МИНОБРНАУКИ РОССИИ САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ «ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)

Кафедра математического обеспечения и применения ЭВМ

() 4F,

по научно-исследовательской работе в осеннем семестре 2022 года
Тема: Разработка сервиса электронного протоколирования совещаний и
зашит

Студент гр. 7303	 Петров С.А.
Руководитель	 Шевская Н.В

Санкт-Петербург

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
ПОСТАНОВКА ЗАДАЧИ	4
РЕЗУЛЬТАТЫ РАБОТЫ В ОСЕННЕМ СЕМЕСТРЕ	5
ОПИСАНИЕ ПРЕДПОЛАГАЕМОГО МЕТОДА РЕШЕНИЯ	11
ПЛАН РАБОТЫ НА ВЕСЕННИЙ СЕМЕСТР	13
СПИСОК ЛИТЕРАТУРЫ	14

ВВЕДЕНИЕ

Зачастую в академической и деловой сферах проводятся разнообразные встречи и звонки, как в очном, так и электронном виде (с помощью средств электронной телефонии, таких как Zoom [1], Google Meet [2], и т. д.). Данные встречи, звонки и собрания необходимо протоколировать как для статистики, так и для последующей обработки для заполнения разнообразных документов и форм. В настоящее время для подобного заполнения сотрудникам приходится вручную прослушивать аудиозаписи встреч для получения необходимой информации (например, для заполнения протоколов защит бакалавров). Следовательно, есть необходимость в автоматическом протоколировании аудиозаписей звонков/собраний/встреч с транскрипцией [3], статистикой и фиксацией определенных моментов.

постановка задачи

Целью данного исследования является создание программы для автоматического протоколирования аудио собраний/встреч/звонков.

Предполагаемый функционал:

- Транскрипция аудио;
- Составление статистики по участникам (сколько задали вопросов, как участвовали в дискуссии, как быстро отвечали и т. д.);
 - фиксация вопросов и задач;

Для достижения цели были поставлены следующие задачи:

- 1. Провести обзор технологий и библиотеку для транскрипции речи по аудиоданным.
 - 2. Разработать метод определения вопросов в распознанном тексте.
 - 3. Спроектировать и разработать программу для автоматического заполнения документов.
 - 4. Реализовать и проанализировать разработанное решение Объектом исследования является автоматическое протоколирование. Предметом исследования является обработка аудиоданных собраний.

Практическая ценность работы: разработанное решение позволит сотрудникам вузов или компаний автоматизировать протоколирование рабочих совещаний (например, очных защит дипломов студентов и аспирантов).

РЕЗУЛЬТАТЫ РАБОТЫ В ОСЕННЕМ СЕМЕСТРЕ

Результаты на момент окончания предыдущего семестра.

В предыдущем (весеннем) семестре 2022 года было разработано веб-приложение (прототип) на языке программирования Python 3.10 [4] и библиотеки для распознавания речи Vosk [5], для расшифровки потокового звука с микрофона и определения вопросов (на русском языке).

Функции:

- запись аудиоданных спикеров и подготовка пула спикеров
- онлайн-обработка потока микрофона с транскрипцией и обнаружением говорящего
- экспорт записанных аудиосессий (метаданные + обнаруженные вопросы)

Ссылка на репозиторий - https://github.com/moevm/zoom-transcriber.

План на осенний семестр.

- Работа над устранением сильной зависимости качества распознавания от качества звука (микрофона), необходимость дополнительной программная обработка на клиентской и/или серверной стороне;
- Работы над устранением зависимости качества распознавания от быстроты и качества речи;
 - Распознавание английских терминов в русской речи;
- Отображение плохо распознаваемых слов в клиентской части приложения.

Результаты осеннего семестра.

Работа велась в том же репозитории, в котором велась работа в весеннем семестре 2022 года.

Устранение зависимости качества распознавания от качества входного звука.

В рамках работы по данному направлению было решено добавить средство для очистки входного звука от шума на серверной стороне клиентского приложения. Таким образом, если раньше звуковой отрезок подавался напрямую в библиотеку для распознавания речи, то теперь он будет пропускаться через программную очистку шума.

В качестве программного решения для очистки было принято решение взять Python-библиотеку noisereduce [6] [7]. Данная библиотека вычисляет спектрограмму сигнала (и шумового сигнала, если он подан на вход) и оценивает шумовой порог для каждой полосы частот этого сигнала/шума. Этот порог используется для вычисления маски, которая блокирует шум ниже порога изменения частоты [6].

После внедрения в решение данной библиотеки была проведена серия из 4-х экспериментов.

Каждый эксперимент включал в себя следующие шаги:

- 1. Выбор исходного отрывка из записи аудио-конференции;
- 2. Ручная разметка (распознавание) отрывка для получения текстового эталона;
- 3. Подача звукового отрывка на вход разработанному решению в различных режимах для получения результирующих текстов с распознанной речью
- 4. Сравнение сходства распознанных отрывков с эталоном;
- 5. Получение текстовой разницы распознанных отрывков и текста-эталона с помощью сторонних решений.

