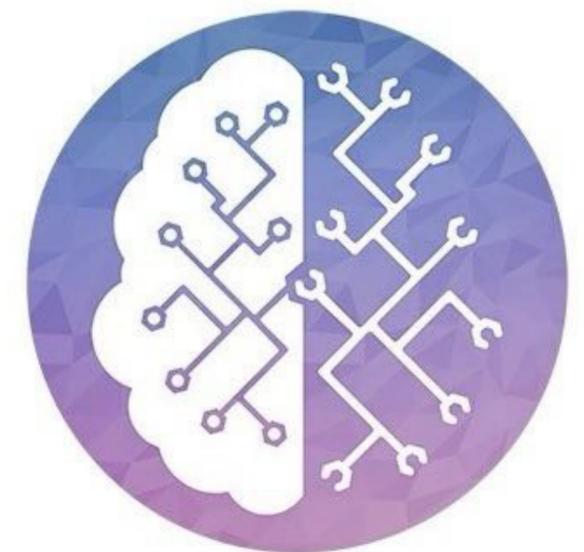


Повышение точности субтитров в художественных фильмах: технологии и опыт

Шатнёв Дмитрий, 2024

О себе

- Инженер машинного обучения в Inventos
- Студент ОГУ им. И.С. Тургенева
- Сферы научных интересов:
 - Обработка естественного языка
 - Речевые технологии



ИПАИТ
ОГУ ИМЕНИ
И.С. ТУРГЕНЕВА

План доклада

- Введение в задачу 4-6
- Первая итерация 7-14
- Препроцессинг 16-19
- Работа с моделью 20-27
- Постпроцессинг 28-29
- Итоги 30-32

Webcaster - отправная точка



<https://webcaster.pro/>

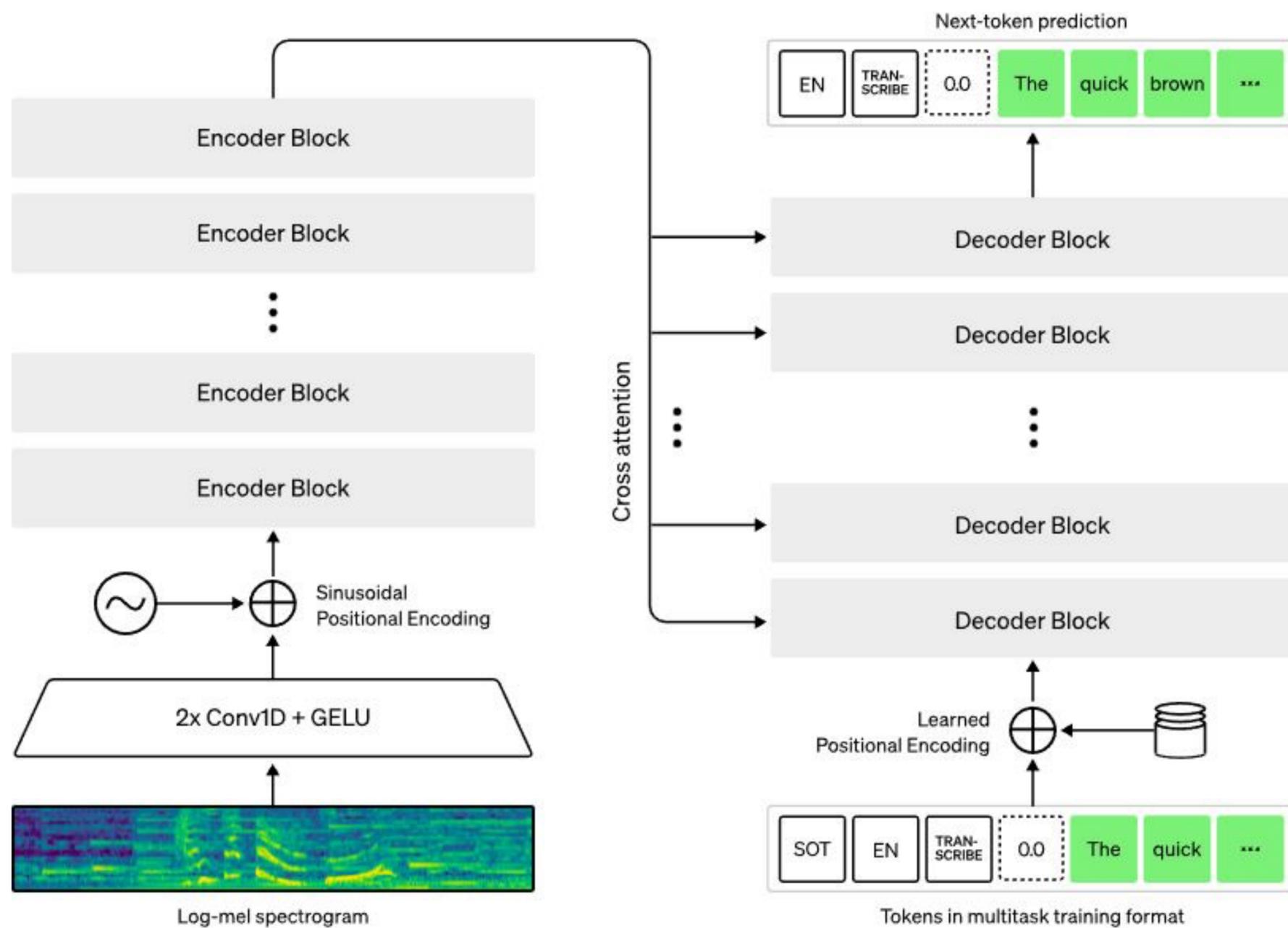
Особенности контента

- Существенная длительность
- Сложная звуковая картина
- Мультиязычность
- Аудиовизуальная корреляция



Action Vectors by Vecteezy

Модель Whisper



Первая итерация



Использование Whisper из коробки:

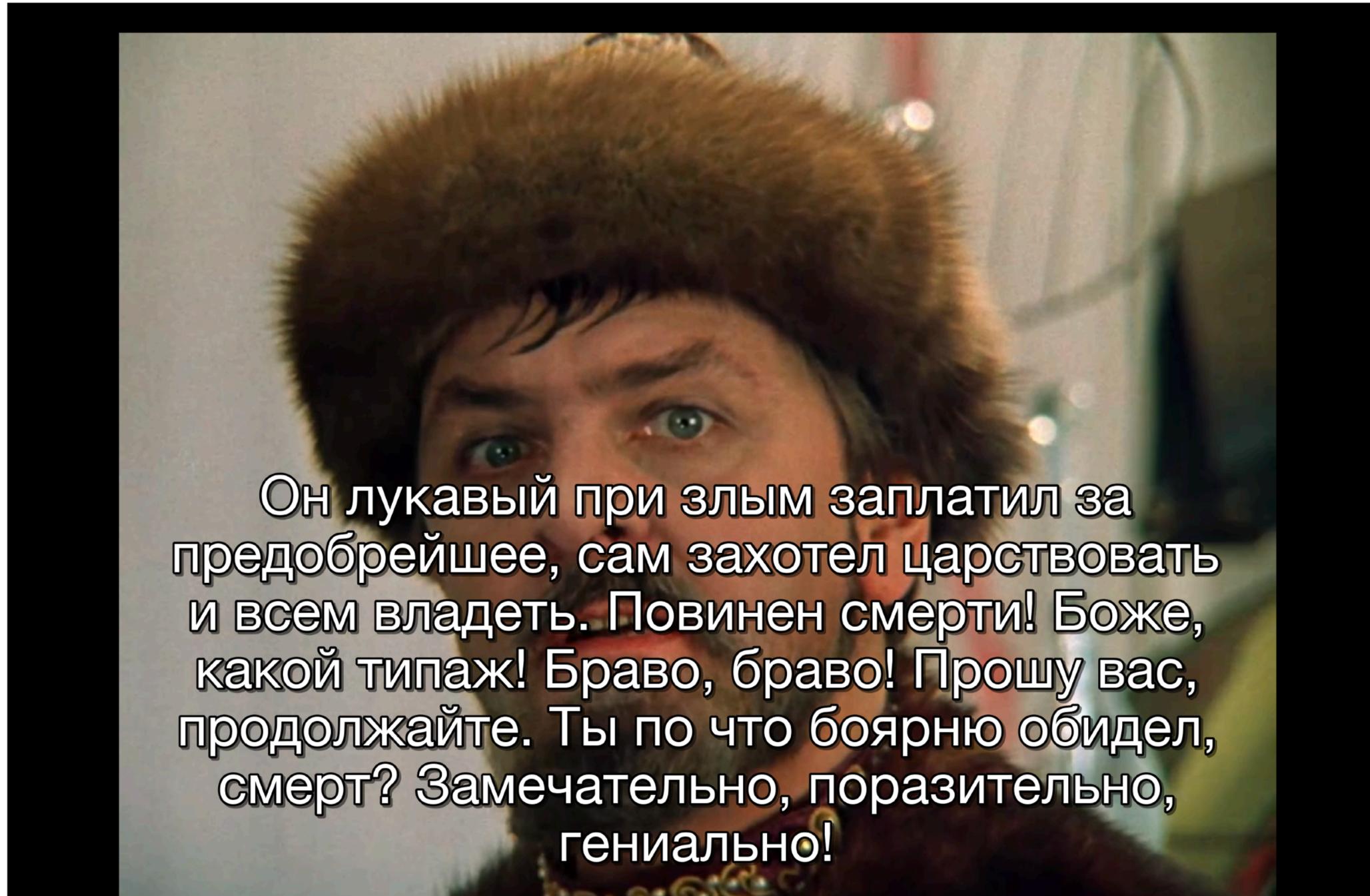
- Максимально простой пайплайн
- Минимум усилий

