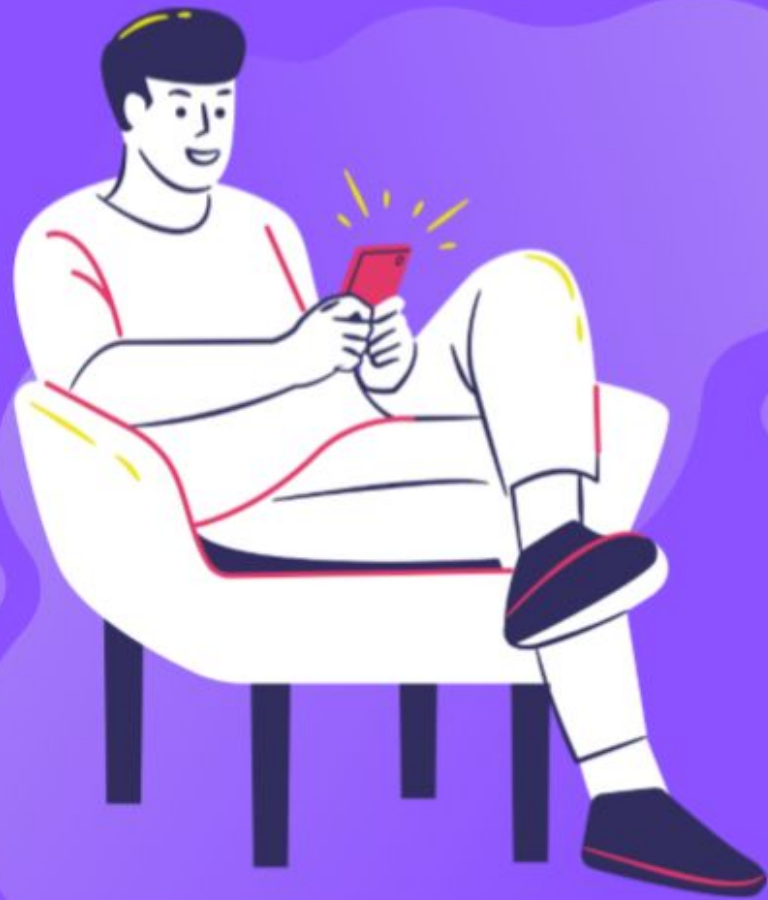


# Разбор кейса

## В поисках интересного



## Задача

Разработать концепцию архитектуры сервиса по выделению наиболее интересного фрагмента трека для пользователя с целью сокращения длительности аудио контента.

## Бизнес-метрики

- Индекс удовлетворенности CSI
- Количество пользователей за период



## Основная гипотеза для MVP

Наиболее популярные жанры музыки имеют структуру треков с повторяющимися частями, в которых сконцентрированы самые запоминающиеся ходы и приемы (хуки) — в припевах. При этом зачастую в музыке непосредственно перед припевом имеются гармонические ходы с нарастанием напряжения перед последующим разрешением и наиболее эмоционально заряженные строчки в песне.

Будем считать часть перед припевом и припев самым интересным местом.

Таким образом, задача сводится к Music Track Structure Segmentation/Classification и поиску повторяющихся частей песен с последующей разметкой снippets.

Статистика по популярности жанров: 1, 2, 3

## Papers and Technologies Research

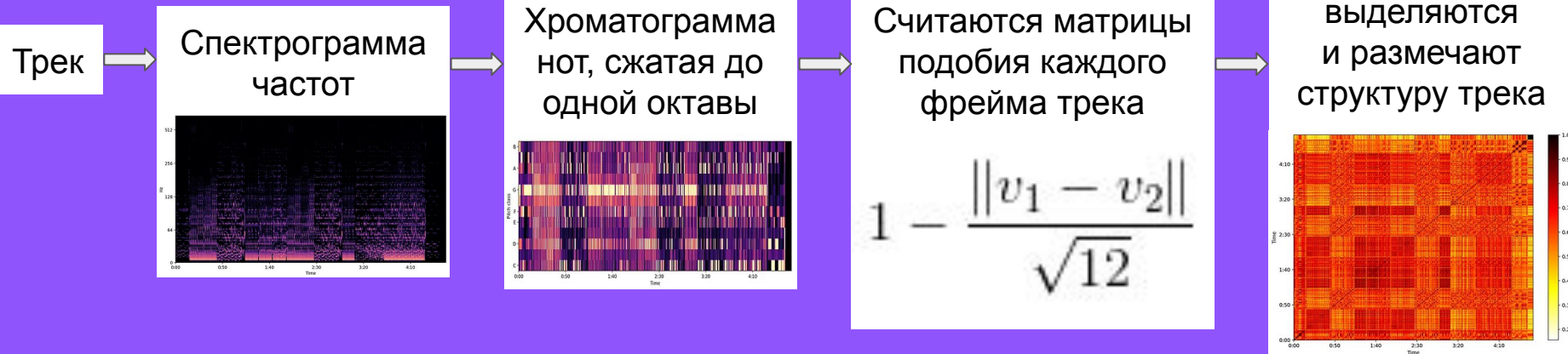
- Finding Choruses in Songs with Python на основе статьи
- Audio-Based Music Structure Analysis: Current Trends, Open Challenges, and Applications
- Music Segmentation PhD Defense
- Automatic Structural Segmentation of Music
- Youtube Most Replayed Part Prediction
- и так далее...



