

# Перформанс метрики



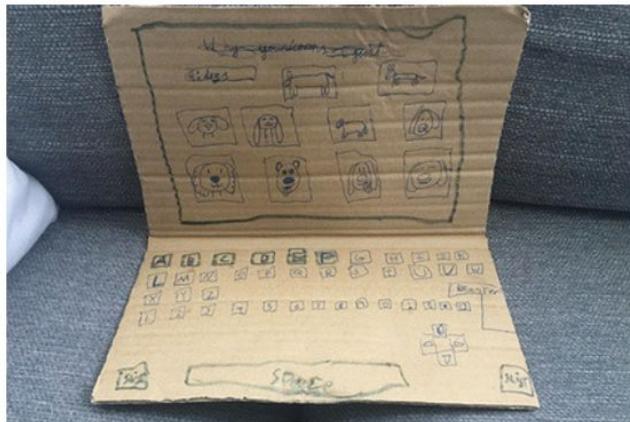
# Что сегодня будет

Расскажу про

- Клиент-серверная система
- latency/throughput/bandwidth
- median, p75, p90
- Браузерный перформанс
- Как на все это смотрит пользователь



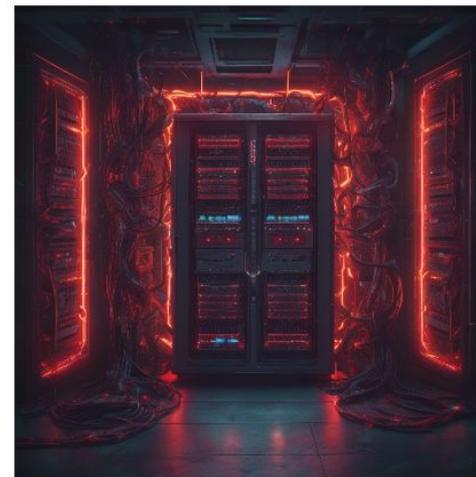
# Клиент-Серверная инфраструктура



Клиентская машина



Наша великая инфраструктура



Наш мощный сервер

# Наши приложения - "Data Intensive"



Клиентская Машина 1



Клиентская Машина 2



Новый провайдер



Новый Мощный сервер



Новая Великая сеть датацентра



Мощный сервер БД



Мощный сервер Чего То еще

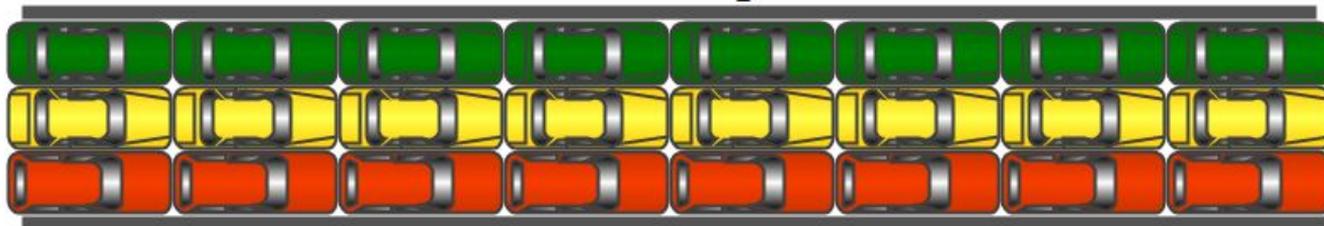
Наши сервисы большую часть времени проводят в ожидании ответа от других сервисов