Bad case #1



<https://www.mosfilm.ru/>

Bad case #2



<https://www.mosfilm.ru/>

Первые проблемы

- Аномалии распознавания
- Рассинхронизация по временным меткам
- Большой объем текста в блоке субтитров
- Сложность объективной оценки

Основная метрика качества

$$\text{Word Error Rate (WER)} = \frac{\text{Insertions} + \text{Deletions} + \text{Substitutions}}{\text{Number of words in reference}}$$

$$\text{Accuracy} = 100\% - \text{WER}$$

Эталон: Жизнь как **коробка** шоколадных конфет.

Распознано: Жизнь **коровка** шоколадных конфет. **Пока!**

Дополнительные метрики качества

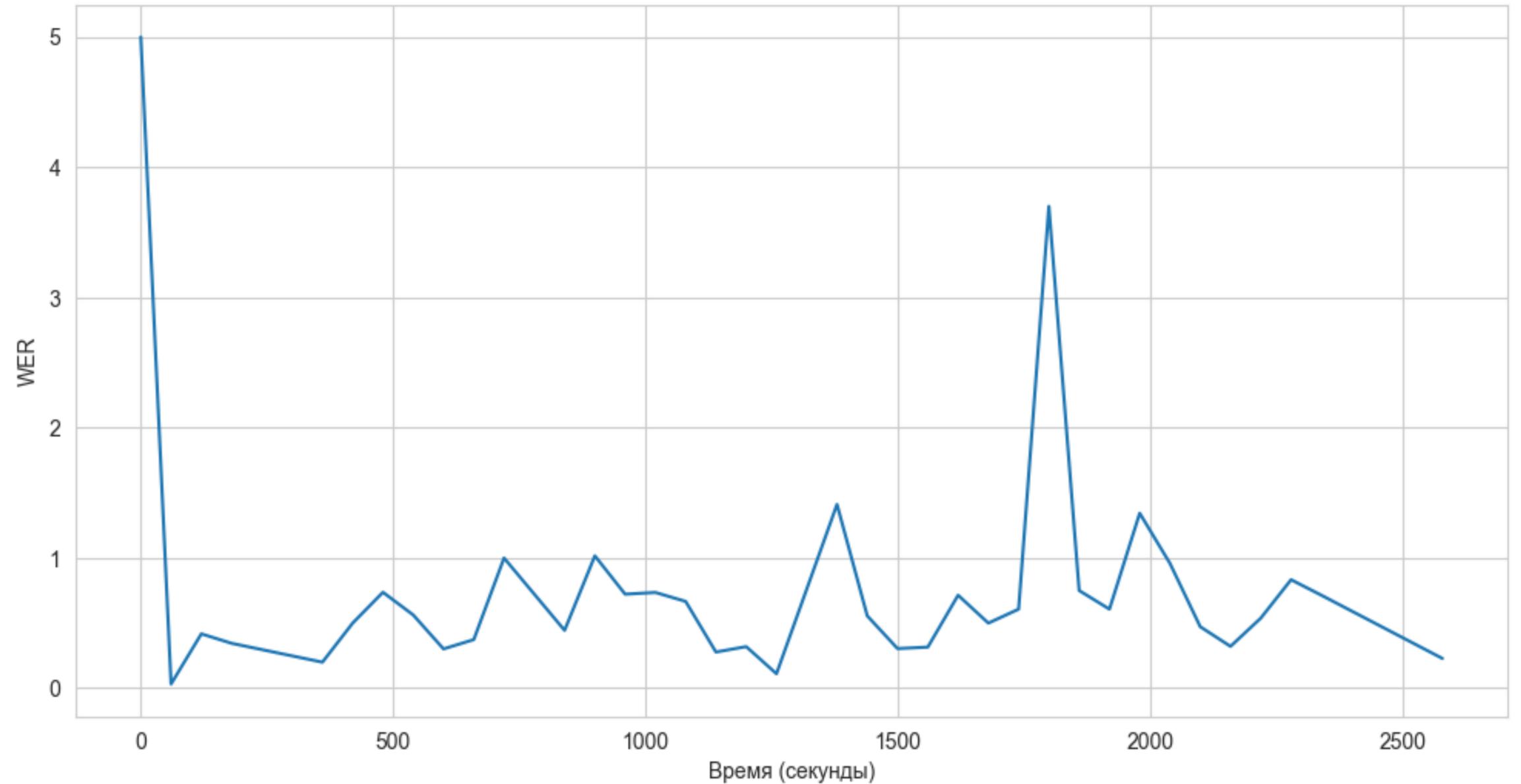
- CER (Character Error Rate)
- WIP (Word Information Preserved)
- WIL (Word Information Lost)



Metrics Vectors by Vecteezy

Локализация падения метрик

Window: 1 min.
Step: 1 min.
Audio duration: 45 min.

