

# Эффективные сервисы ML Inference нейросетей в Яндекс рекламе

Дмитрий Ульянов

Сентябрь 2024

The Yandex logo, featuring a red 'Y' followed by the word 'andex' in black.The Yandex logo, featuring a red 'Y' followed by the word 'andex' in black.

# Overview

- 1 О нас
- 2 Что такое Inference Server
- 3 Какую задачу мы решаем в рекламе
  - Наши масштабы
  - Open source и почему решили делать своё решение
- 4 Как устроен наш сервис
- 5 Оптимизации
  - Сплит архитектура моделей
  - Кэширование
  - Батчевание
  - GPU Inference
- 6 Итоги
- 7 Планы

О нас

Мы — команда разработки сервисов моделей в рекламе

- Отвечаем за внедрение ML моделей в продакшн рекламы:
  - ▶ Процесс внедрения
  - ▶ Сервисы применения нейросетей
- Делаем из нашего сервиса базовую технологию, чтобы
  - ▶ Внедрять модели быстрее в рекламе
  - ▶ Помогать в этом остальной части Яндекса

Реклама — огромный рекомендательный сервис.

Действующие лица:

- Пользователь
- Рекламодатель
- Площадка
- Яндекс

Задача отбора рекламы — максимизировать счастье каждого актора.

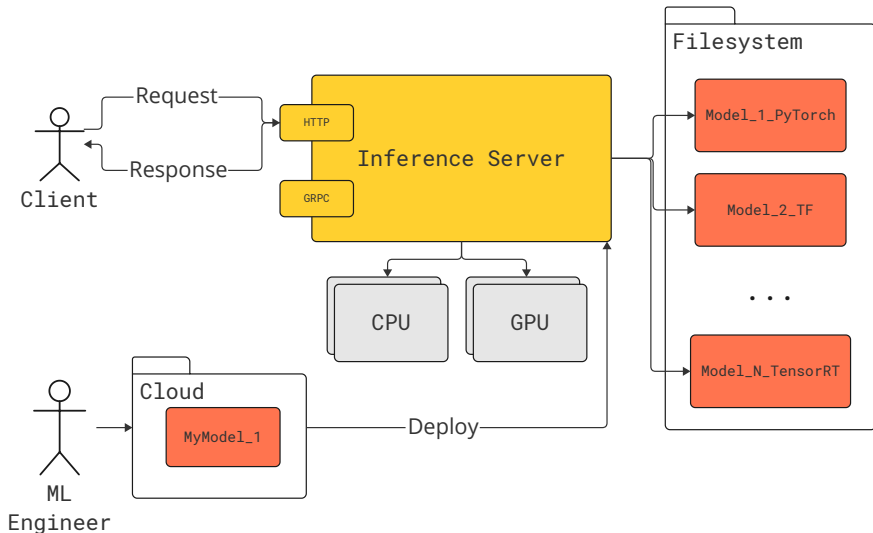
- Кандидатогенерация:
  - ▶ Отбирает релевантные пользователю и запросу баннера
  - ▶ Например, KNN индекс баннеров (HNSW)
  - ▶ Фильтрация
- Двухстадийное ранжирование баннеров:
  - ▶ Фильтрация на нижней стадии
  - ▶ Финальное ранжирование (кликочные, конверсионные прогнозы)
- Автобюджет:  
Расчёт ставок для эффективных трат рекламной компании в определённый срок
- Предсказание стоимости показа баннера на определённой позиции
- Реклама в телеграмме
- etc.

# Что такое Inference Server

- **Inference** ML модели/нейросети — применение на конкретных входах/запросе
- **Вход на Inference** модели — именованный набор тензоров
- **Inference Server** — сервис с некоторым API, который позволяет деплоить ML модели и отвечать на запросы Inference от пользователя
- **Деплой ML модели:**
  - ▶ Развернуть новый Inference сервис в облаке с конкретной моделью
  - ▶ Выкатить модель под уже существующий Inference сервис
- **Фреймворк инференса / Inference Backend** — формат модели и код применения.  
Примеры: PyTorch/TensorFlow/ONNX/TensorRT



# Схема Inference Server



# Какую задачу мы решаем в рекламе

# Какую задачу мы решаем в рекламе

Разрабатываем ML инфраструктуру, в частности сервисы ML Inference.