# Великолепная тройка

Ширина канала

Bandwidth

24 Cars per second

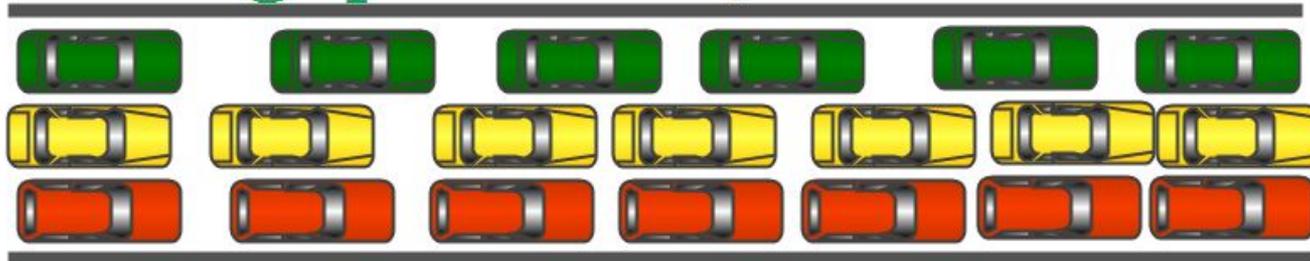


= RPS

Пропускная способность

Throughput

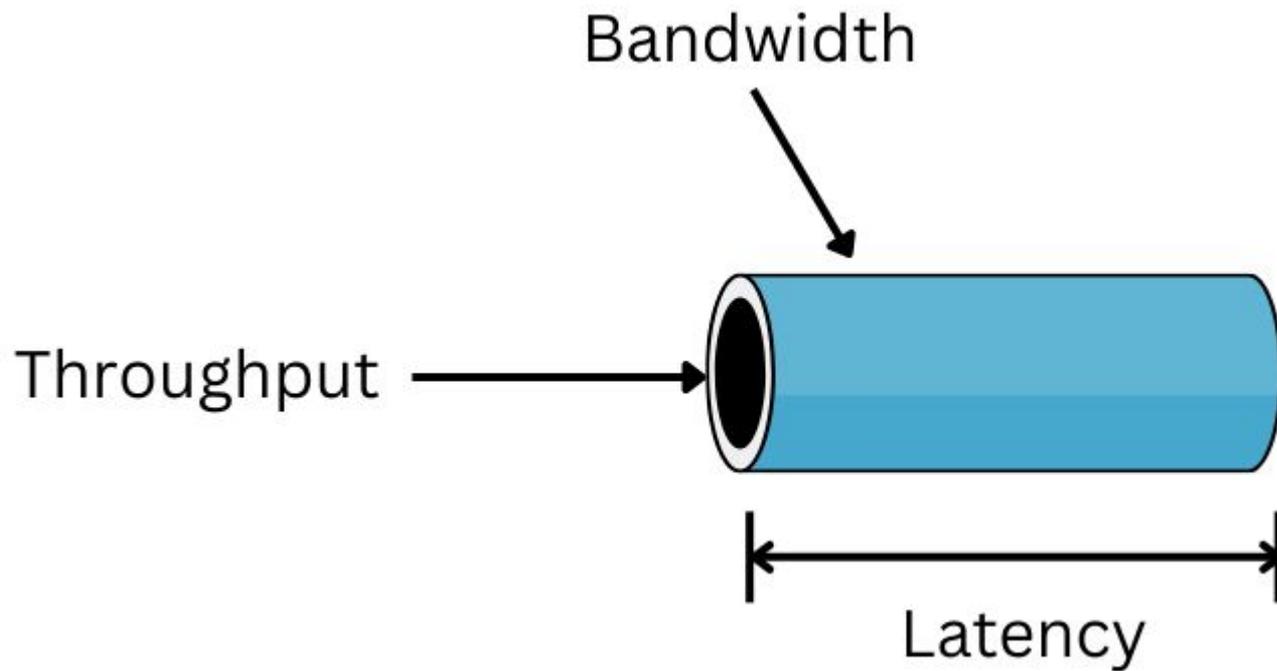
20 Cars per second

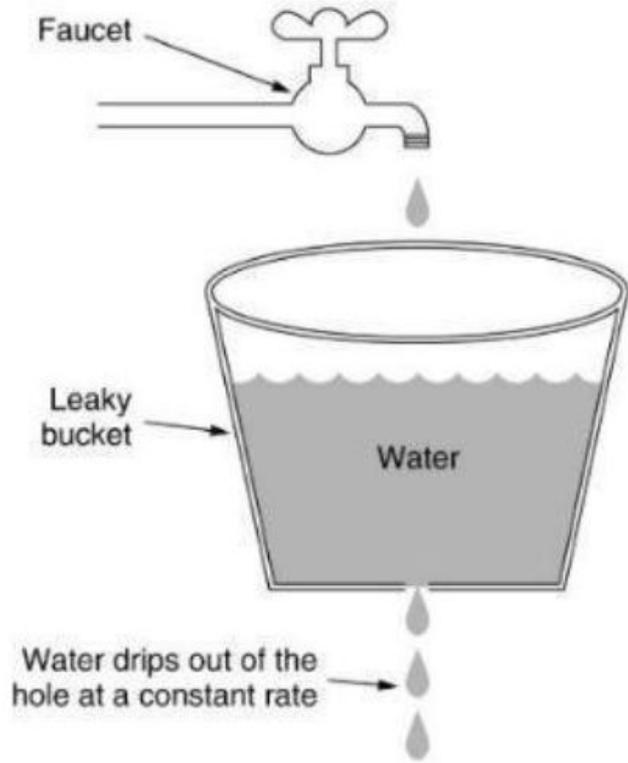


+

Latency - Длина "пути"