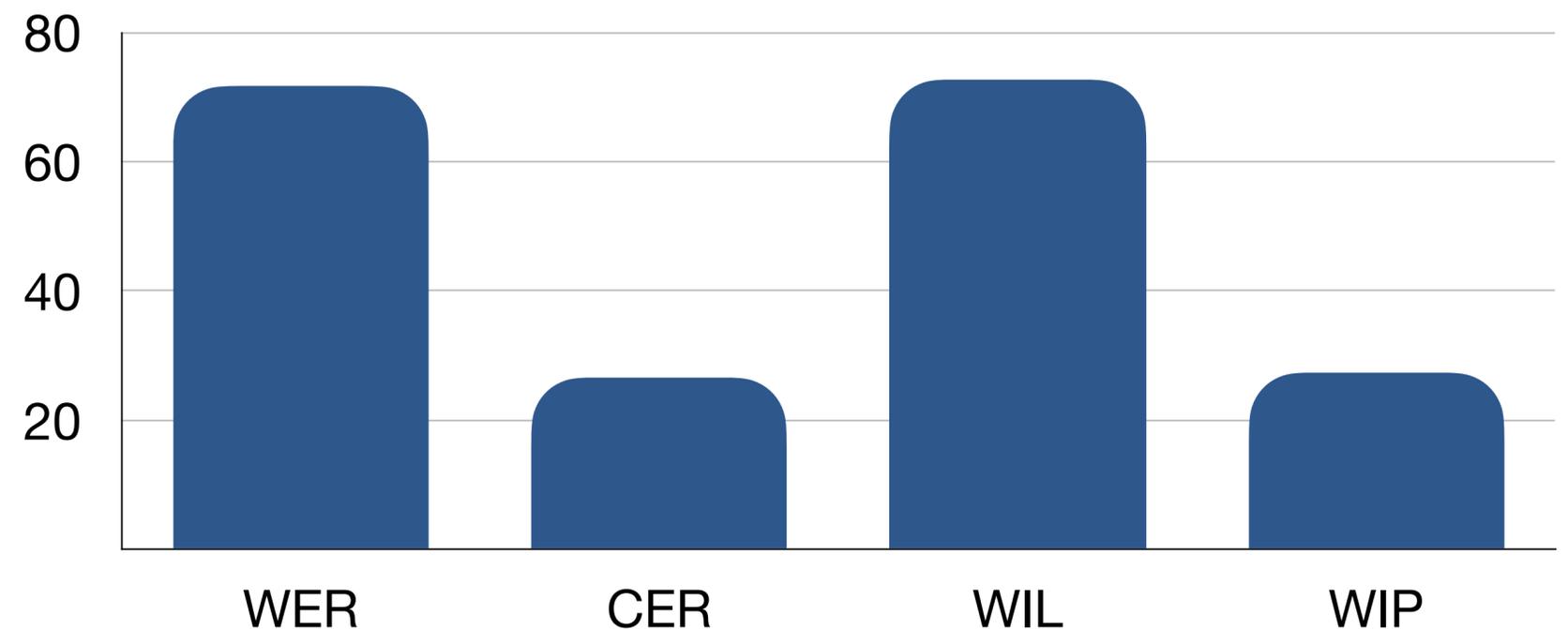


Промежуточные выводы

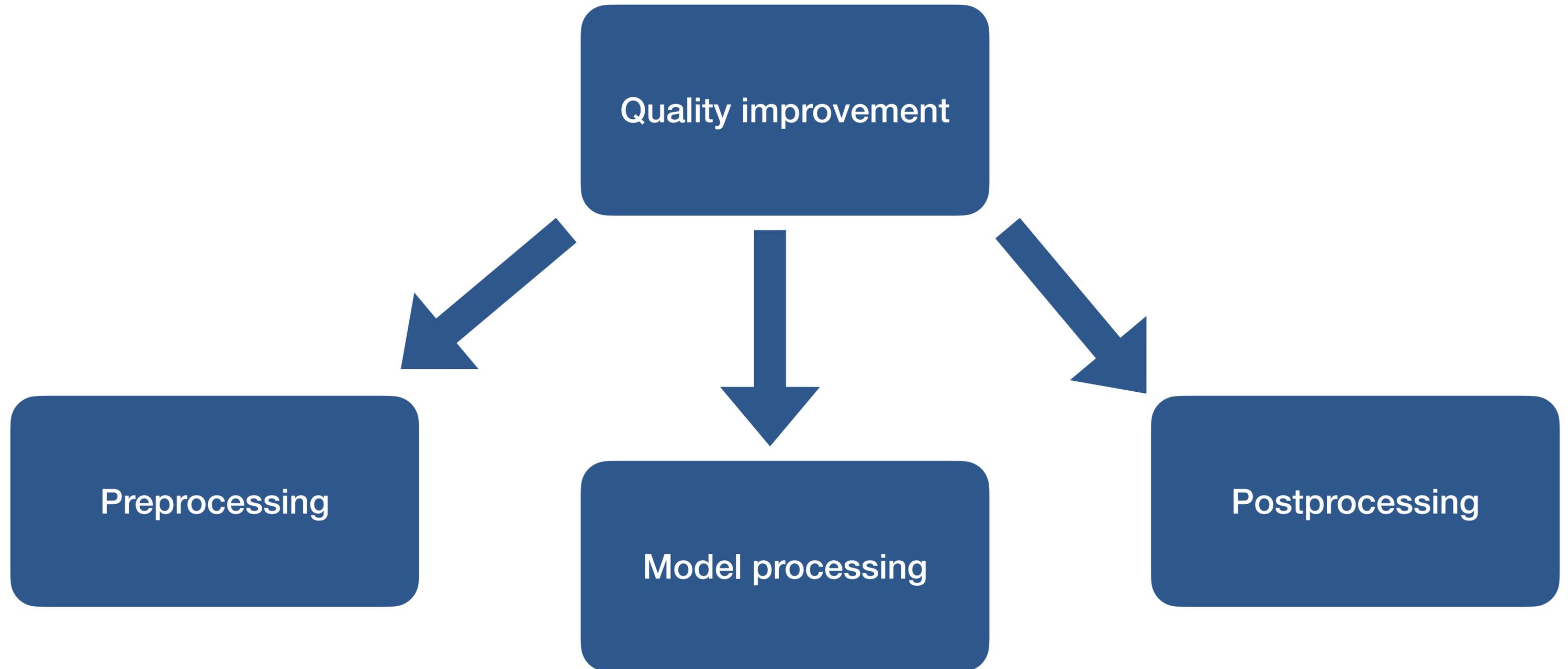
Текущие показатели свидетельствуют о **низкой** точности распознавания.

Необходимо комплексно оптимизировать систему, для достижения приемлемого качества субтитров.

