

Министерство науки и высшего образования Российской Федерации
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА

**РАЗРАБОТКА АЛГОРИТМОВ РАБОТЫ С ФОРМАЛЬНОЙ МОДЕЛЬЮ
ДИАЛОГОВ, ПРЕДСТАВЛЕННЫХ В ВИДЕ ГРАФОВ**

Автор: Савон Юлия Константиновна _____

Направление подготовки: 01.03.02 Прикладная
математика и информатика

Квалификация: Бакалавр

Руководитель ВКР: Ульяновцев В.И., к.т.н. _____

Санкт-Петербург, 2020 г.

Обучающийся Савон Юлия Константиновна
Группа М3437 Факультет ИТиП

Направленность (профиль), специализация
Математические модели и алгоритмы в разработке программного обеспечения

Консультанты:

а) Ступаков И.М., канд. тех. наук, доцент

ВКР принята «_____» _____ 20__ г.

Оригинальность ВКР _____%

ВКР выполнена с оценкой _____

Дата защиты «26» июня 2020 г.

Секретарь ГЭК Павлова О.Н.

Листов хранения _____

Демонстрационных материалов/Чертежей хранения _____

Министерство науки и высшего образования Российской Федерации
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

УТВЕРЖДАЮ

Руководитель ОП
проф., д.т.н. Парфенов В.Г. _____
« ____ » _____ 20__ г.

ЗАДАНИЕ
НА ВЫПУСКНУЮ КВАЛИФИКАЦИОННУЮ РАБОТУ

Обучающийся Савон Юлия Константиновна

Группа М3437 **Факультет** ИТиП

Квалификация: Бакалавр

Направление подготовки (специальность): 01.03.02 Прикладная математика и информатика

Направленность (профиль): Математические модели и алгоритмы в разработке программного обеспечения

Тема ВКР: Разработка алгоритмов работы с формальной моделью диалогов, представленных в виде графов

Руководитель Ульяновцев В.И., к.т.н., доцент факультета информационных технологий и программирования Университета ИТМО

2 Срок сдачи студентом законченной работы до: «20» июня 2020 г.

3 Техническое задание и исходные данные к работе

Требуется провести исследование и разработать набор алгоритмов для выявления отвлечений в графовой модели для голосовой диалоговой системы. Алгоритм принимает набор диалогов, в размере нескольких тысяч. Предварительно выстроенный граф и кластеризацию для фраз оператора.

На выходе ожидается получить набор отвлечений и перестроенный граф. В качестве метрики качества будет использоваться сравнение с уже существующими графами, которые создавались вручную.

4 Содержание выпускной квалификационной работы (перечень подлежащих разработке вопросов)

Пояснительная записка должна описывать предметную область диалогов представленных в виде графов. Так же формулировать цель и задачу выделения отвлечений, содержать описание алгоритмов их поиска. Должны быть описаны сложности и методы их разрешения, если они возникали. Кроме того должны быть приведены примеры работы алгоритмов и сравнение с существующими решениями. Кроме того пояснительная записка должна содержать описания задач из смежных областей и их то, как эти задачи связаны с задачей решаемой в работе.

5 Перечень графического материала (с указанием обязательного материала)

Графические материалы и чертежи работой не предусмотрены

6 Исходные материалы и пособия

- а) Среда разработки Visual Studio Code;
- б) ГОСТ 7.32–2001 «Система стандартов по информации, библиотечному и издательскому делу. Отчет о научно-исследовательской работе. Структура и правила оформления».

7 Дата выдачи задания «01» сентября 2019 г.

Руководитель ВКР _____

Задание принял к исполнению _____ «01» сентября 2019 г.

Министерство науки и высшего образования Российской Федерации
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

АННОТАЦИЯ
ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ

Обучающийся: Савон Юлия Константиновна

Наименование темы ВКР: Разработка алгоритмов работы с формальной моделью диалогов, представленных в виде графов

Наименование организации, где выполнена ВКР: Университет ИТМО

ХАРАКТЕРИСТИКА ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ

1 Цель исследования: Разработка алгоритма выделяющего отвлечение в диалоге, представленном в виде графа.