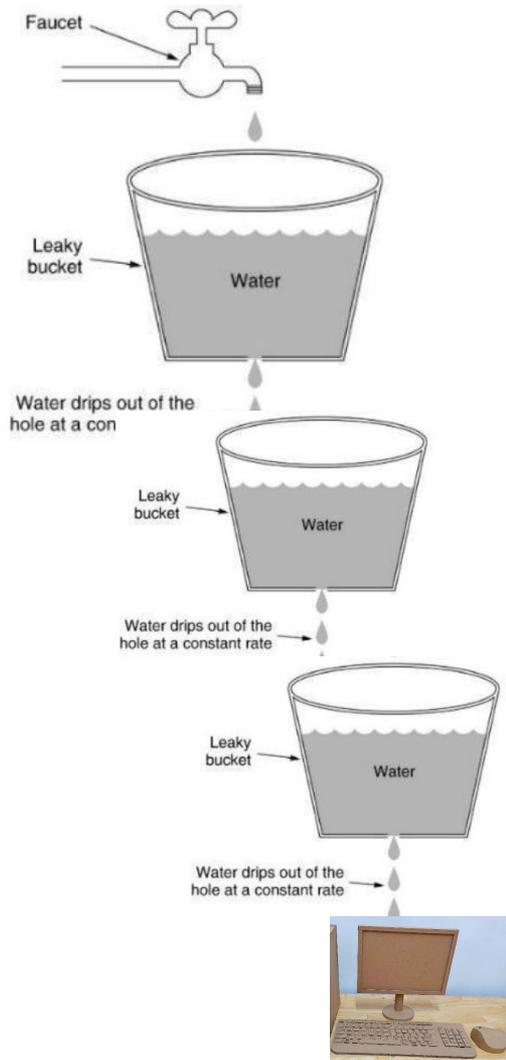
# Аналогия с шлангом





Входной поток  $\neq$  Скорость  
выполнения запроса

Вывод – Где то есть очередь  
обработки



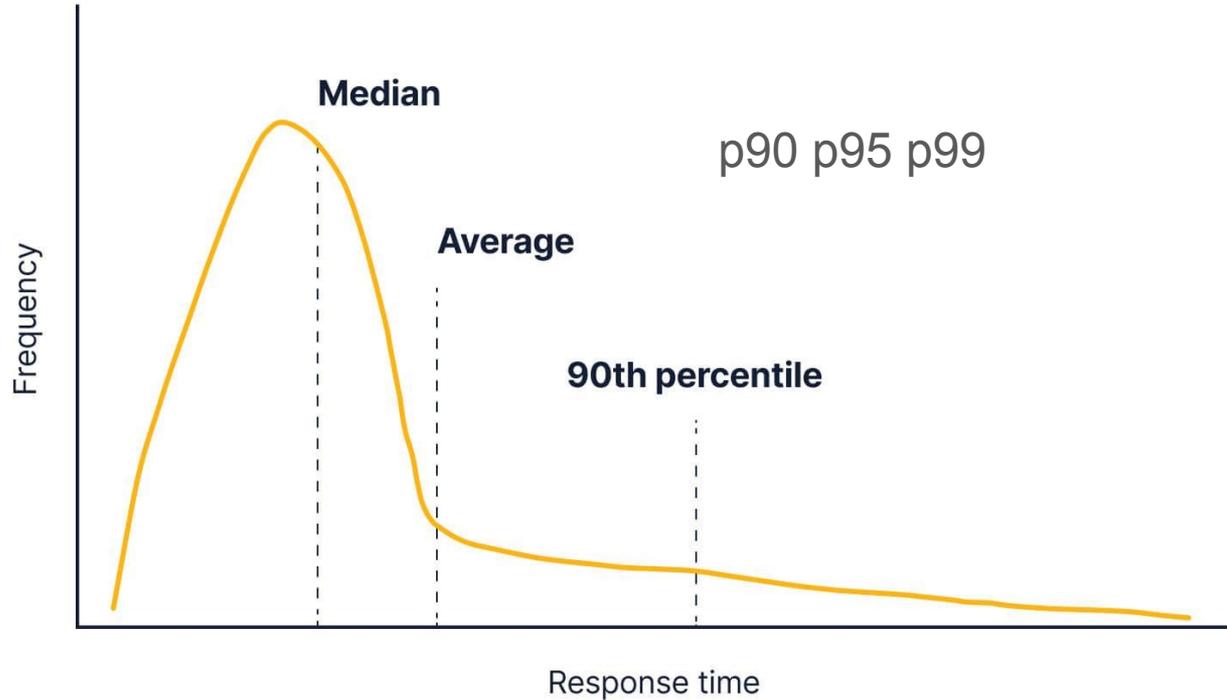
# На практике

Очередей целая цепочка

Некоторые “бочки” - самые **медленные**  
(бутылочные горлышки)

