

Алгоритм для поиска предложенных скидок в телефонных разговорах с клиентами

Команда



Федор Жилкин

Капитан

tg: @feodoros

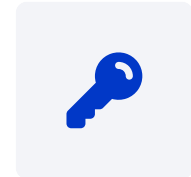


Александр Савельев

Инженер

tg: @a1arick

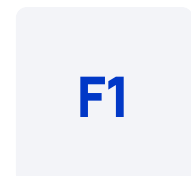
Преимущества решения



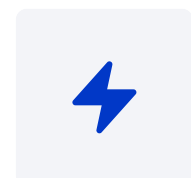
Мы использовали только открытые open-source модели



Мы натренировали на синтетической разметке маленькую быструю BERT-модель для определения наличия скидки по локальному контексту



Метрика F1 – 0.80



Скорость работы всего решения на ЦПУ (Intel(R) Core(TM) i7-10870H CPU @ 2.20GHz) – 10 телефонных разговоров в секунду

Подготовка данных

Шаг 1: Обучение классификатора на оригинальной разметке

Описание датасета

Позитивы: Чанки с ключевыми словами, которые есть в B-discount (1800)
Негативы: Чанки с ключевыми словами, которых нет в B-discount (600)

Аугументация: Back translation

Метрики

deepvk/deberta-v1-base

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.65 | 0.23 | 0.35 |
| 1 | 0.78 | 0.94 | 0.85 |

MyMeetLLM

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.67 | 0.30 | 0.42 |
| 1 | 0.80 | 0.94 | 0.87 |

DeepPavlov/ruBert

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.63 | 0.22 | 0.33 |
| 1 | 0.75 | 0.91 | 0.82 |

cointegrated/rubert-tiny

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.60 | 0.23 | 0.34 |
| 1 | 0.75 | 0.92 | 0.83 |

Подготовка данных

Шаг 2: Обучение классификатора на ответах Chat-GPT 3.5

Описание датасета

Позитивы: Положительно размеченные чанки с ключевыми словами (800)
Негативы: Негативно размеченные чанки с ключевыми словами (1600)

Аугументация: Back translation

Метрики

deepvk/deberta-v1-base

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.42 | 0.42 | 0.42 |
| 1 | 0.79 | 0.79 | 0.79 |

MyMeetLLM

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.37 | 0.99 | 0.54 |
| 1 | 0.99 | 0.39 | 0.55 |

DeepPavlov/ruBert

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.43 | 0.43 | 0.43 |
| 1 | 0.80 | 0.80 | 0.80 |

cointegrated/rubert-tiny

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.41 | 0.41 | 0.41 |
| 1 | 0.78 | 0.78 | 0.78 |

Подготовка данных

Шаг 3: Обучение классификатора на ответах MyMeetLLM

Описание датасета

Позитивы: Положительно размеченные чанки с ключевыми словами (1700)
Негативы: Негативно размеченные чанки с ключевыми словами (700)

Аугументация: Back translation

Метрики

deepvk/deberta-v1-base

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.94 | 0.95 | 0.94 |
| 1 | 0.83 | 0.85 | 0.84 |