2 Задачи, решаемые в ВКР:

- а) Разработать алгоритмы выделения отвличий;
- б) Реализовать описанные алгоритмы;
- в) Перестроить граф в соответствии с используемой моделью в компании;
- г) Проанализировать результаты работы алгоритмов;
- д) Интегрировать разработки в инфраструктуру компании.

3 Число источников, использованных при составлении обзора: 0

4 Полное число источников, использованных в работе: 11

5 В том числе источников по годам:

Отечественных			Иностранных		
Последние 5 лет	От 5 до 10 лет	Более 10 лет	Последние 5 лет	От 5 до 10 лет	Более 10 лет
0	2	1	2	2	4

6 Использование информационных ресурсов Internet: да, число ресурсов: 1

7 Использование современных пакетов компьютерных программ и технологий:

Пакеты компьютерных программ и технологий	Раздел работы
Интегрированная среда разработки PyCharm	Глава 2.2
Программное обеспечение для автоматизации развертывания и управления приложениями в средах с поддержкой контейнеризации Docker	
Распределённая система контроля версий Git	Глава 2.2

8 Краткая характеристика полученных результатов

9 Гранты, полученные при выполнении работы

10 Наличие публикаций и выступлений на конференциях по теме выпускной работы
По теме этой работы был сделан доклад на Конгрессе Молодых Ученых.

Обучающийся Савон Ю.К. _____

Руководитель ВКР Ульянцев В.И. _____

« _____ » _____ 20__ г.

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	5
1. Обзор диалоговой модели	7
1.1. Разговорная диалоговая система.....	7
1.2. Диалоговый менеджер.....	7
1.3. Анализ задачи выделения отвлечений	10
1.4. Постановка задачи	11
1.5. Выводы по первой главе.....	11
2. Описание алгоритмов поиска отвлечений и рекластеризации	12
2.1. Процесс разработки	12
2.2. Задача выделения отвлечений	12
2.3. Подход поиска отвлечений с большим количеством рёбер входящих в одну вершину	13
2.4. Подход поиска отвлечений с поиском циклов.....	15
2.5. Выборка хороших кластеров	16
2.6. Функция сравнения сообщений	17
2.7. Слияние кластеров операторских сообщений.....	18
2.8. Оценка для сравнения графов	19
3. Результаты экспериментов.....	21
3.1. Используемые данные.....	21
3.2. Результаты улучшения кластеризации	23
3.3. Результаты работы алгоритма поиска отвлечений по рёбрам.....	25
3.4. Результаты алгоритма поиск циклов.....	25
3.5. Программная реализация.....	29
ЗАКЛЮЧЕНИЕ	30
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	31
ПРИЛОЖЕНИЕ А. Исходный код на языке Python	33
ПРИЛОЖЕНИЕ Б. Диплом Конгресса Молодых Учёных	49

ВВЕДЕНИЕ

Эта работа посвящена исследованию в области диалоговых систем.

Диалоговая система – это алгоритм, который умеет принимать участие в диалоге на естественном языке и использует правила общения между людьми.

В качестве примера диалоговых систем можно привести:

- Чат-ботов
- Голосовых помощников
- Автоответчиков в колл-центрах

Такие диалоговые системы могут быть как довольно простыми (например, чат-бот отвечающий на заранее известный набор команд), так и сложными (бот, отвечающий на вопрос, сформулированный на естественном языке и возвращающий некоторую информацию из базы знаний).

В последнее время набрали популярность технологии распознавания и генерации речи, которые позволили создавать диалоговые системы, ведущие голосовой разговор. Данная работа рассматривает алгоритмы именно для голосовых диалогов. Пример решения такой диалоговой системы можно рассмотреть в статье [10].