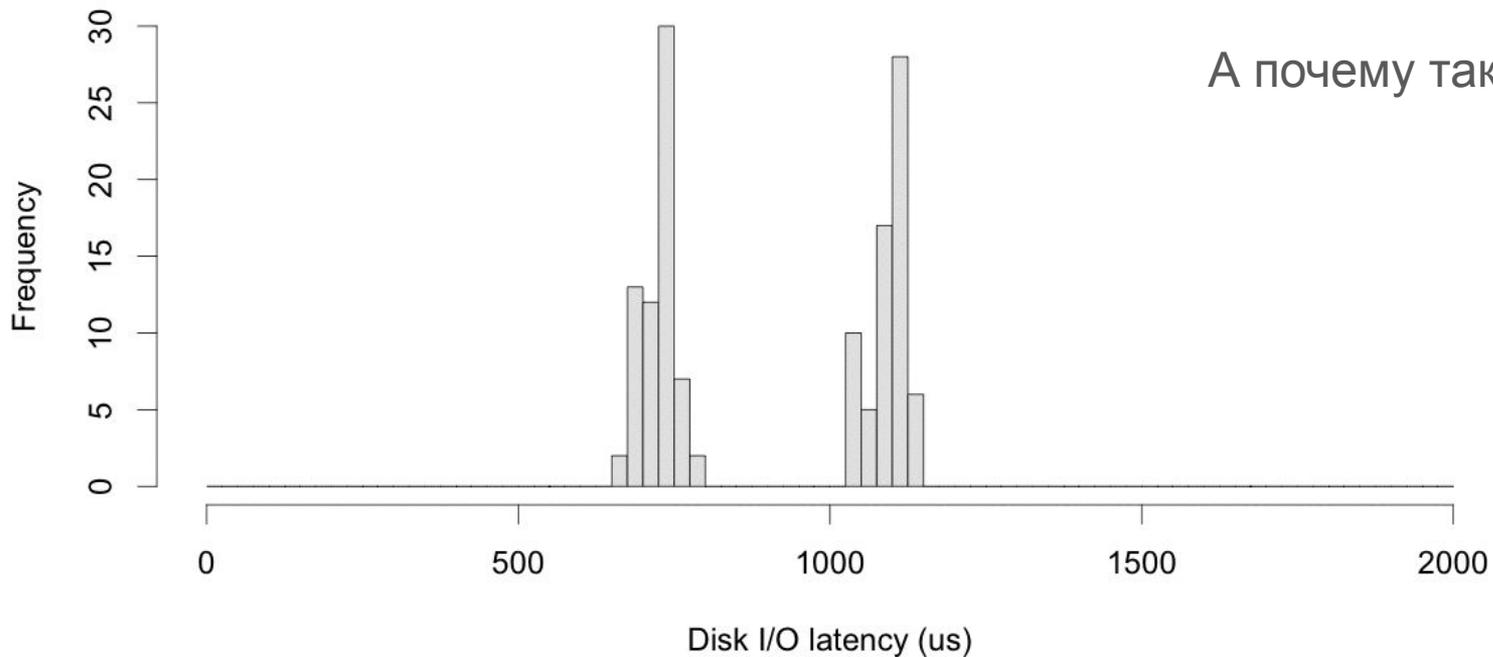
Обработка запросов моделируется хорошо  
системой очередей