Наши сервисы должны:

- Быть **эффективными**
- Иметь хороший time-to-prod
- Быть надёжными

## disclamer

В этом докладе в **не** будем рассматривать MLOps процессы, а сосредоточимся на **оптимизациях**, которые мы внедряли.

Какую задачу мы решаем в рекламе:

Наши масштабы

- Поисковая реклама:
  - ▶ ~ Десятки тысяч RPS
  - ▶ ~ Десять тысяч CPU
  - ▶ 4 сервиса
- Рекламная сеть Яндекса:
  - ▶ ~ **Сотни тысяч RPS**
  - ▶ ~ **Десятки тысяч CPU**
  - ▶ 4 сервиса
- Контент система рекламы:
  - ▶ Десятки миллиардов баннеров в день ~ 10 – 100 тысяч RPS
  - ▶ ~ Десятки тысяч CPU
  - ▶ 5 сервисов
- Другое (реклама в телеграме, генерация баннеров, ...)
- **Много GPU**
- Время генерации ответа рекламы — 200 ms Q99

Какую задачу мы решаем в рекламе:

Open source и почему решили делать своё решение

# Triton Inference Server

- <https://github.com/triton-inference-server>
- Даёт только C++ библиотеку:
  - ▶ Можно развернуть сервис в kubernetes со статическим набором моделей
  - ▶ Систему деплоя моделей в сервис и многое другое нужно настраивать самому

## Плюсы:

- Готовое достаточно эффективное решение
- Поддерживает большинство фреймворков инференса (TensorRT, ONNX runtime, ...)
- Широко используется известными компаниями
- Поддерживает переподгрузку моделей без рестарта сервиса

## Минусы:

- Недостаточно эффективный на большом количестве маленьких входов
- Налаживать инфраструктуру нужно самому (шардирование моделей по подам, autoscaling):
  - ▶ Шардирование моделей по подам
  - ▶ Autoscaling
  - ▶ Hot reload: для переподгрузки моделей без рестарта есть только API, нужно разработка для полной поддержки

Другие:

## 1. BentoML

Написан на питоне. Можно использовать как самостоятельно, так и как надстройку над triton.

Плюсы: Гибкий, простой

Минусы:

- ▶ Не подходит из-за производительности
- ▶ Всё ещё требует интеграции в яндекс инфру

## 2. Seldon core, Ray, ...

Имеют примерно те же плюсы и минусы, что BentoML

## 3. AMD inference server/kserve/torch serve/... — являются очень похожими по смыслу на triton.

- ▶ Как правило, ничем не лучше triton
- ▶ Поддерживают меньше фреймворков инференса
- ▶ Имеют те же недостатки, что triton

## 4. TLDR; Посмотрите замечательный доклад:

PyterPy 2023, «Model serving. Какой фреймворк выбрать?»

<https://piterpy.com/archive/2023/talks/f0472db53d144a81af68e952783ee756/>



# Почему решили развивать свой Inference Server

Итог по open source технологиям: Вероятно одно из самых лучших решений — взять triton и строить вокруг него свою инфру.

Почему свой:

1. Легче разрабатывать **свои модули** к сервису
2. Проще и быстрее делать **интеграцию в инфраструктуру** Яндекса
3. Верим, что можем сделать **оптимальнее**
4. Широко используем **ragged тензоры**.

Пример: история тайтлов баннеров пользователя.

issue на гитхабе про поддержку нескольких ragged размерностей висит с декабря 2022

# Как устроен наш сервис

- **InferenceBackend.**

Позволяет подключать разные форматы модели к сервису.