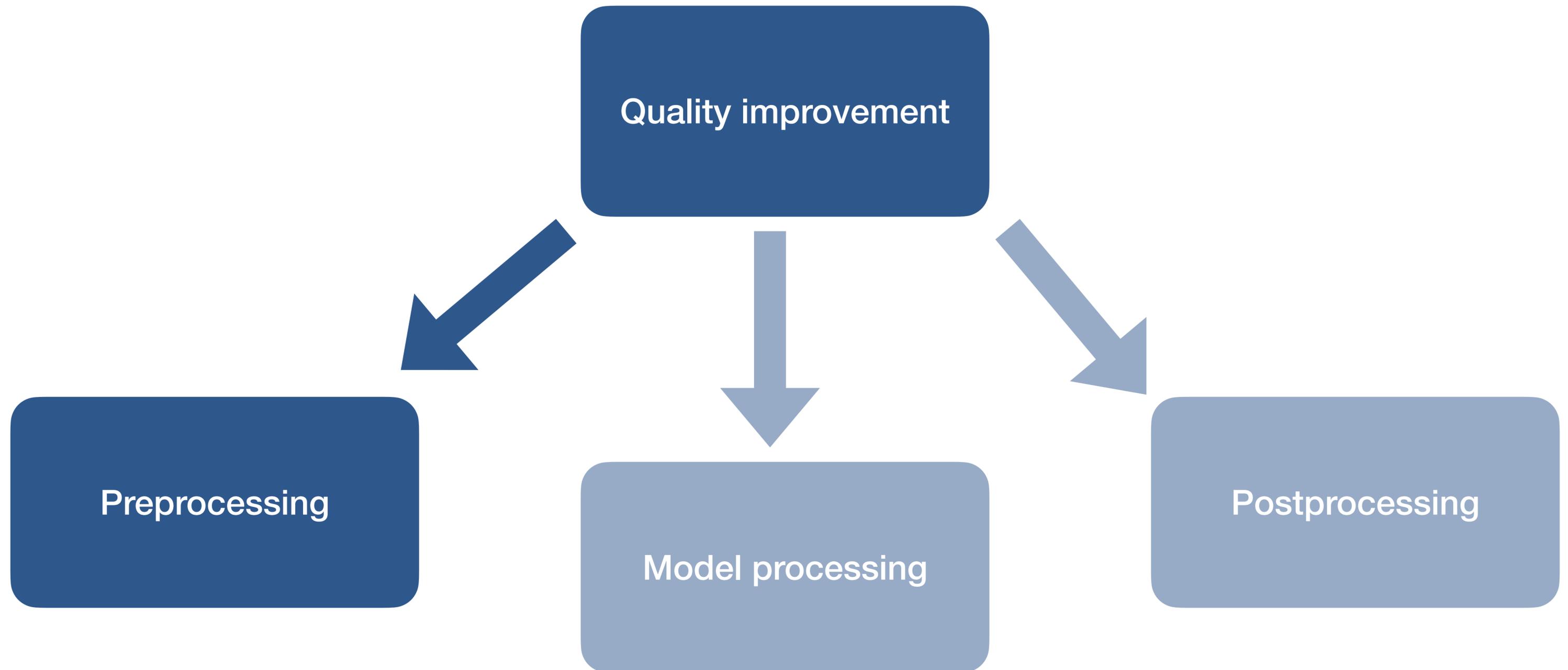
Средние метрики качества



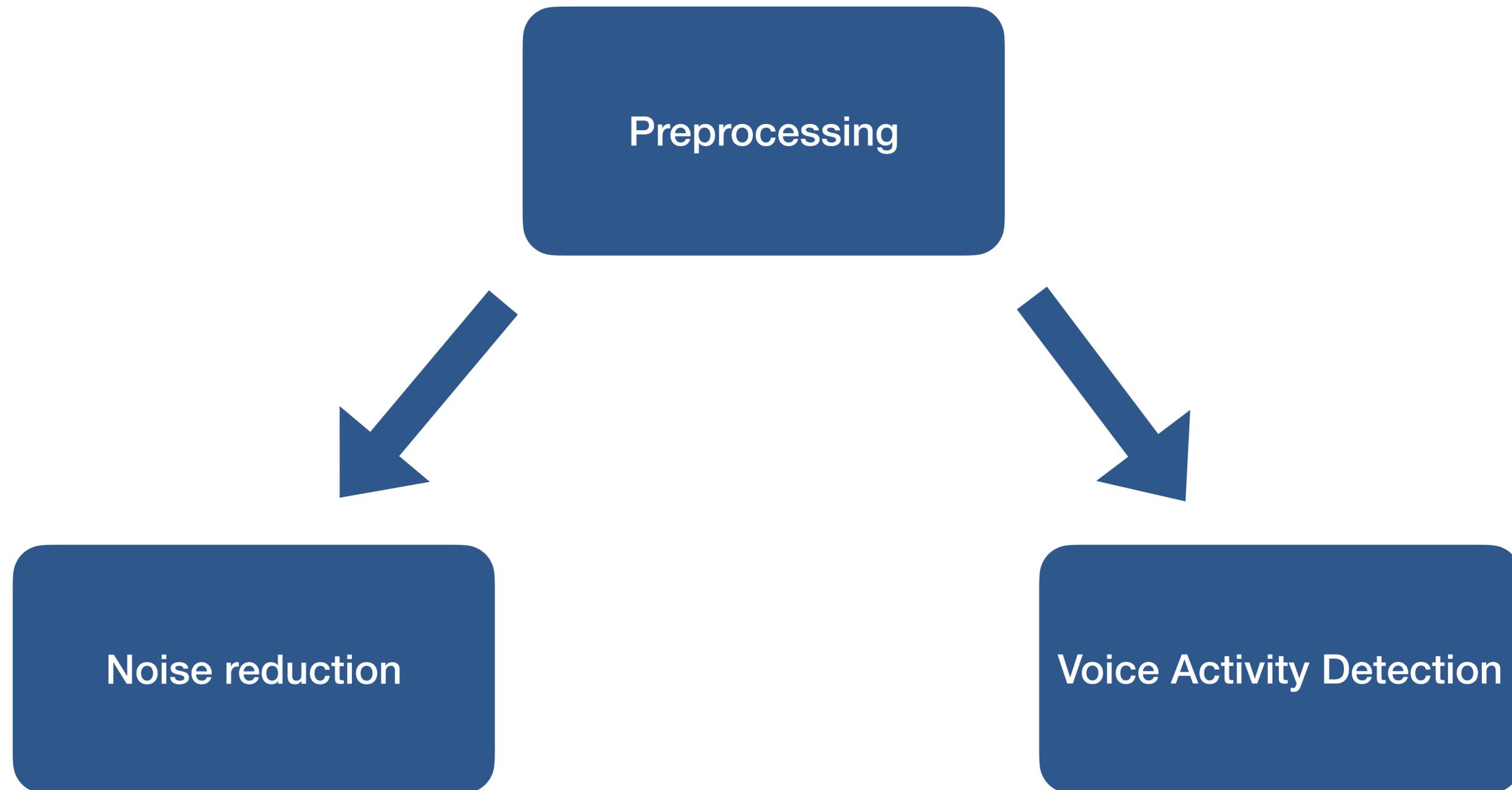
Способы улучшения качества



Препроцессинг



Вариации препроцессинга



Шумоподавление

В теории:

Улучшает качество
входного аудио

Облегчает выделение
речевого сигнала

Уменьшает количество
ошибок распознавания

На практике:

Удаляет важные звуковые
детали

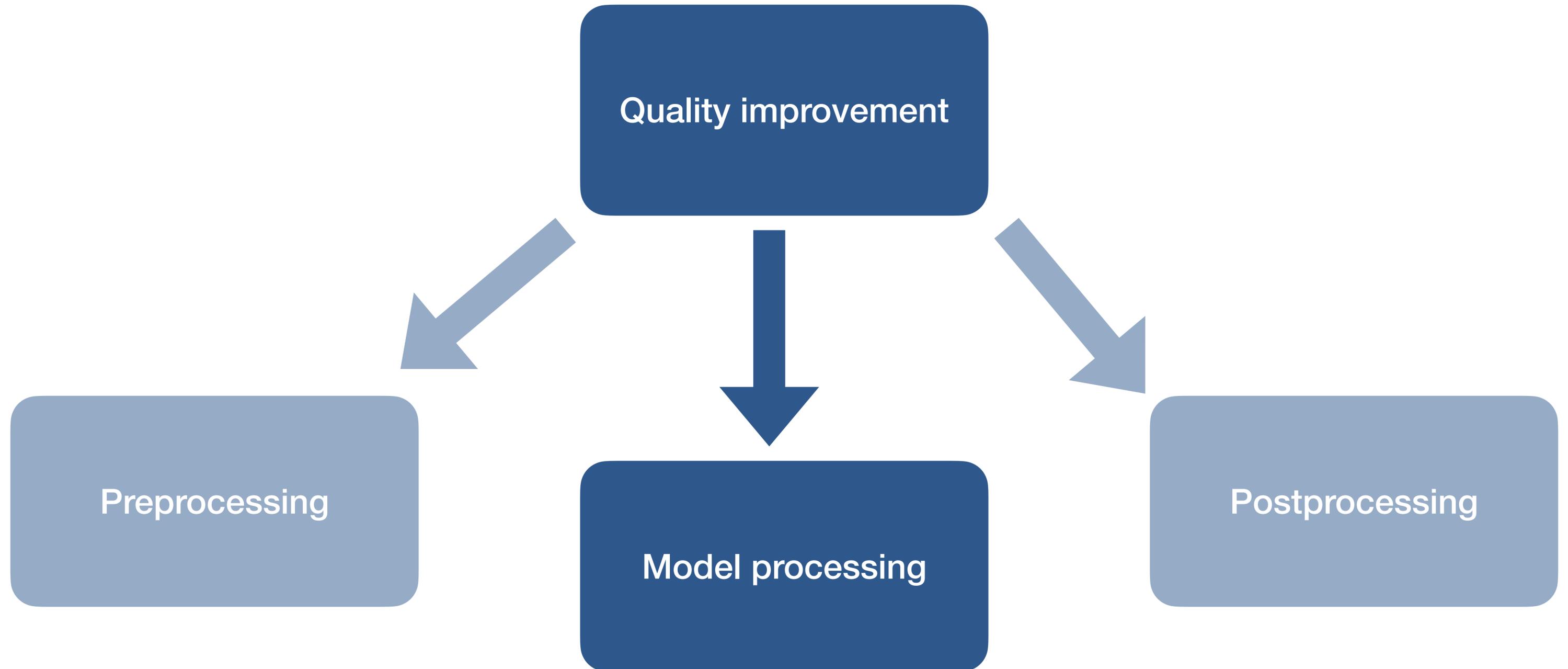
Ухудшает распознавание в
сценах с наложением
голосов