Среднее никому не интересно, всем интересен хвост



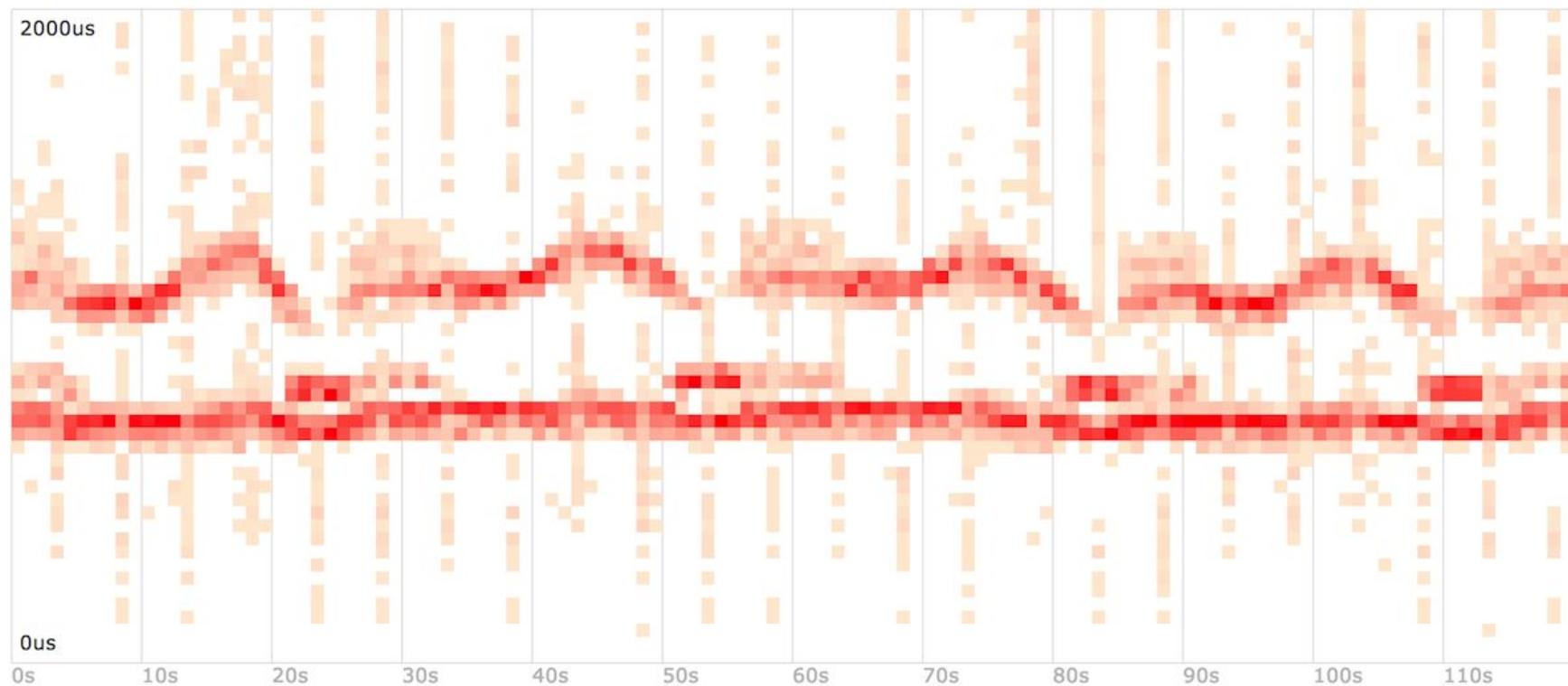
# Latency на практике

Latency Distribution, per second



# Тепловая карта

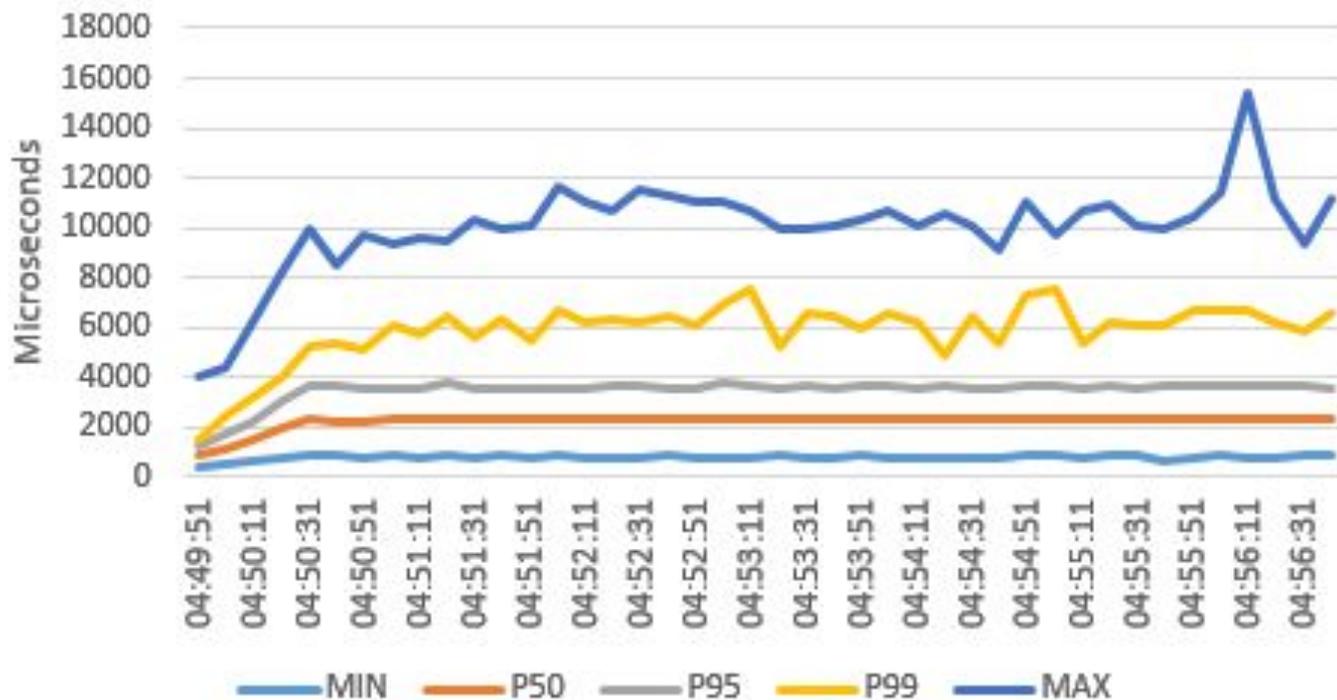
Latency Heat Map



time 54s, range 840-880us, count: 34

Time

# Latency на графиках



# Конфликт требований к ОПТИМИЗАЦИИ

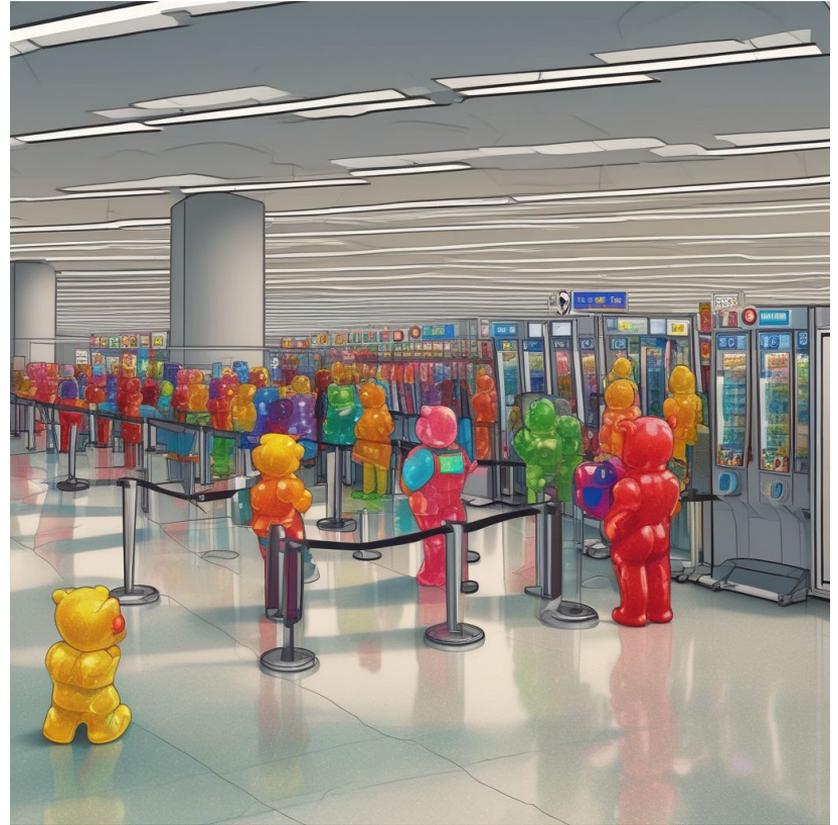
- Для того, чтобы **снизить latency**, необходимо **уменьшить среднюю длину очереди (очередей)** в системе
- Для того, чтобы **увеличить RPS (throughput)**, необходимо **увеличить среднюю длину очередей в системе**, чтобы избежать простоев на фоне неравномерной нагрузки
- Для равномерного распределения нагрузки в распределённой системе, необходимо, чтобы очереди к разным узлам были примерно одного размера

# Под нагрузкой все становится медленным



# Как сделать быстрее

- Ищи **узкое** горлышко
- Под **нагрузкой** в какой то момент все становится сильно хуже, делай **нагрузочные**
- Не используй **среднее**
- Смотри на распределение
- **Оптимизировать** под latency и throughput это разные вещи
- **Умей** бенчить