# Non-ML Baseline

5

**Идея:** Разметка структуры трека на основе анализа повторения нот — находим припев с помощью матриц подобия и размечаем необходимый отрывок. [\[источник\]](#)



## Pros

- Простота и дешевизна: не нужны вычисления на видеокартах и разметка

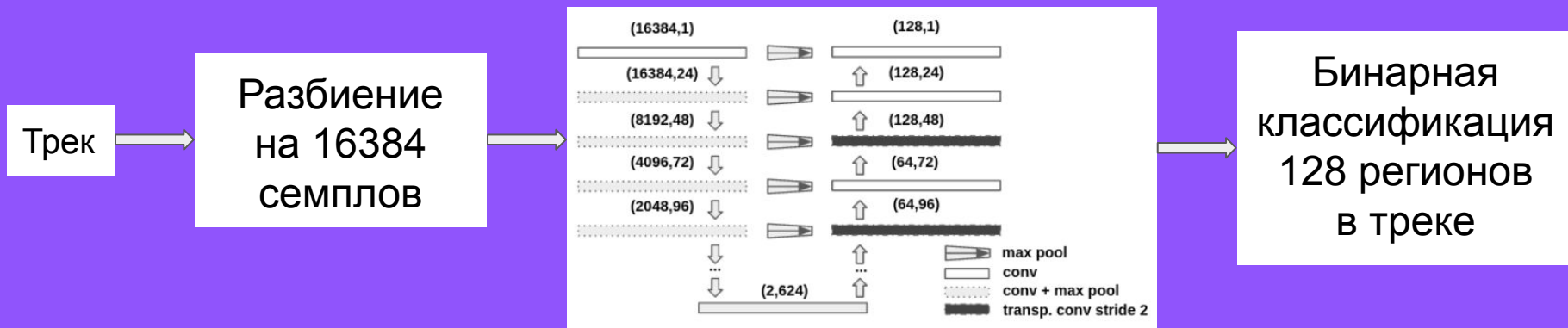
## Cons

- Возможна низкая точность, т.к. выделяются повторенные один в один части трека
- Низкая скорость инференса

# CV Solution - Hook-net

6

**Идея:** После выделения фичей из аудио, обучаем сверточную нейросеть на основе U-net — Hook-net, которая классифицирует отрезки треков по частям и проводит сегментацию. [\[источник\]](#)



## Pros

- Учитывается структура трека

## Cons

- Накладные расходы выше Non-ML Baseline

# RNN Solution - LSTM

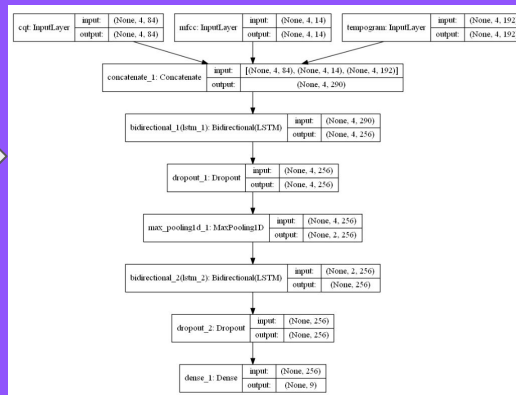
7

**Идея:** аудиоинформация и, в частности, музыка может быть обработана не только как изображения, но и как последовательность — возможно применение NLP подходов — Bi-Directional LSTM сеть так же справляется с задачей сегментации и анализа структуры треков. [\[источник\]](#)

Трек

Получение признаков методами DSP:

Constant-Q Transform  
Mel-Frequency Cepstral Coefficients  
Tempogram



Мультиклассовая  
классификация каждого  
фрейма  
последовательности

Pros

- Теоретически может учитывать контекст повторяющихся структур внутри трека при поиске самого интенсивного припева, а не просто размечать структуру