Такие звонки, с одной стороны, должны быть не отличимы от звонков человека, с другой – они должны придерживаться некоторого сценария.

Сценарий диалога с оператором – некоторый алгоритм, предоставленный человеку, который звонит по некоторому набору номеров (либо же принимает входящий звонок или общается посредством программного обеспечения для аудиосвязи). Целью сценария обычно является получение или донесение до клиента информации.

Несмотря на то, что под сценарием диалога понимается некоторый алгоритм, необходимо понимать что для человека и диалоговой системы это принципиально разные сущности. Между скриптом¹ и алгоритмом, с точки зрения набора действий для машины есть большая разница.

Для человека это скорее структурированный список вопросов которые он должен задать и некоторая дополнительная база знаний с ответами на нестандартные вопросы. Кроме того человек может помнить некоторые фак-

¹Здесь и далее в тексте слова «скрипт» и «сценарий» будут использоваться как синонимы

ты и выдавать их дополнительно в зависимости от контекста. Он сам умеет обрабатывать ситуации, такие как отвлечение от основных вопросов или переспрашивание.

Для скрипта же, реакцию на любую фразу человека надо прописывать, все возможные данные хранить и обновлять. Кроме того, есть требование поддерживать этот скрипт доступным для восприятия человеком (например лингвистом), поскольку возникает необходимость в ручном анализе и редактировании. В данном случае такой скрипт будет представлен в виде графа с дополнительной информацией.

Голосовая диалоговая система — программа, которая используя сценарий умеет проводить диалог с клиентом, интерпретировать и записывать информацию полученную от клиента, а так же состояния завершенного разговора. Кроме того, она умеет отвечать на заранее прописанный в скрипте набор вопросов и возвращаться обратно к диалогу.

На данный момент существуют графы для диалогов, которые создаются вручную. Но писать их долго, а продумывать все важные случаи реакций сложно и трудозатратно.

Помимо этого, хочется иметь возможность усложнять вариативность диалогов, делая их более похожими на процесс общения двух людей. В связи с этим, ставится глобальная задача по созданию графа из набора проведенных диалогов.

Достаточно часто возникает ситуация, когда некоторый общий вопрос может возникнуть в любом месте диалога (например, «А что это за компания?»). Таким образом, было решено разделить их в отдельные подграфы и сделать возможность переходить в них при некоторых условиях из каждой вершины. В дальнейшем такие случаи будут называться **отвлечениями**.

В работе будут рассмотрены различные алгоритмы автоматического поиска таких отвлечений.

ГЛАВА 1. ОБЗОР ДИАЛОГОВОЙ МОДЕЛИ

1.1. Разговорная диалоговая система

Данная работа выполнялась в компании Dasha.AI Inc. В ней будут рассматриваться интеллектуальные диалоговые системы. Введем несколько определений.

Разговорная диалоговая система – система, позволяющая пользователю-человеку получать доступ к информации и сервисам доступным на компьютере или в сети Интернет, используя разговорный язык как средство взаимодействия, согласно [7].

Разговорные диалоговые системы используют распознавание речи для преобразования фраз человека в формат, удобный для использования диалоговым менеджером.

1.2. Диалоговый менеджер

Диалоговый менеджер – компонент диалоговой, системы ответственный за текущее состояние и ход диалога, согласно[9]. В основном, диалоговые менеджеры оперируют предобработанными текстами. У диалогового менеджера так же есть состояние, которое может хранить в себе например последний неотвеченный вопрос или историю диалога.