# Графана, и что в ней встретишь

- метрики **сервисов**,
- метрики **серверов**
- иногда **бизнес метрики**

## *Полезное*

- Latency (p95)
- RPS = *Throughput*
- **Thread pool size**
- CPU load = *Bandwidth*
- *Длина/Лаг очереди*



# Типичный график дневной нагрузки

- Когда меряешь перформанс нужно учитывать что люди не всегда на **работе**

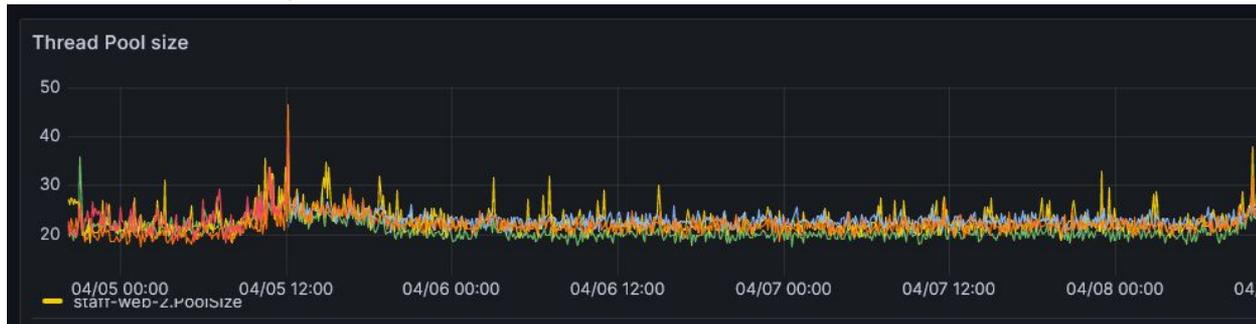


- Когда сервис не нагружен, перформанс оценить **нереально**
- Лучше всего, сравнивать с **прошлым** аналогичным периодом

# Thread Pool size

- Очень часто сервис ждет ответа от кого то еще

(например от базы данных)



Если сервис написан криво поток запроса лочится и ждет ответа

Ради чего придумали (async/await) - dotnet поднимает доп потоки если все чем то заняты (или ждут)

- На графике можно увидеть проблему с блокировками, с кривым async-await

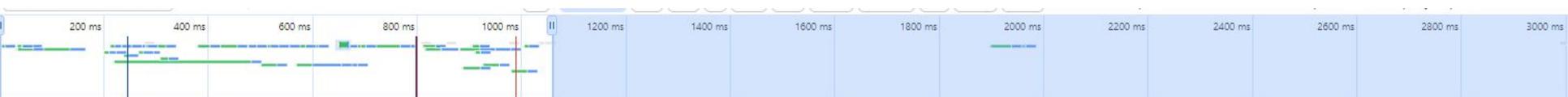
# GC Time

- Важно для сервисов где есть уборка мусора
- Требуется отдельной заботы

# Браузерный перформанс

**Кривой** фронт может похоронить все оптимизации на сервере и в инфраструктуре





Типичный кривой фронт

Или, зачем я оптимизировал бек?