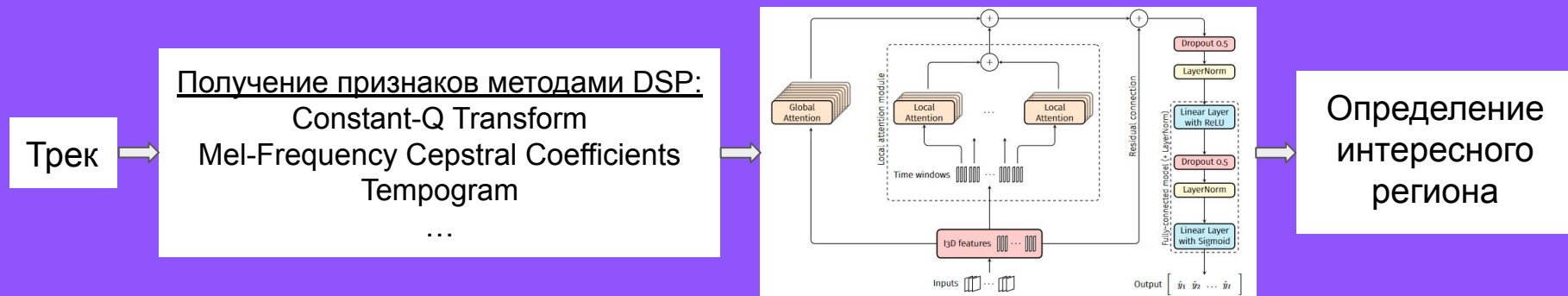
Cons

- Более тяжелая архитектура сети

# Transformer Solution - Attention-based PGL-SUM

8

**Идея:** Задача выделения наиболее просматриваемого фрагмента была решена для видео с помощью архитектуры PGL-SUM, основанной на двух контурах Attention — глобальной и локальной. Возможно данная архитектура справится с аналогичной задачей для аудио. [\[источник\]](#)



## Pros

- Двойной контур Attention должен улучшить качество разметки по сравнению с LSTM

## Cons

- Еще более тяжелая архитектура сети



## Метрики



### ML-метрики для оценки качества и мониторинга:

1. IoU — по перекрытию ground truth сегмента и предсказания
2. Accuracy
3. Precision, Recall, F1-score

### Полезные показатели для мониторинга при A/B-тестировании гипотез:

1. Mean Opinion Score — формировать по опросам пользователей
2. Conversion Rate — конверсия прослушивания трека после прослушивания снippets (возможно потребуется учитывать предпочтения пользователя для повышения адекватности показателя)

## Датасет

1. Открытые датасеты
2. Разметка на Я.Толоке
3. Youtube scraping — возможно получение данных о Most Reviewed Part через API [[Stackoverflow thread](#)]
4. Генерация примеров
5. Аугментация для увеличения датасета:
  - a. Добавление шумов и глитчирование
  - b. Изменение тональности
  - c. Перемешивание частей трека (в случае полностью размеченной структуры)

## Доступные датасеты

11

Датасет	Количество треков	Features	Примечание
<a href="#"><u>FMA</u></a>	106,574	DSP фичи, полученные с помощью librosa	Необходима разметка, 30-секундные треки
<a href="#"><u>SALAMI</u></a>	1359	Иерархические аннотации структуры треков	Ограничен по исполнителям и жанрам
<a href="#"><u>The Harmonix Set</u></a>	912	DSP фичи, метадата, структура	Ограничен по исполнителям и жанрам
<a href="#"><u>SPAM</u></a>	50	DSP фичи, метадата, структура	Ограничен по исполнителям и жанрам
<a href="#"><u>RWC</u></a>	300	DSP фичи, метадата	Платный, нет разметки структуры
<a href="#"><u>Isophonics</u></a>	300	DSP фичи, метадата, структура	Ограничен по исполнителям и жанрам



1. Разметка структуры всего трека скорее всего будет требовать толokers высокого уровня и займет не только больше времени, но и будет больше стоить. Возможно рациональней будет упростить задачу и просить толokers определить только наиболее интересный регион трека по их мнению.
2. При этом необходимо оставить возможность для выбора у толкера — какие-то жанры могут быть нелюбимыми и разметка интересного региона в таком случае будет необъективной. Вопрос объективности при опросах возможно решить методами статистики [[Пример](#)].
3. Пулы следует формировать с большим количеством перекрытий, чтобы агрегация региона после разметки была наиболее объективна.
4. Набор треков должен быть сбалансирован по жанрам и возможно составлен из жанровых топ-листов для упрощения задачи толкерам.