Для сравнения сходства текстов были выбраны 4 метрики:

- метрика Python-библиотеки для обработки естественных языков Spacy
 [8] (какой конкретно метод используется в библиотеке, не указано в документации);
- метрика сходства Жаккара отношение длины пересечения множеств к длине объединения множеств. В качестве множеств в данном случае выступают отдельные слова. В применении к текстам была применена в 2-х видах: с первичной лемматизацией слов (jaccard_raw), и без лемматизации (jaccard_process). Выражается в общем виде следующей формулой:

$$J(A,B)=rac{|A\cap B|}{|A\cup B|}=rac{|A\cap B|}{|A|+|B|-|A\cap B|}.$$

• метрика сходства на основе евклидового расстояния между двумя векторами — выражается формулой $\frac{1}{e^{\sqrt{\sum (x_i-y_j)^2}}}$. Для применения в данной

задаче тексты были векторизованы с помощью библиотеки Spacy на основании метода Word2Vec [9]. Таким образом чем ближе расстояние между векторами, тем более 2 текста похожи друг на друга;

• метрика сходства по косинусу – вычисляется по формуле косинуса угла между двумя векторами, подходит в качестве метрики так как угол между двумя параллельным/сходными векторами равен 0. Тексты были векторизованы так же, как и для метрики на основе еквлидового расстояния, с помощью метода Word2Vec;

Эксперименты проводились на записях защит июня 2022 года, предоставленных кафедрой MO ЭВМ.

Результаты проведенных экспериментов не дали возможность выявить значительную разницу между распознаванием речи без очистки от шума и с включенной очисткой от шума, из-за минимального наличия шумов на

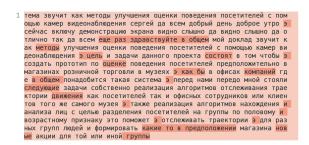
представленных записях. Однако, проведенные эксперименты позволили выявить наиболее эффективную конфигурацию библиотеки для очистки шумов (неправильные конфигурации ухудшали качество распознавания при отсутствии шумов). При наличии шумов данная конфигурация должна обеспечить наилучший результат.

Таблица 1 - результаты метрик для одного из экспериментов

Эксперимент/Метрика	spacy	jaccard_raw	jaccard_proc ess	euclidean	cosine
result_only_medfilt	0.991160949	0.694736842	0.761904761	0.716367811	0.991436979
	9029664	1052632	9047619	4524413	4842874
result_no_medfilt_no_sample_nft	0.991737237	0.752688172	0.797619047	0.729512428	0.991863927
t4096_timemask64	4288371	0430108	6190477	0737039	7942686
result_no_medfilt_with_sample_n	0.991473991	0.726315789	0.788235294	0.725283639	0.991560848
ftt4096_timemask64	6524232	4736842	117647	3861609	3163058
result_no_medfilt_no_sample_de fault	0.991274413	0.708333333	0.773809523	0.717483152	0.991236572
	1028703	3333334	8095238	8613965	174088
result_no_medfilt_with_sample_d efault	0.991090301	0.708333333	0.773809523	0.714997942	0.991210154
	4787334	3333334	8095238	2902976	6229997
result_raw	0.991737237	0.752688172	0.797619047	0.729512428	0.991863927
	4288371	0430108	6190477	0737039	7942686

Из таблицы видно, что наивысшие коэффициенты сходства по метрикам наблюдаются у результатов без очистки от шума (results_raw) и с определенной конфигурацией (result_no_medfilt_no_sample_nftt4096_timemask64).

Также были получены текстовые сравнения отрывков по всем экспериментам.



1 тема звучит как методы улучшения оценки поведения посетителей с пом ощью камер видеонаблюдения сергей да всем добрый день доброе утро с ейчас включу демонстрацию экрана видно слышно да видно слышно да от лично так да всем через здравствуйте потому что мой доклад звучит к ак метода улучшения оценки поведения посетителей с помощью камер ви деонаблюдения цели и задачи данного проекта состоит в том чтобы соз дать прототил по для оценки поведения посетителей и предположительн о в магазинах розничной торговли в музеях кабы в офисах компании гд е понадобится такая система перед нами передо мной стояли следующей задачи собственно реализация алгоритмов отслеживания траектории передвижения как посетителей так и офисных сотрудников или клиентов то с же самого музея также реализация алгоритмов нахождения анализа л иц с целью разделения посетителей на группы по половому возрастному признаку это поможет отслеживать траектории для разных групп людей и формировать какие—то в принципе предположения магазина новой акци и для той или иной

Рисунок 1 – пример текстового сравнения

Результаты по всем экспериментам и ipynb-ноутбук с реализацией метрик сходства и генерацией результатов были загружены на Google Drive и доступны по следующей ссылке - результаты экспериментов по очистке шума.

Отображение плохо распознаваемых слов в клиентской части приложения.

Из диалога с научным руководителем было принято расширить функционал клиентской части путем добавления выделения слов, которые были "плохо" распознаны используемой библиотекой распознавания речи.