Требует тонкой настройки
для каждого фильма

Voice Activity Detection

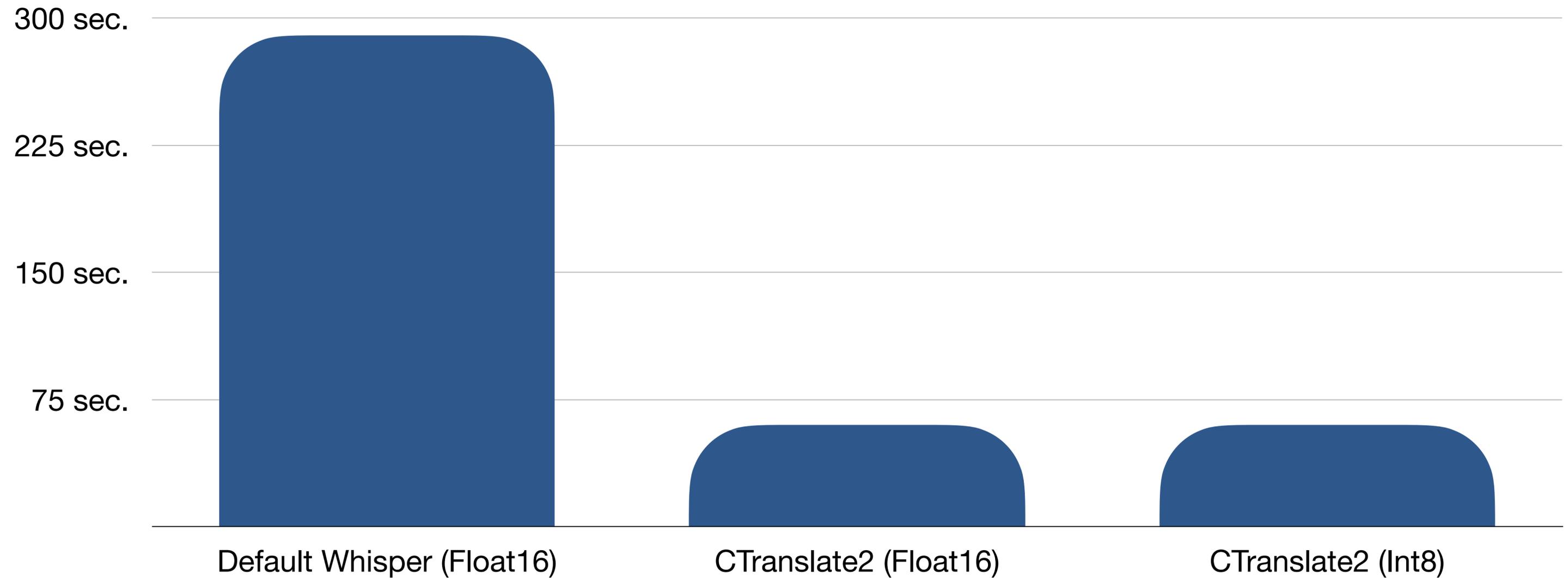
- Сосредоточивает вычислительные ресурсы только на значимых сегментах аудио
- Минимизирует аномалии на фоновых шумах и музыке
- Повышает робастность системы к сложным акустическим условиям
- Упрощает интеграцию с системами анализа и идентификации говорящего

Работа с моделью



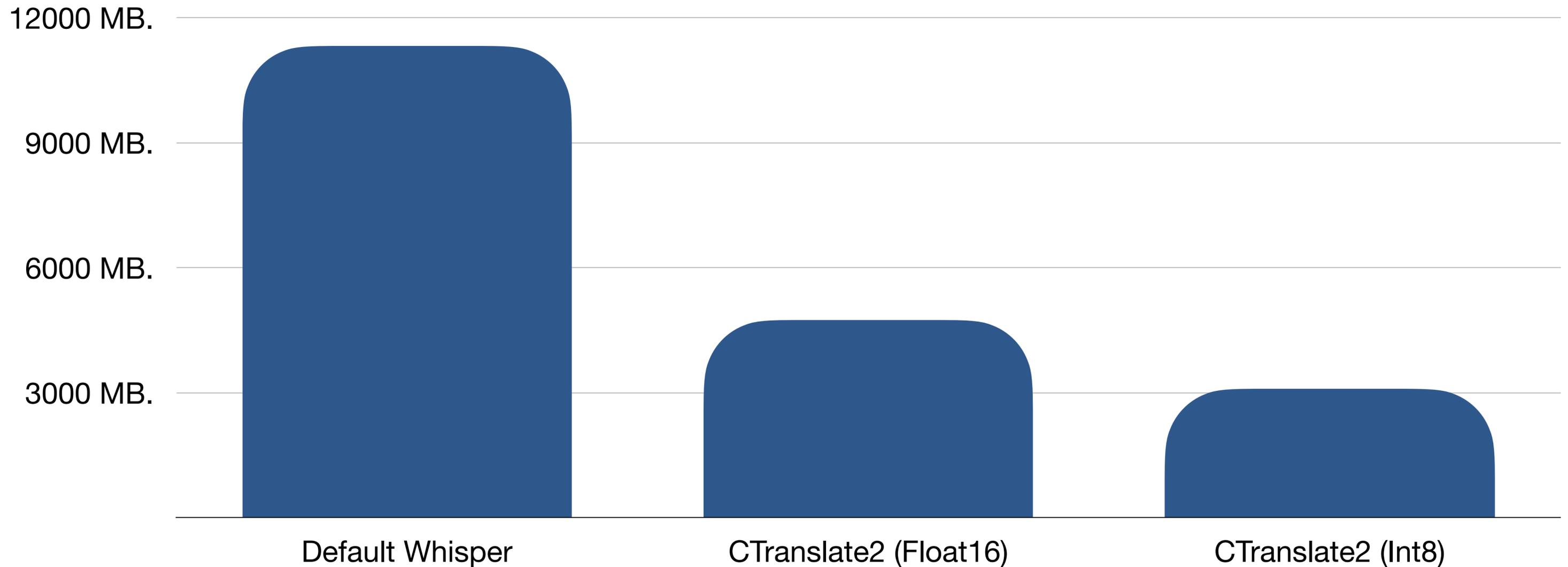
Оптимизация по скорости

Model Large-v2, audio duration ~15 min.