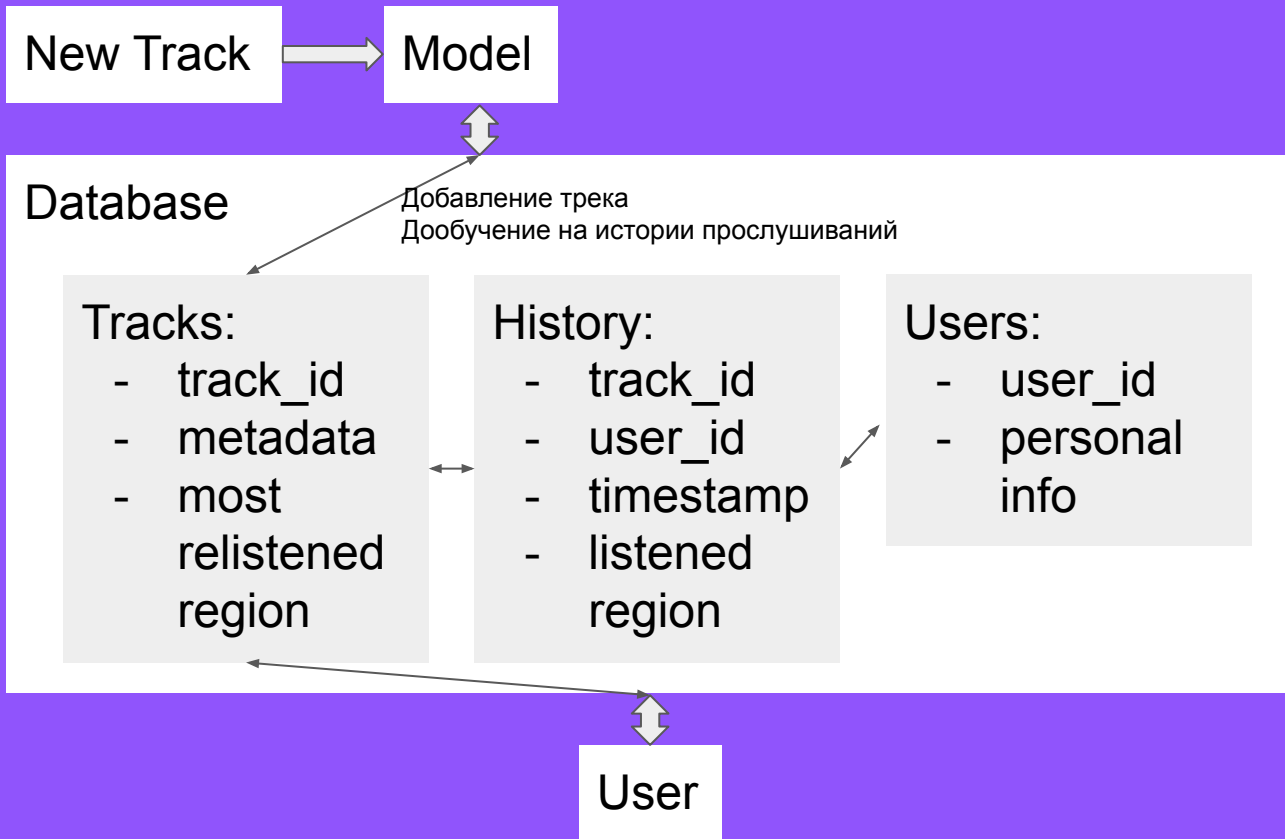
## Возможные трудности

13

Описание	Вариант решения
Жанры со сложной структурой	Техно, джаз, академические жанры имеют структуру отличную от популярных и выделение фрагмента на основе повторяющихся сегментов может вызвать трудности. Для подобных жанров можно выделять фрагмент на основе простых эвристик — начало\середина трека, например, пока не будет статистики их реального прослушивания для дообучения модели.
Различная длина треков	Для треков короче выделяемого снippets — выделяется весь трек. В остальных случаях для подачи в сеть требуется последовательность фиксированной длины. Тогда принимается величина продолжительности трека, к которой будут нормализоваться все отправляемые нейросети треки. При этом треки короче этой величины будут иметь частоту высокую семплирования, а длинней, соответственно, низкую. Для очень длинных треков, возможно, проще будет разбивать трек на части и искать по ним, или же использовать иной подход.

# Архитектура системы

14



## Технологии

librosa  
Pytorch  
PostgreSQL  
Airflow  
MLflow  
DVC



## Оценка сроков реализации

1. Сбор, подготовка и очистка датасета: ~1 месяц
2. Реализация и эксперименты с архитектурами моделей: ~2 месяца
3. Создание MVP для Demo Review: ~1 месяц

Итого: ~4 месяца на реализацию MVP

## Оценка ресурсов

1. Аренда Яндекс.Облако для обучения\инференса модели: ~75к/мес
2. ФОТ: ~100-150к/мес на сотрудника + отчисления
3. Оборудование (ноутбук для удалённой работы): ~100-150к на сотрудника

Спасибо!