Диалоговый менеджер в том числе можно представить в виде автомата, пример которого можно увидеть на Рисунке 1:

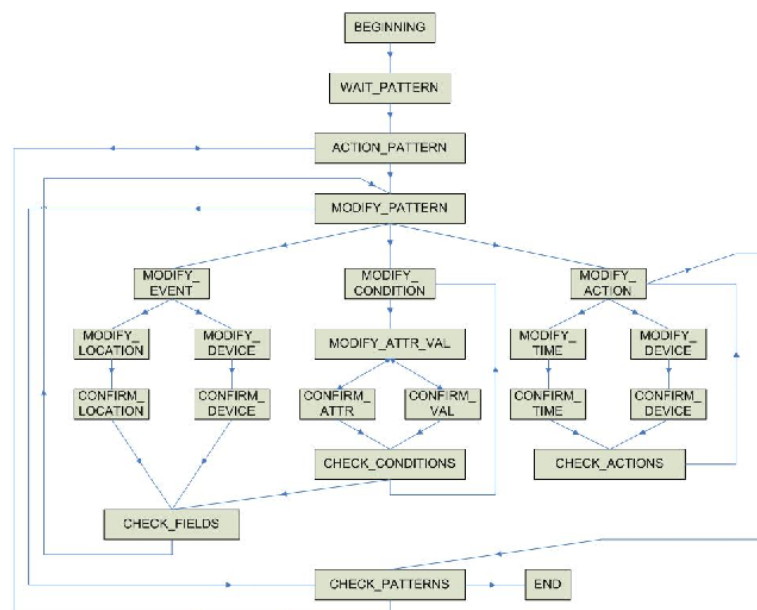


Рисунок 1 – Пример автомата для диалогового менеджера

В этом примере, в зависимости от текущего состояния и информации, предоставленной пользователем – система пытается принять решение о том, какое действие хочет выполнить пользователь. Эта система использует не только текущий ответ, а так же информацию о контексте. Кроме того, если системе не хватает информации, она запрашивает её у пользователя. Рисунок и описание взяты из [4].

Основные способы создания моделей диалоговых менеджеров:

- Составление вручную.
- Генерация модели менеджера с помощью машинного обучения.

Плюсами первого подхода можно назвать прозрачность процесса перевода имеющегося скрипта в модель. Серьёзным минусом создания модели диалогового менеджера вручную является необходимость каждый раз полностью создавать структуру диалога, что очень ресурсозатратно.

Плюсом генерации модели менеджера с помощью машинного обучения является невысокая трудозатратность на каждую новую систему. Минусом такого подхода можно назвать необходимость в большом количестве данных для обучения и зависимости от качества исходных данных.

Пример модели диалогового менеджера, улучшенного с помощью методов машинного обучения представлен в статье [6]. В указанной работе байесовский подход к изучению пользовательской модели применялся одновременно с политикой менеджера диалогов. Данное улучшение позволило организовать управление инвалидным креслом с помощью голоса, посредством инструкций более похожих на естественный язык.

В основном для промышленного использования диалоговых систем в компаниях используются диалоговые менеджеры на основе правил, составленных вручную. Создание одного такого диалогового менеджера может требовать много ресурсов, которые при таком подходе нельзя значительно оптимизировать.

В этой работе рассматривается построение на основе графов.

Задача, решаемая командой состоит в следующем: необходимо по набору диалогов восстановить диалоговую систему. Данная работа решает одну из подзадач.

Диалоговая система в данном случае представлена в виде графа, но так же существуют и другие представления. На Рисунке 2 представлен пример упрощенной модели, с несколькими вариантами ответов человека.