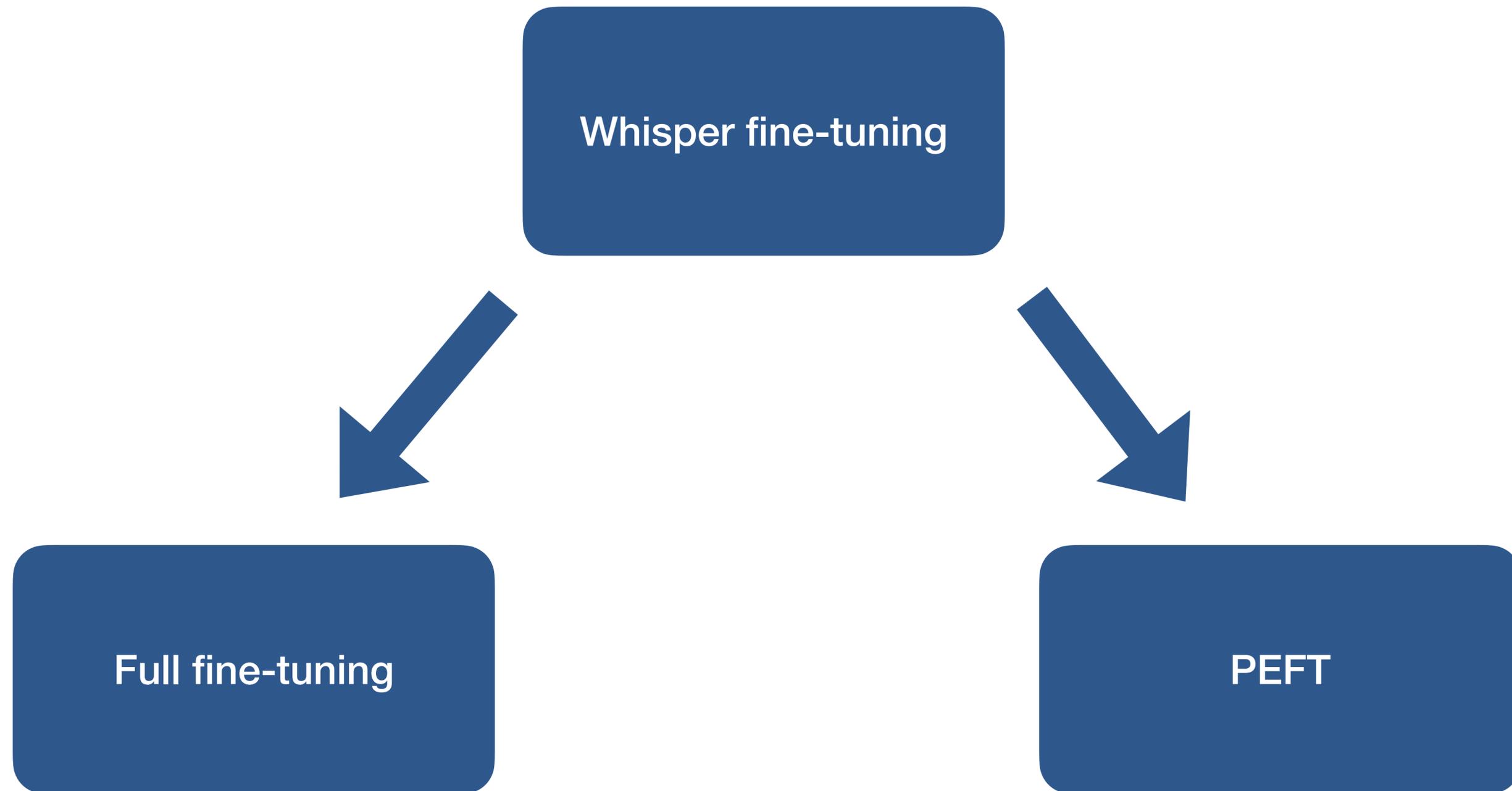


Оптимизация по ресурсам

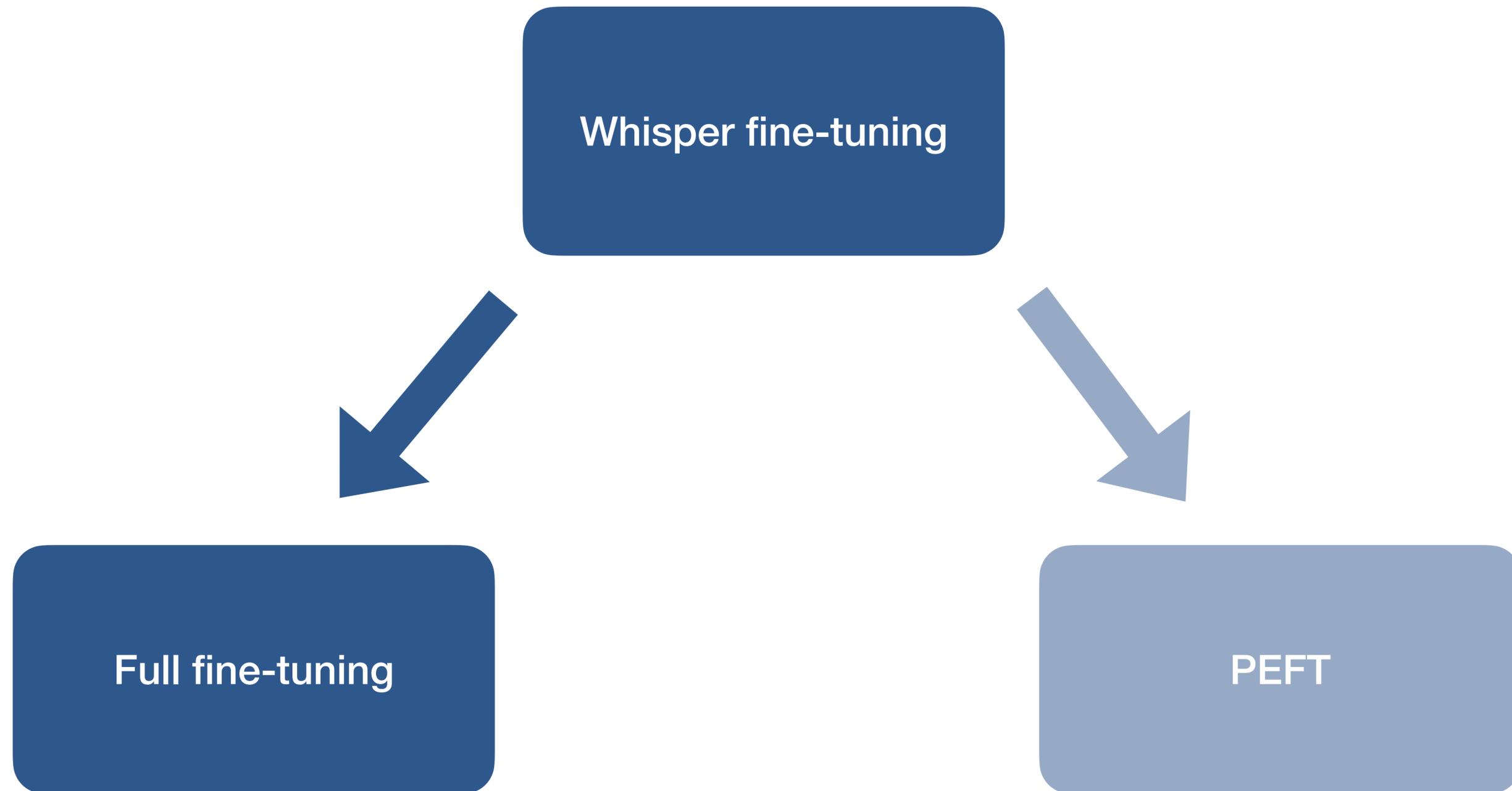
CUDA 12.1 on a NVIDIA Tesla T4 (16 VRAM)