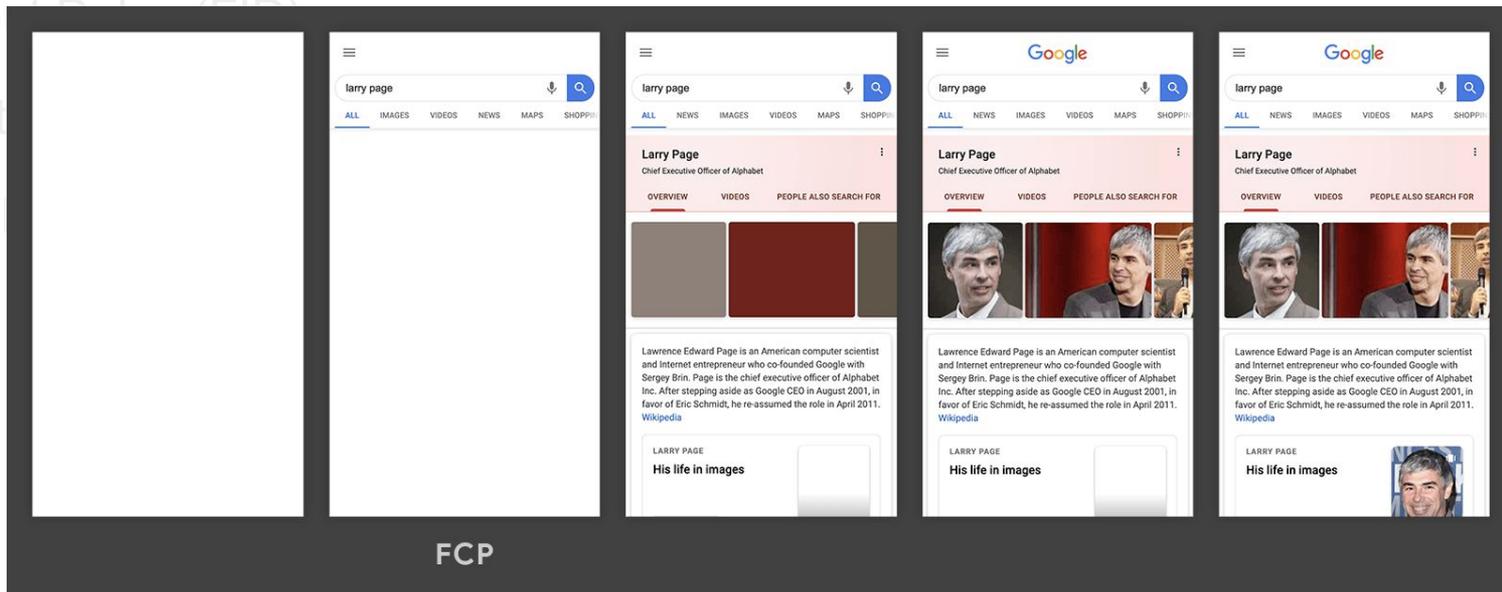
# Как измеряют перформанс в браузерах

## First Contentful Paint (FCP)

First Input Delay (FID)

Cumulative

Time to



# Как измеряют перформанс в браузерах

First Contentful Paint (LCP)

First Input Delay (FID)

Время реакций на ввод, инпут лаг

Cumulative Layout Shift (CLS)

Time to First Byte (TTFB)