Библиотека в качестве результата выдает коэффициент уверенности для каждого распознанного слова. Данный коэффициент не использовался предыдущем семестре, слова просто конкатенировались в клиентской части без дополнительной смысловой нагрузки. В текущей работе было принято решение использовать данные коэффициенты — слова отправляются на клиент вместе с соответствующими коэффициентами. Клиент (html + javascript + css) рассматривает каждое слово в отдельности вместе с его коэффициентом, и, если коэффициент не достигает настроенного порога (используется значение 0.7), использует html-разметку чтобы выделить слово для пользователя.



Рисунок 2 — отображение распознанной речи на клиенте с выделением слов с низким коэффициентом распознавания

Другие результаты.

Был расширен формат экспортируемых аудио-сессий, путем добавления колонки с "сырыми" данными распознанных слов (слово + коэффициент распознавания) в формате JSON, которые можно загрузить программными средствами и использовать для дальнейшего анализа.

ОПИСАНИЕ ПРЕДПОЛАГАЕМОГО МЕТОДА РЕШЕНИЯ

Предполагаемое решение представляет из себя клиент-серверное приложения. Сервер представляет из себя веб-приложение на языке Python, клиент – статические файлы HTML + JavaScript/CSS. Данная архитектура упрощает развертывание приложения, а также позволяет сосредоточиться на непосредственно разработке метода распознавания вместе траты времени на разработку клиентских приложений под разные платформы. Передача аудиоданных с клиента осуществляется при помощи языка JavaScript и WebSockets [10], поток аудиоданных с микрофона получается при помощи Web Audio API [11]. Серверная часть принимает поток аудиоданных и обрабатывает его: очищает от шума библиотекой noisereduce и посылает для распознавания другому развернутому веб-серверу, который использует библиотеку для Vosk распознавания речи ДЛЯ транскрибации полученных Использование программного решения для распознавания речи в качестве отдельного компонента архитектуры позволяет легко изменить/заменить его в случае необходимости, без изменения кода клиента и основного приложения.

Транскрибированные фразы в виде текста вместе с метаданными (временные метки, метка вопроса, коэффициент точности распознанной фразы) сохраняются в NoSQL-базу MongoDB [12]. База была выбрана из-за отсутствия необходимости хранить нормализованные данные.

Архитектура приложения позволяет легко изменять и расширять правила для определения вопросов, а также дает возможность внедрения дополнительной обработки аудиоданных и распознанного текста.

Приложение позволяет осуществлять экспорт записанных и обработанных аудиоданных. При экспорте пользователь получает zip-архив, содержащий 2 csv-файла:

• records.csv – содержит распознанные фразы спикеров, с временными метками, показателем точности, меткой вопроса, данными по отдельным словам в формате JSON (упакованном в строку);

• metadata.csv — метаданные обработанного аудио/конференции (название, временные метки, распознанные спикеры, и т.п.)

Выгруженный на Google Drive пример экспортированных результатов обработки можно посмотреть здесь.

ПЛАН РАБОТЫ НА ВЕСЕННИЙ СЕМЕСТР

Направления для улучшения/дальнейшего исследования:

- исследование подходов для улучшения распознавания вопросов спикеров, расширение/улучшение набора правил для распознавания вопросов;
- добавление возможности для редактирования получившейся транскрипции/расшифровки совещания/конференции, предложение альтернатив для распознанных слов с низким коэффициентом распознавания;
- оптимизация разработанного решения, в частности взаимодействия между основным сервером и сервером распознавания речи;

СПИСОК ЛИТЕРАТУРЫ

- 1. Официальный сайт Zoom: [сайт]. URL: https://zoom.us
- Обзор возможностей Google Meet: [сайт]. URL:
 https://apps.google.com/intl/ru/intl/ru_ALL/meet/how-it-works/
- 3. Кибрик А. А., Подлесская В. И. К созданию корпусов устной русской речи: принципы транскрибирования //Научно-техническая информация (серия 2). 2003. Т. 6. С. 5-11.
- 4. Официальный сайт языка программирования Python: [сайт]. URL: https://www.python.org/
- 5. Официальная документация библиотеки для распознавания речи Vosk: [сайт]. URL: https://alphacephei.com/vosk/
- 6. Репозиторий библиотеки noisereduce: [сайт]. URL: https://github.com/timsainb/noisereduce
- 7. Sainburg T., Thielk M., Gentner T. Q. Finding, visualizing, and quantifying latent structure across diverse animal vocal repertoires //PLoS computational biology. 2020. T. 16. №. 10. C. e1008228.
- 8. Документация библиотеки Spacy для обработки естественных языков: [сайт]. URL: https://spacy.io/
- 9. Goldberg Y., Levy O. word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method //arXiv preprint arXiv:1402.3722. 2014.
- 10. Руководство по технологии WebSockets: [сайт]. URL: https://developer.mozilla.org/ru/docs/Web/API/WebSockets_API
- 11.Документация Web Audio API: [сайт]. URL: https://developer.mozilla.org/en-US/docs/Web/API/Web_Audio_API
- 12.Официальная документация MongoDB: [сайт]. URL: https://docs.mongodb.com/manual.