Примеры: Pytorch, TensorRT, TensorFlowLight, ONNX Runtime

- **ModelInstance.**

Метод Infer принимает на вход и возвращает именованный набор тензоров:

$$\text{Infer} :: \left[ (\text{InputName}, \text{Tensor}) \right] \rightarrow \left[ (\text{InputName}, \text{Tensor}) \right]$$

- **Runtime.**

Вся бизнес логика

- **Server.**

Предоставляет API (Http/Grpc/WebSocket) Server

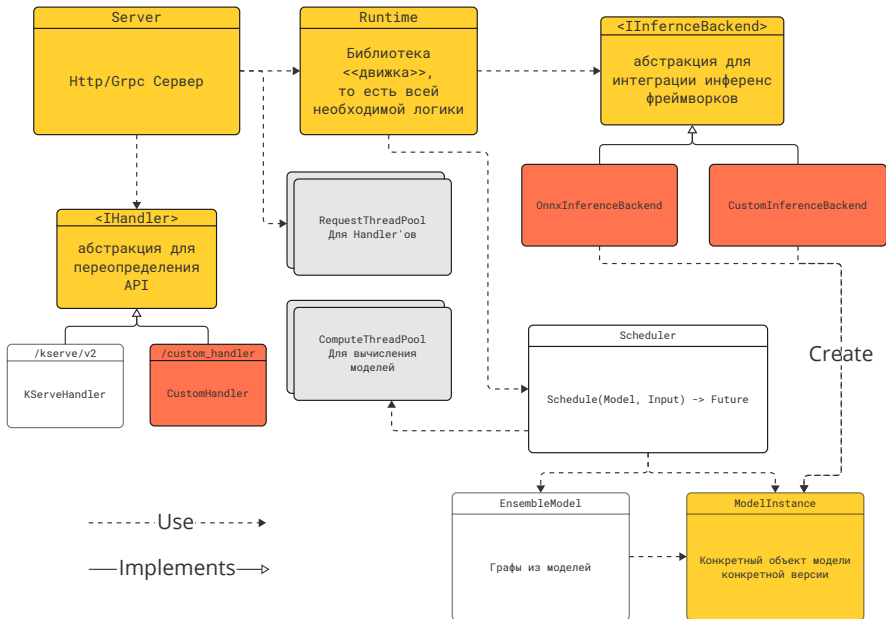


Figure: Схема основных компонент сервиса

# Примеры конфигов нашего сервиса

## application.txt.pb

```
ProgramConfig: {
  ServerConfig: {
    IoServerThreads: 2
    RequestThreads: 2
    Port: 80
  }

  TracingConfig {
    ServiceName: "my_gpu_inference"
    SampleRatio: 0.1
  }
}

RuntimeConfig: {
  InferencePoolSize: 20
  SharedBatcherConfig: {
    BatchingThreads: 2
  }
  ModelRepositoryConfig: {
    Path: "/workloads/ml_inference/repo"
  }
}

ManagedByReloader: true
```

## my\_model.yaml

```
backend_config:
  backend: unitednn
  custom_params:
    - key: lock_memory
      value: true
    - key: subnetwork
      value: _head
    - key: gpu
      value: '0'
runtime_config:
  batching_config:
    max_batch_size: 8
    batch_timeout_mcs: 2000
```

# Оптимизации

Оптимизации:

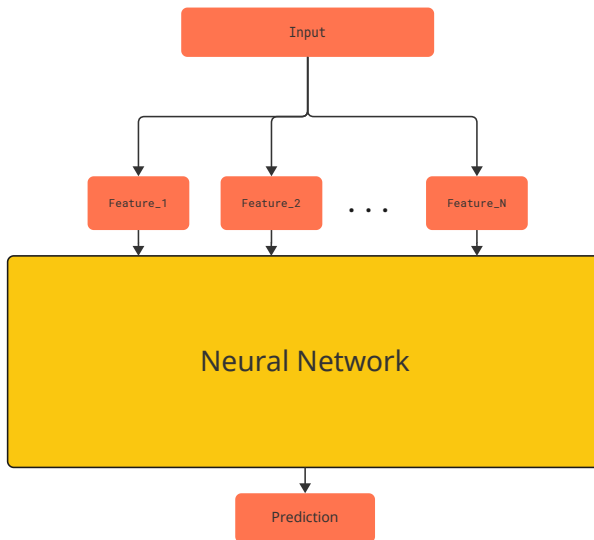
Сплит архитектура моделей

# Проблема 1: Количество кандидатов

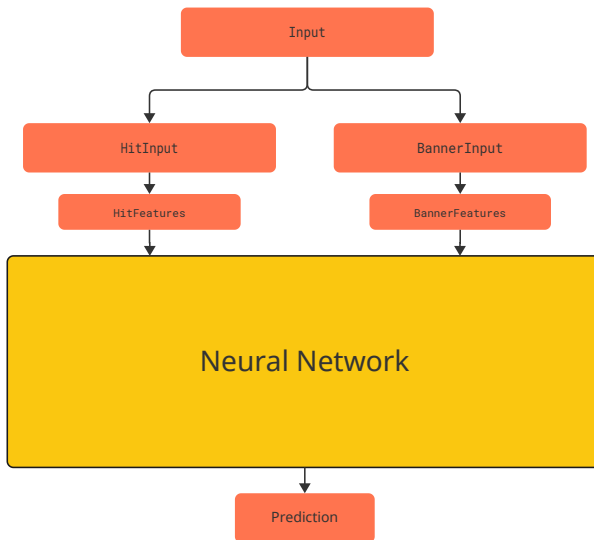
1. Возьмём Рекламную Сеть Яндекса
2. 300 000 RPS
3. Несколько сотен баннеров-кандидатов на каждый запрос
4.  $\implies$  60 миллионов документов в секунду
5. Считать нейросети на таком потоке *очень* дорого
6. Поэтому большинство моделей имеют архитектуру позднего связывания (или split)



# Нейросетевая модель в общем случае

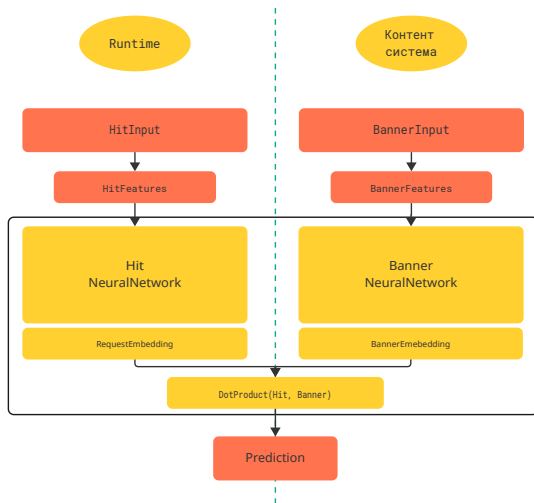


# Нейросетевая модель в рекламе

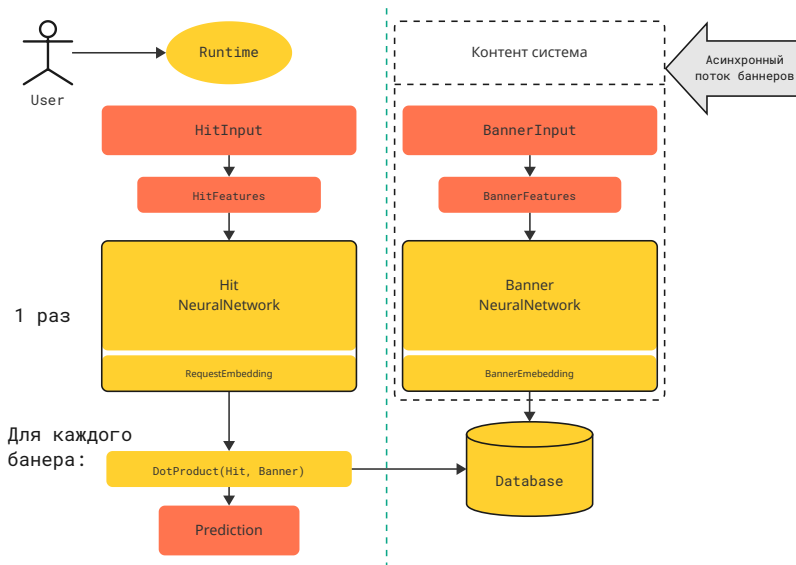


# Нейросетевая split модель в рекламе

- Учим модель специальной архитектуры (похоже на DSSM)
- Hit и Banner слои связываем в конце через скалярное произведение
- После обучения делим модель на две части
- Hit часть считаем в рантайме
- Banner часть считаем в контент системе
- В рантайме достаём эмбединг баннер части из базы
- Predict считаем как DotProduct



