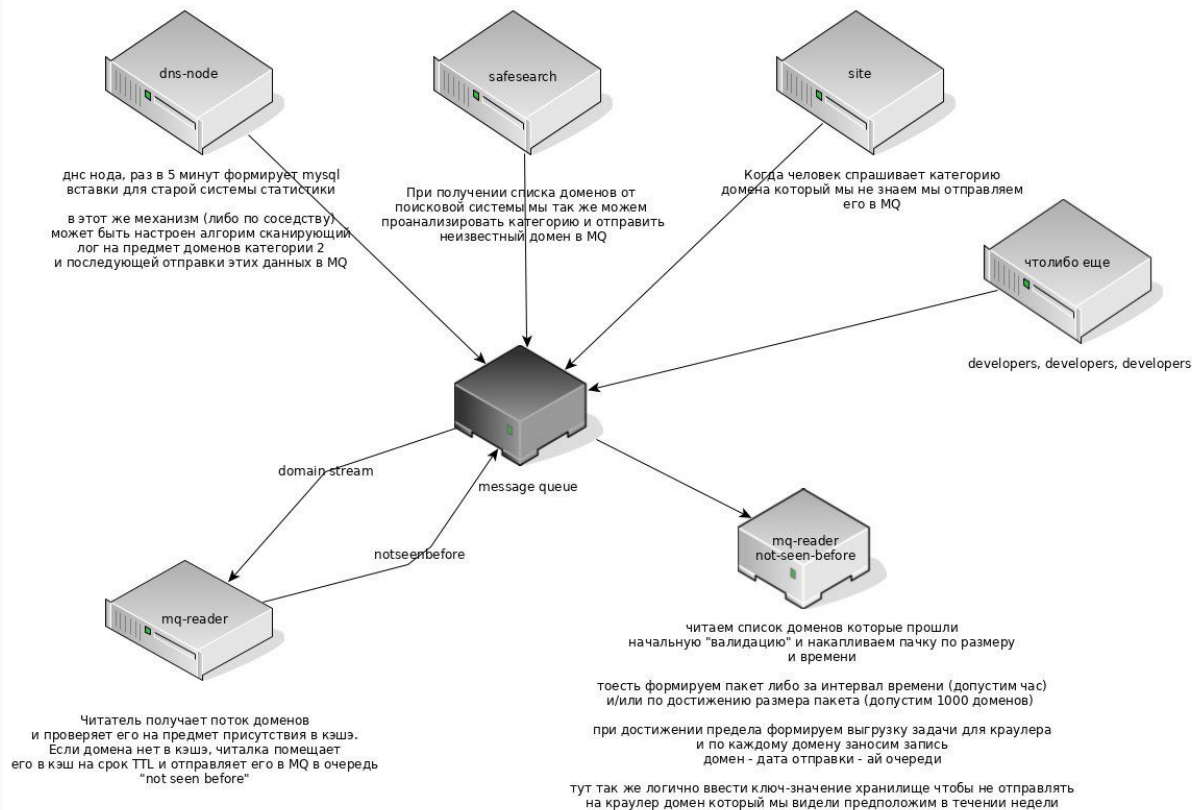
Есть новые домены

Надо собрать контент

Сбор контента - дорого

Вычитать недавно  
виденные

Отправлять собранный  
контент на  
классификацию



# Рализация кэша

Два уровня

Цикл разработки

Цикл клиентов

Реализация на инмемори базе данных.



# Увеличение эффективности

При использовании двухуровневого кеша

после первого уровня режется 30% входящих данных,

После второго уровня - до 50% первоначального значения.

До 10-15 млн неизвестных доменов из 25-30

Смогли избежать дублей, и сильной загрузки серверов.

# Спасибо за внимание

Надеюсь этот небольшой  
обзор был вам полезен