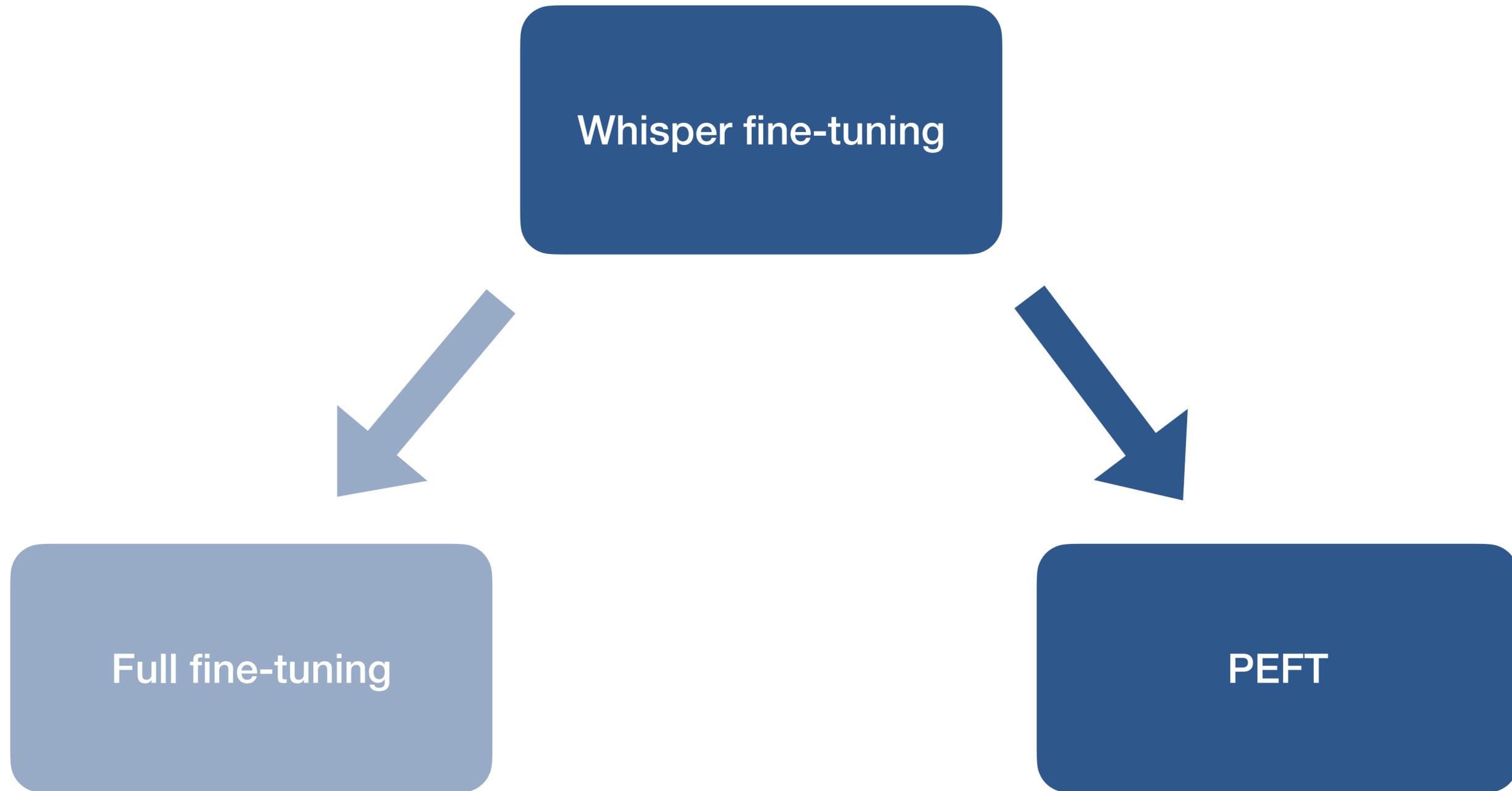
Вариации дообучения



Full fine-tuning

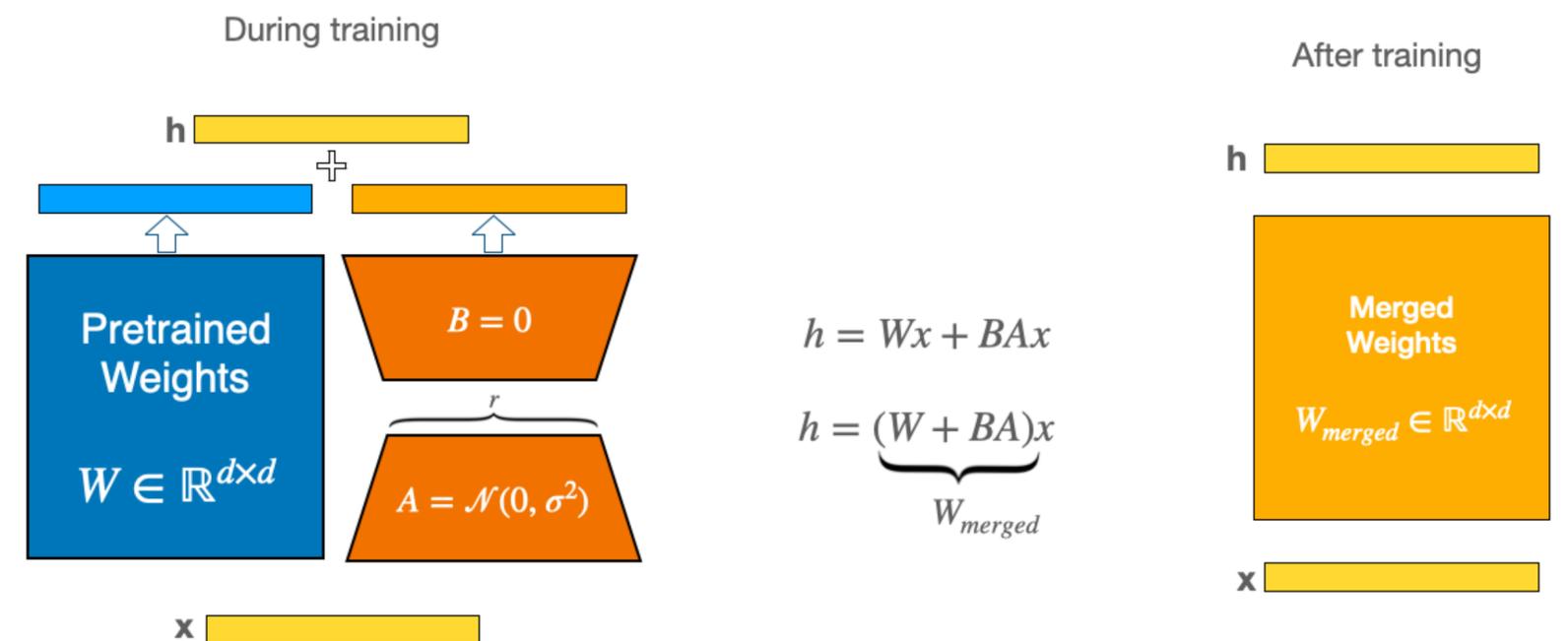


PEFT



LoRA powered by PEFT

- Ускорение времени обучения
- Снижение затрат на вычисления и хранение
- Меньший риск переобучения
- Преодоление катастрофического забывания
- Удобство развертывания