# Нейросетевая split модель в рекламе



## Плюсы:

1. Кратно оптимизируется нагрузка (сотни раз)

## Минусы:

1. Требует изменений архитектуры
2. Можно потерять в качестве.  
Эксперименты в рекламе не доказали этого (но и не опровергли)
3. Требует разработки:
  - ▶ Сервис в контент системе
  - ▶ Доставка эмбедингов в базу
  - ▶ Лукап в базе и подсчёт DotProduct

Оптимизации:

Кэширование

# Кэширование слабого человека

Добавление кэширования выглядит просто:

- Пишем или берём готовый кэш
- Profit

На практике:

- Наивное кэширование по всем входам даёт очень маленький кэшхит
- Кэшируем по подмножеству входов
- Входы по User кэшируем по UserId на TTL (TTL = 120) секунд

Пример:

```
Input {
  UserProfile {
    UserID: 42, // участвует в ключе кэша
    ...       // не участвует в ключе кэша
  }
  HitContext {
    PageID: 1 // номер площадки рекламодателя, участвует в ключе кэша
    PageTitle: "Super Title" // Зависит от PageID, не участвует в ключе кэша
    ...
  }
}
```

**Итог: 10-20% кешхит** -> такая же экономия CPU .

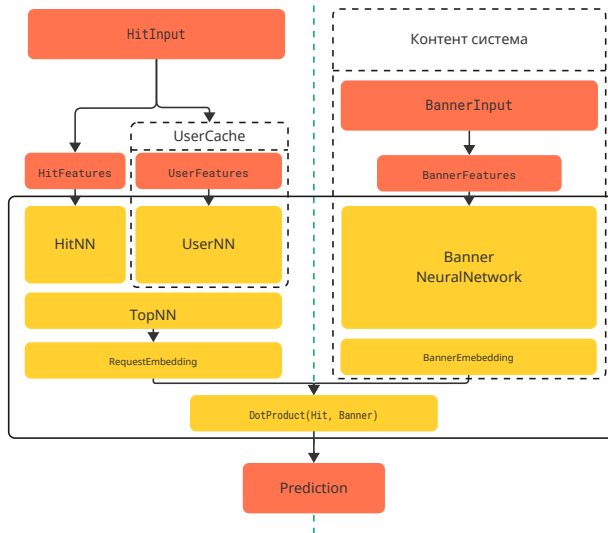
# Кэширование сильного разработчика

## Алгоритм:

- Выделим ещё кусок у модели — UserNN
- Сделаем по нему отдельный кеш

## Профит:

- Кешхит по User равен 60%.
- Допустим CPU (UserNN) = CPU (HitNN + TopNN)
- Экономия =  $60\% \cdot 0.5 = 30\%$