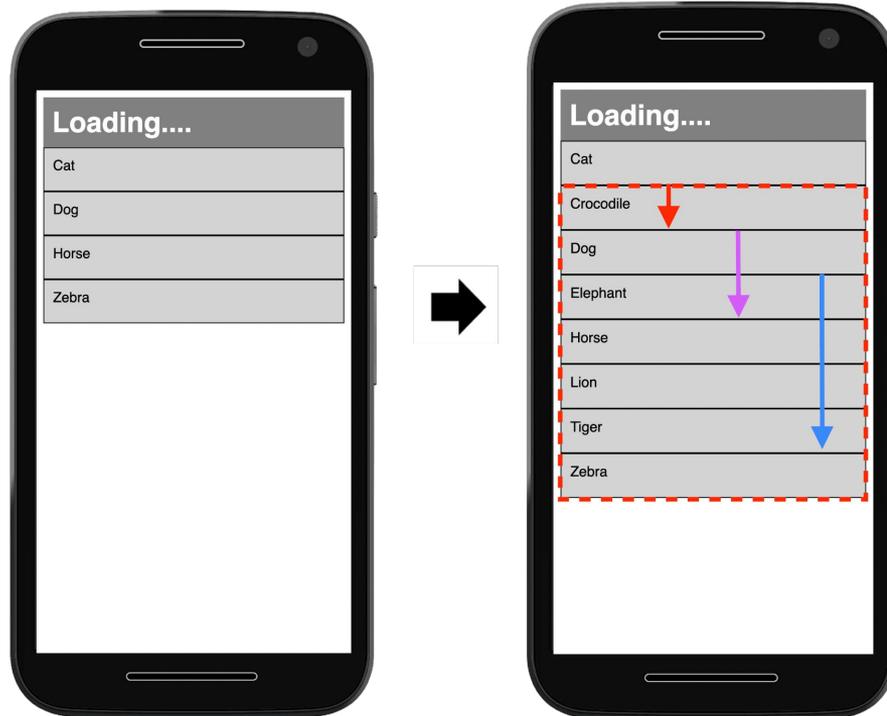
# Как измеряют перформанс в браузерах

First Contentful Paint (LCP)

First Input Delay (FID)

Cumulative Layout Shift (CLS)

Time to First Byte (TTFB)



Когда сайт грузится, интерфейс скачет

# Как измеряют перформанс в браузерах

First Contentful Pain (LCP)

First Input Delay (FID)

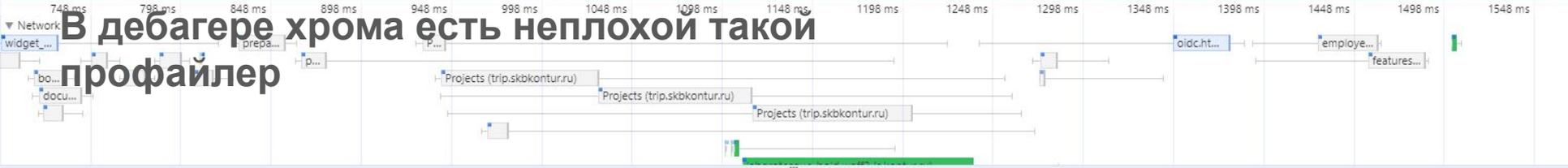
Cumulative Layout Shift (CLS)

Time to First Byte (TTFB)

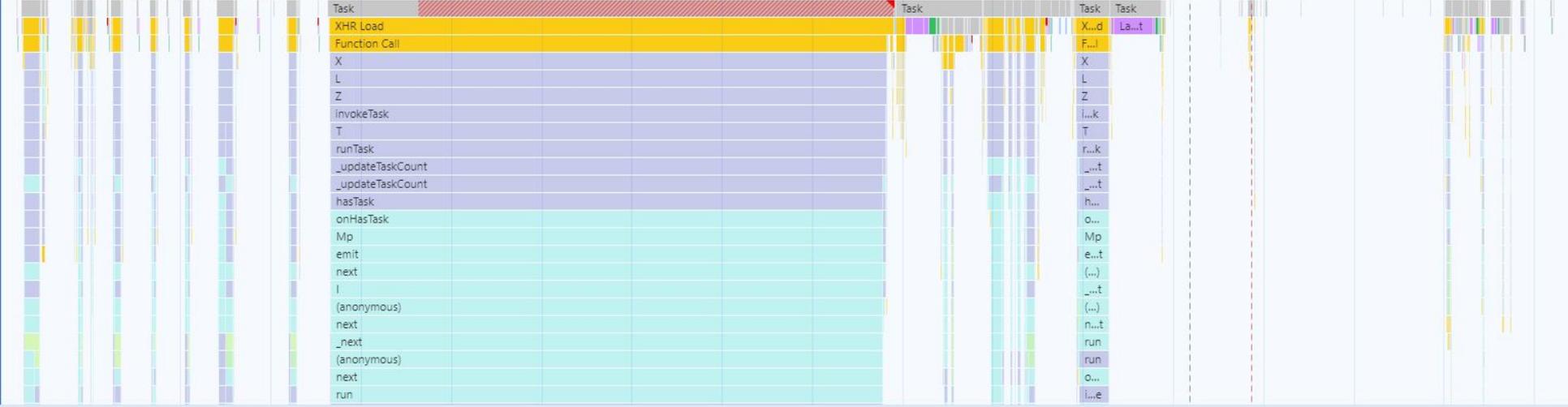
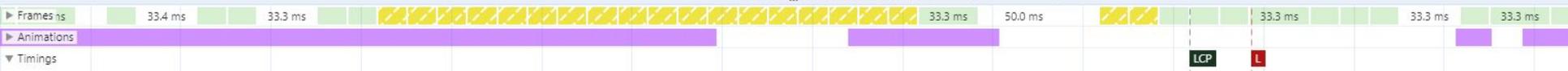
## Waterfall View



Latency со стороны клиента  
время до первого байта от сервера



В дебагере хрома есть неплохой такой профайлер



# А Пользователь?

- Больше 300мс заметно
- Делай **лоадер** если дольше, но лучше чини
- Не у всех хороший инет, **лоадер** нужен даже если сервер сегодня быстрый

Конец

Вопросы?