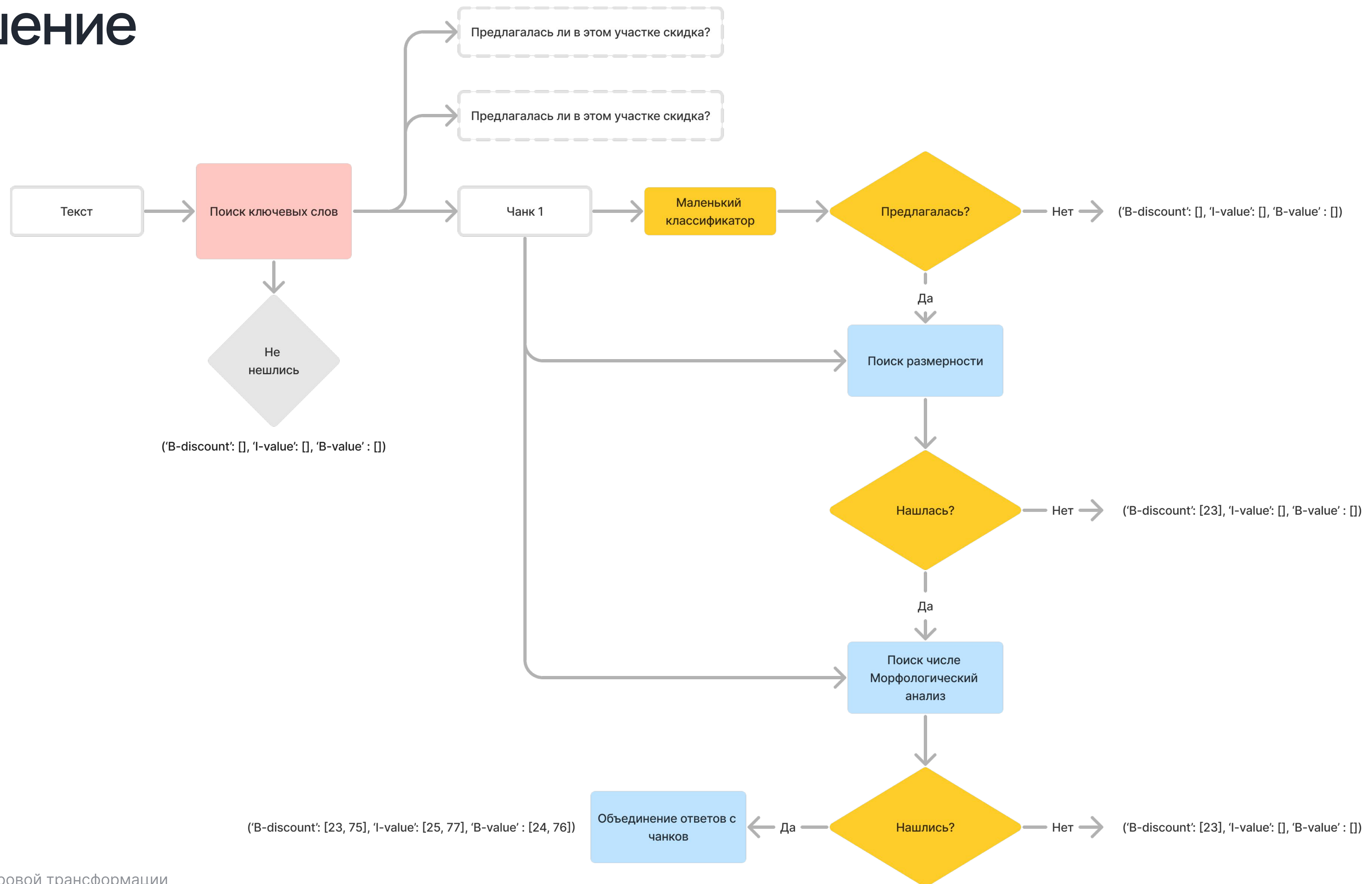
DeepPavlov/ruBert

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.93 | 0.95 | 0.94 |
| 1 | 0.81 | 0.84 | 0.82 |

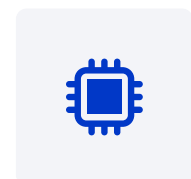
cointegrated/rubert-tiny

| | precision | recall | f1-score |
|---|-----------|--------|----------|
| 0 | 0.92 | 0.91 | 0.91 |
| 1 | 0.78 | 0.81 | 0.79 |

Решение



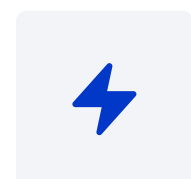
Требовательность к ресурсам



1 vCPU



2 GB RAM



Скорость работы всего решения на ЦПУ – **10 телефонных разговоров в секунду**

Пример работы

Intel(R) Core(TM) i7-10870H CPU @ 2.20GHz

Входной текст:

добрый день NAME меня зовут игнать артема агенты недвижим по ней города хотим завтра с клиентом к вам приехать с утречком в городе парк так с секунда секундан так дикту девятьсот двадцать пять ноль тридцать четыре пятьдесят пять семь два да селен да и еще такой такси а дом один да совершенно связь ну вот как раз здесь если по идее мы должны успеть во сколько там ждать если что тяжело будет наверно получше ждать ну внеси вести есть место ел ввести да ну такси добрый день меня зовут NAME конечно назовите пожалуйста номер клиента NAME северно горький пар такси могу сейчас вызвать а можно будет завтра позвонить тоже по горячей линии мало ли вы разберетесь до метро именно клиент доберется до метро а оттуда хочется на такси доехать вас правильно понимаю все смотрите так я вам назначу встречу завтра а во сколько будет удобно смотрите у нас по времени там определенное время поэтому если придете раньше позже придется подождать в живую очереди но недолго готового времени я не знаю какого как менеджер освободиться в ближайшее то есть ну максимум минут минус финан в десяти утра да записала подскажите поедете наличные автомобиле на общем метро выше я вышел вам адрес клиенту точнее на номер телефона вот а тогда как подъездить к метро набрать на горячую линию и скажите то что вы назначены у вас назначена встреча вот и так смотрите так как вы записались сегодня подойдете завтра за ближайшее посещение офису у вас лично от меня будет скидка у клиента точнее скидка два процента как за

Результат:

```
{'B-discount': [251], 'I-value': [253], 'B-value': [252]}
```

'B-discount' — скидка

'B-value' — два

'I-value' — процента

Время на ЦПУ:

0.09 сек

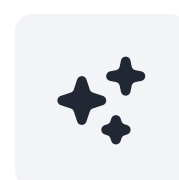


Транскрипт встречи

Разделение на спикеров, главы, пунктуация, имена собственные, англицизмы и 100+ языков

Дообученный Whisper

Самая точная модель для русского языка



Анализ встречи

Протоколирование, выделение задач, кодирование интервью, углубленный анализ речи и семантики

Дообученный Chat-GPT 3.5/4

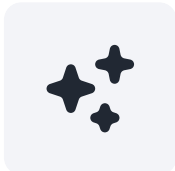
Локальная предобученная LLM модель



Интеграции

Автоматическая запись встреч в крупнейших платформах ВКС и Телеграм-бот





Оркестр ИИ-решений

Whisper ☒

Speaker diarization ☒

neural-chat ☒

Mistral ☒

ChatGPT ☒

LLaMA ☒



Точность транскрибации

96.7%



Более 100 языков

English

Türkçe

Русский

Français

Deutsch

Українська

Español

Italiano

日本語

汉语

Português

+74 други

Спасибо за внимание!

Контакты, Федор Жилкин:
hello@mymmeet.ai
телеграм @feodoros