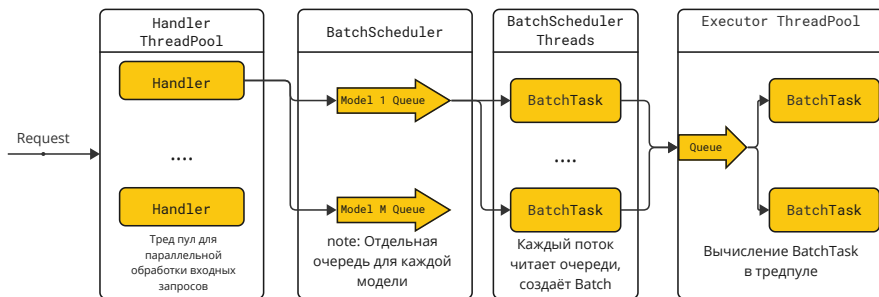


# Оптимизации:

## Батчевание

# Батчевание — схема

- Складываем запросы в очередь и отдельными тредами выгребаем элементы из очереди батчами:
  - ▶ Параметры: `BatchTimeout` и `MaxBatchSize`
  - ▶ `MaxBatchSize` — максимальный размер батча
  - ▶ Ждём накопления батча не более `BatchTimeout` времени
- must-have фича на GPU сервисах
- Оказывается, экономит много даже на CPU сервисах



# Почему экономится CPU ?

Представим, что:

- Вход модели – вектор размера  $m$
- Модель — умножение вектора на матрицу
- В батче  $n$  элементов

Соображения:

- Умножать матрицу  $n \times m$  на матрицу можно асимптотически быстрее, чем  $n$  раз вектор на матрицу  
Модели используют суперумную библиотеку Eigen для работы с матрицами
- При умножении батча входов на матрицу вы меньше ходите в память (за весами матрицы)
- Скорее всего, помогают кеша процессора на обращения в память

Профит:

- Экономия не менее 20+% CPU сервиса
- Пример: Потребление CPU одного из сервисов в момент включения батчевания: 6800 -> 5150 ядер

Минусы:

- Требуется настройка параметров батчевания
- Ухудшает latency (время ответа)

# Оптимизации:

## GPU Inference

Код сервиса:

- Сделали некоторый Preprocessing — CPU работа
- Пошедулили вычисление модели на GPU и дождались
- Вернули результат

```
GpuServiceHandler
- IScheduler Scheduler;

modelInput = Preprocess(request);
modelResult = scheduler.Schedule(modelName, modelInput);
modelResult.Wait();
return modelResult.Get();
```

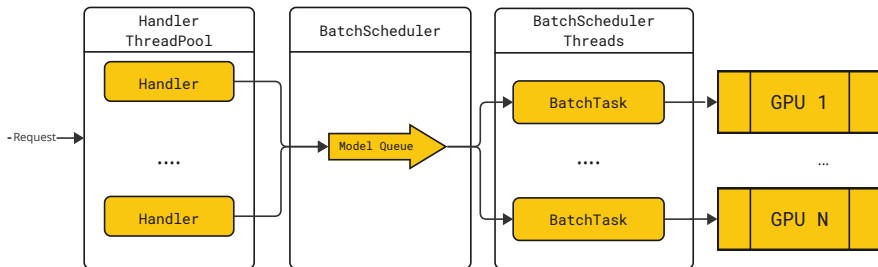


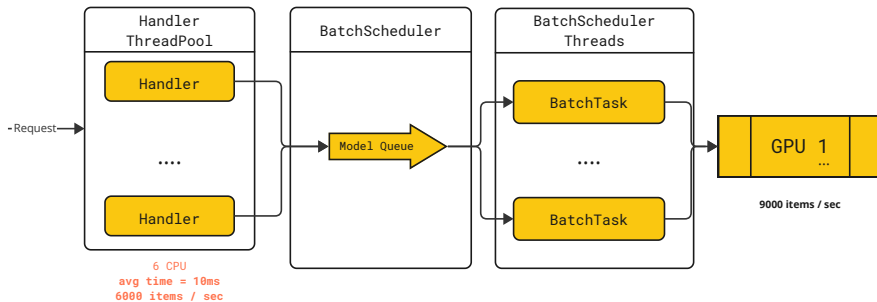
Figure: Схема запроса в GPU сервисе с одной моделью

```
GpuServiceHandler  
  
- IScheduler Scheduler;  
  
modelInput = Preprocess(request);  
modelResult = scheduler.Schedule(modelName, modelInput);  
modelResult.Wait();  
return modelResult.Get();
```