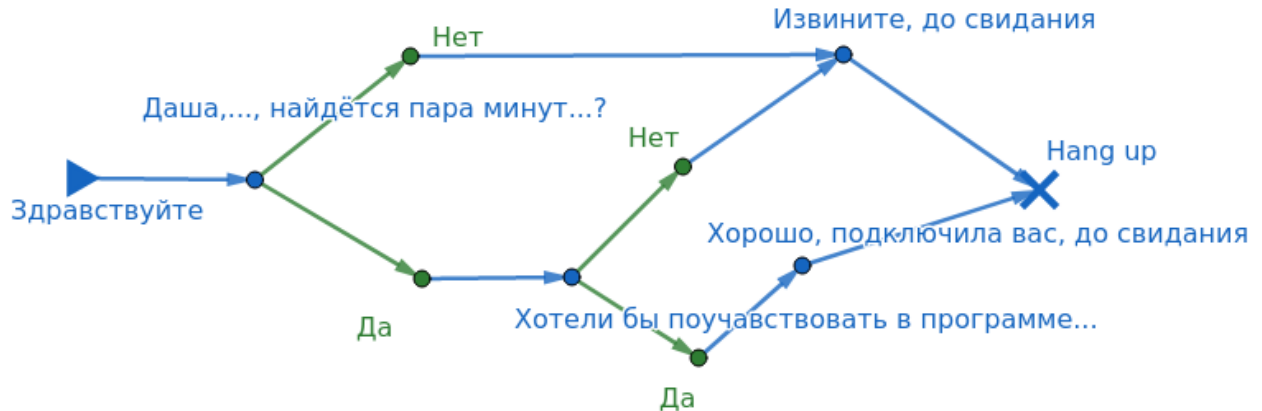


Рисунок 2 – Простой пример модели

На данный момент, кроме восстанавливаемой модели – разработаны графовые модели, которые составляются вручную. В упрощенном виде, в этих моделях вершинами являются фразы диалоговой системы, а рёбра ассоциированы с группами возможных ответов человека. То есть, при различных ответах человека происходит переход в разные вершины.

Особенностью данной структуры, которую важно отметить, является то, что помимо обычных переходов существуют также скрытые переходы, которые по умолчанию могут встретиться в любом месте. В качестве примера можно привести вопрос: «А откуда вы взяли мой номер телефона?». При звонке люди могут задать этот вопрос не сразу.

Для того, чтобы не рассматривать каждый такой случай при переходе из каждой вершины, выделяются **отвлечения**.

Отвлечение – это вопрос, который может быть задан в любом месте диалога. Кроме того, к одному из подтипу отвлечений относятся случаи, когда оператор не услышал фразу или вопрос и вынужден переспросить. В зависимости от модели, это может выделяться в специальную сущность или являться отвлечением.

В задаче восстановления графа по набору диалогов используется видоизменённая модель. В ней вершинами являются, как кластеры фраз оператора,

так и кластеры фраз человека. А ребро соответствует наличию перехода между соответствующими кластерами в диалоге.

Кластеризация текстов (фраз)[8] – автоматическая группировка текстовых документов, например веб-страниц, электронных писем, человеческих фраз, основанная на схожести содержимого. В качестве входных данных алгоритмы принимают наборы фраз и количество желаемых кластеров. В качестве выходных данных, алгоритм возвращает фразы, сгруппированные в указанное количество кластеров – кластеры K_1, K_2, \dots .

1.3. Анализ задачи выделения отвлечений

На вход подаётся набор диалогов, по которым нужно получить граф для диалоговой системы, которая бы могла проводить аналогичные диалоги.

Граф, если получать его путём обычной кластеризации операторских фраз, получается избыточно большим. В нём плохо видно структуру, его сложно анализировать.

Конечной целью является научиться строить графы, аналогичные тем, которые создаются другой командой в компании вручную. Исходя из этого, появляется требование привести его в состояние, когда человек мог бы изучить и проанализировать такой восстановленный граф.

То есть, так же как и в графе создаваемом вручную, появляется необходимость выделять отвлечения. В этом случае, его структура становится более удобной для анализа. Алгоритм выделения, в частности, позволяет анализировать особенности графа непосредственно во время процесса выделения отвлечений.

Как следствие, особенности структуры отвлечений позволяют грамотно обрабатывать сценарии, которых не было в изначальном наборе диалогов. Если такие отвлечения не выделить, то в случае вопроса, не предусмотренного в этом месте диалоговая система либо даст некорректный ответ, либо несвоевременно закончит работу.