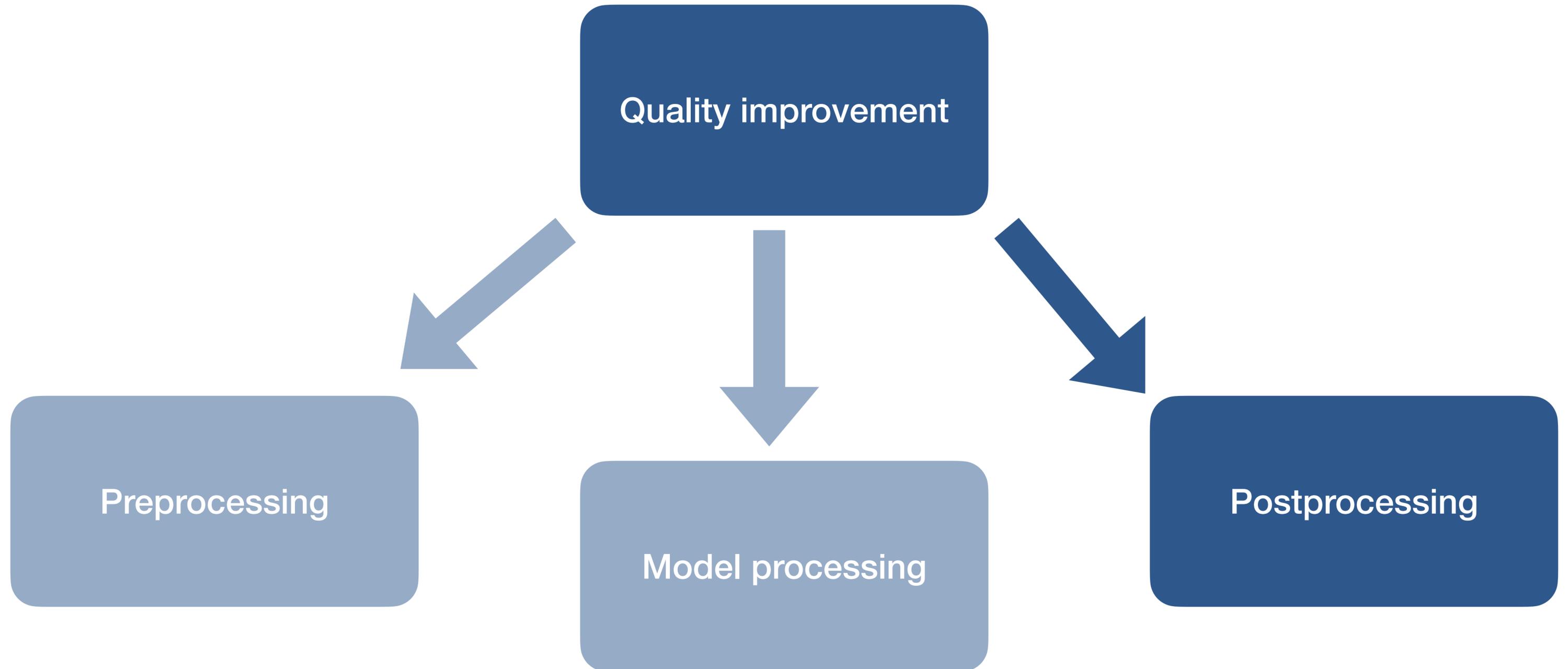


LoRA

Fine-tune pipeline



Постпроцессинг



Alignment pipeline



Good case #1



<https://www.mosfilm.ru/>

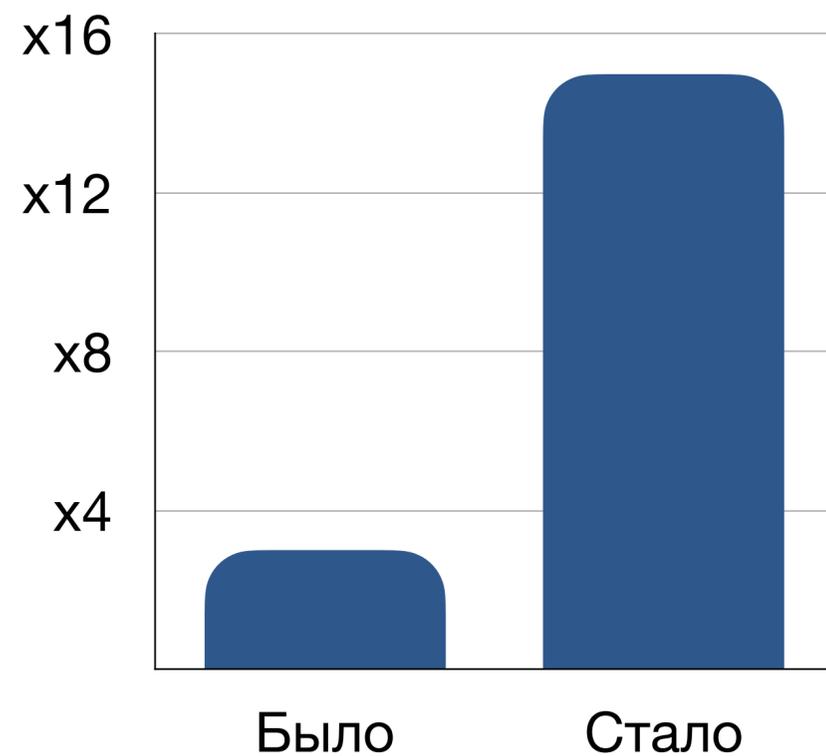
Good case #2



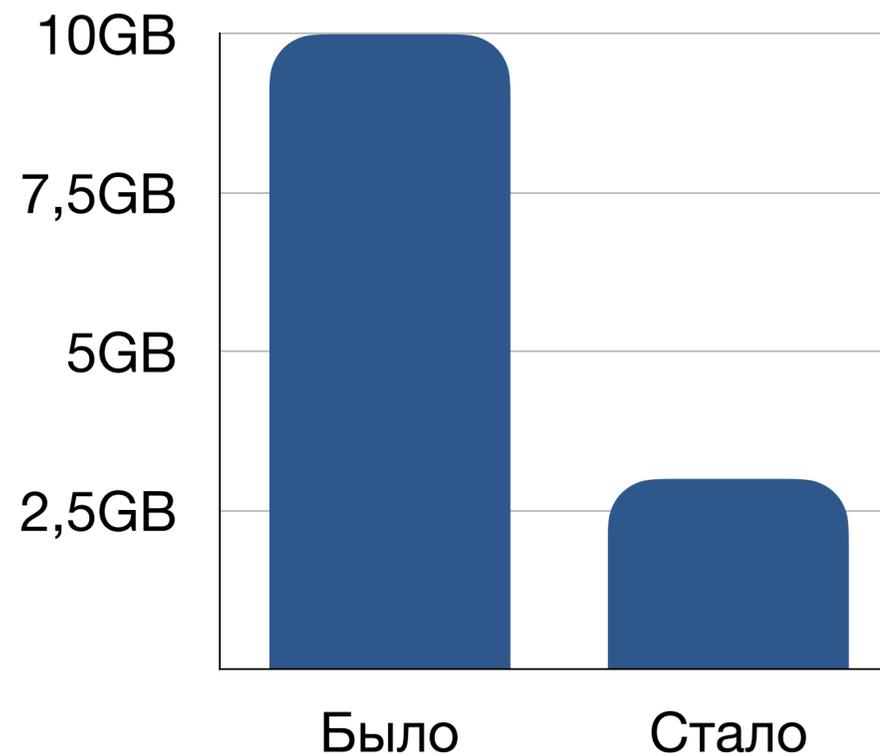
<https://www.mosfilm.ru/>

Итоги

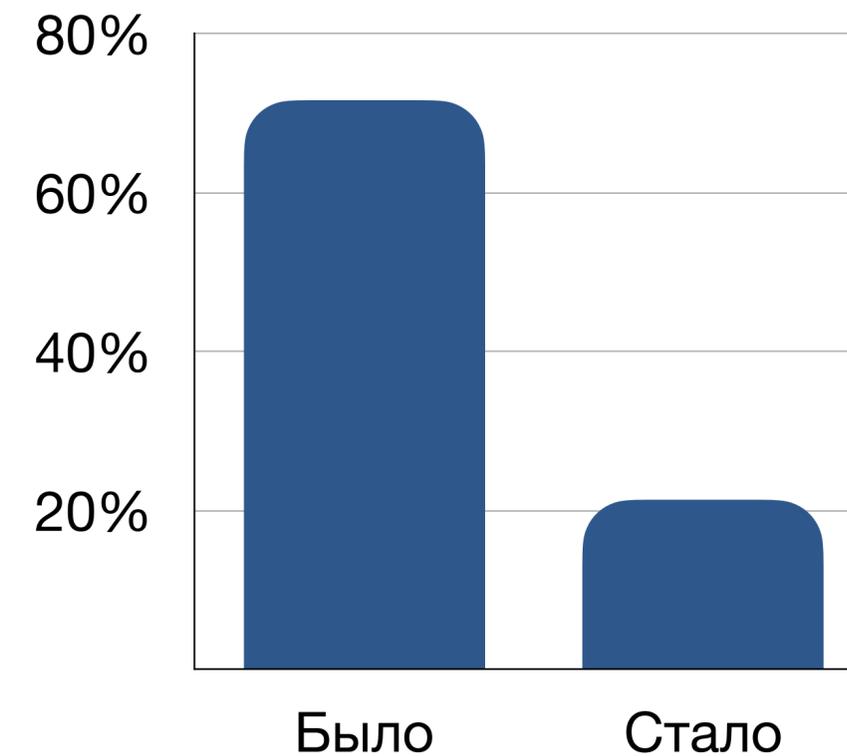
Скорость



VRAM



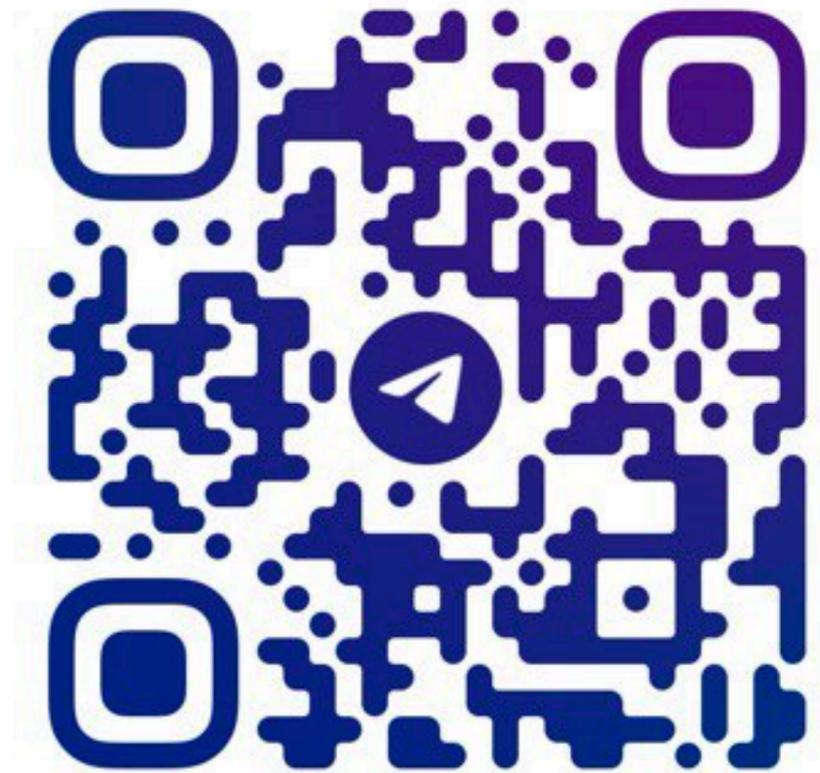
WER



Спасибо за внимание!

Мои контакты:

Почта: dmitry.shatnew@gmail.com



@DMITRII_SHATNEV