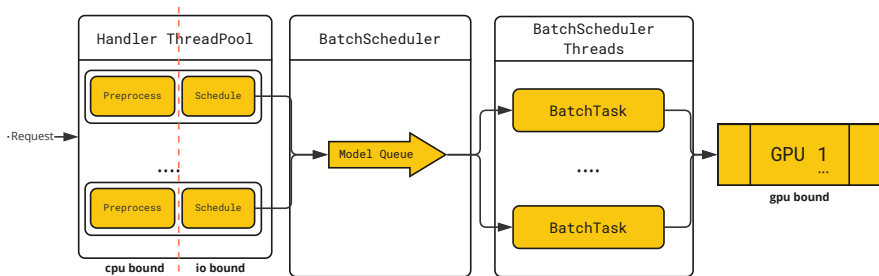
Карточка процессит 9K запросов в секунду.

Конфигурация GPU хоста:

- 56 ядер
- 8 GPU карт



# Заметим что...





# Разделяй и получай профит

Отмасштабируем!

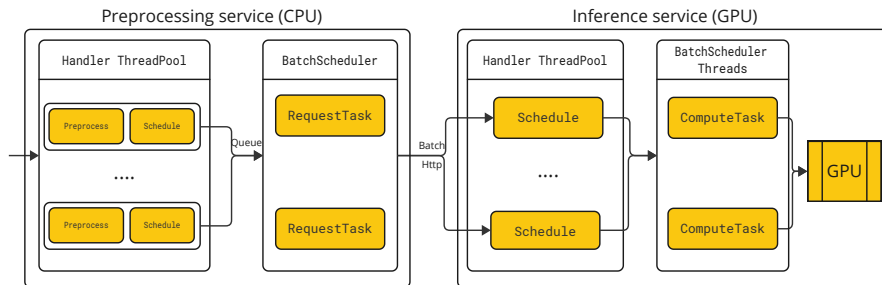


Figure: Схема разделённого по CPU only / GPU подам сервиса

Итог:

- Сэкономили **30%** GPU
- Потратили  $\times 1.5 - 2$  CPU на сервис и перекладывание байт не дорого, по сравнению с экономией GPU

# Итоги

1. Сплит модели (Hit/Banner)  
Экономия: **в сотни раз CPU и GPU**
2. Кэширование v1  
Экономия: **10-20% CPU**
3. Кэширование v2 (сплит User/Hit)  
Экономия: 10-60% CPU , **30% CPU** в среднем
4. Батчевание  
Экономия GPU : в десятки-сотни раз.  
Экономия CPU :  $\geq$  **20%**.
5. GPU split на CPU preprocessing и GPU inference  
Экономия: **30% GPU**

- Оптимизировать можно и нужно
- Оптимизации требуют разработки и часто уместны только начиная с какого-то объёма мощностей
- Большинство крупных оптимизаций продуктовые. Но низкоуровневыми оптимизациями заниматься тоже нужно
- Оптимизации — это tradeoff latency, throughput и качества модели

Мой совет:

- Ищите продуктовые оптимизации

# Планы

Что у нас уже есть:

- Inference Server с основной функциональностью
- Набор инструментов для интеграции в систему деплоя и мониторинга
- Поддержка нескольких InferenceBackend, использующихся в моделях рекламы
- Толпа сервисов в рекламе с немного разными функциональностями и кодом
- Сервисы в маркете и медиасервисах на нашей технологии

Развиваем и обобщаем нашу технологию:

- Продолжаем оптимизировать и разрабатывать фичи:
  - ▶ Кэши
  - ▶ Динамическая балансировка (учитывает throughput пода)
  - ▶ Разные алгоритмы шедулинга вычисления моделей на GPU
  - ▶ etc.
- Инфраструктура
  - ▶ IAC (Infrastructure As a Code)
  - ▶ Автоскейлинг & автошардирование
  - ▶ Inference As A Service
- Инструменты для замера performance моделей

Внедряем:

- Унифицируем сервисы и ускоряем внедрение ML в рекламе:
  - ▶ Амбиция — time-to-exp  $\sim$  2 дня
- Web Поиск
- Различные рекомендательные сервисы Яндекса
- Весь Яндекс

Мы нанимаем ;) )