Хорошая кластеризация текстов является важным предварительным шагом для алгоритма, поскольку любой алгоритм выделения отвлечений непосредственно опирается на информацию, полученную в результате кластеризации. Полезно, при этом учитывать имеющуюся информацию о контексте, то есть фразы до текущей и после.

1.4. Постановка задачи

Целью данной работы является реализация возможности поиска отвлечений в алгоритме восстановления графа по набору диалогов.

Отвлечения могут встретиться в любом месте, таким образом после выделения и перестроения графа, граф будет более структурирован, что позволит строить более сложные конструкции для диалогов, которые всё ещё будут поддаваться анализу и контролю вручную.

Если провести аналогию для человека-оператора, то в его скриптах нет необходимости расписывать вопросы-отвлечения. Человек сам прекрасно понимает когда ответ отклоняет его от скрипта. Но в отличие от оператора, диалоговой системе нужны однозначные инструкции, которые бы покрывали большинство возможных типов вопросов.

Необходимо выделить отвлечения и перестроить граф таким образом, чтобы отвлечения выделились в отдельный подграф, куда можно было бы перейти из любого места диалога. Такой граф будет покрывать большее количество возможных сценариев.

1.5. Выводы по первой главе

В данной главе была рассмотрена предметная область диалоговых систем, введены определения диалоговой системы, менеджера диалогов, отвлечения. Описаны подходы к построению менеджеров диалогов с соответствующими достоинствами и недостатками. Так же описано представление диалоговой системы в виде графа. Разобрана структура графа, используемого в продуктивном окружении, и структура графа для восстанавливаемой модели. Проведён анализ задачи восстановления графа. Поставлена задача выделения отвлечений.

ГЛАВА 2. ОПИСАНИЕ АЛГОРИТМОВ ПОИСКА ОТВЛЕЧЕНИЙ И РЕКЛАСТЕРИЗАЦИИ

2.1. Процесс разработки

Восстановление графа делается поэтапно, в этот процесс включены следующие этапы в соответствующем порядке:

- Распознавание аудио.
- Преобразование аудио в тексты.
- Препроцессинг в виде преобразования и исправления текстов.
- Кластеризация и восстановление графа.
- Выделение отвлечений.
- Конвертация графа для возможности запуска в существующей системе.

Данная работа не будет затрагивать первые 2 этапа.

В процессе разработки были сделаны улучшения для третьего этапа и затронут четвёртый. Основная работа касается этапа выделения отвлечений, кроме того написаны алгоритмы для преобразования, что даёт возможность сравнить графы с существующими.

Важным будет отметить, что поиск отвлечений встроен в алгоритм построения графа и таким образом появляется возможность внести некоторые изменения в его структуру. Одним из таких изменений может быть выделение отвлечений в отдельный подграф, что позволит рассматривать части независимо.

2.2. Задача выделения отвлечений

Поскольку на этапе выделения отвлечений в графе его перестроение всё ещё возможно, а информация о произнесенных человеком фразах является ценной, так как содержит в себе контекст, то фразы человека оставлены в качестве вершин.

Нужно понимать, что изначально кластеризация проводилась только по фразам оператора. Фразы вершин же были разделены на группы, где для каждой группы совпадала предыдущая и последующая вершины оператора. И уже внутри этих групп вершины разделялись на некоторые подгруппы.

Рассмотрим два подхода в решении задачи выделения отвлечений:

- Первый подход заключается в том, чтобы выделить вершины, у которых много входящих ребер. Порог считается функцией от количества кластеров на которые бьются фразы оператора.

Мы предполагаем, что поскольку отвлечение встречается в разных местах, то и рёбра будут идти в него из множества различных вершин. Такая гипотеза обосновывается следующим наблюдением: в части графа без отвлечений большинство вершин имеет рёбра из одной или двух вершин, поскольку в противном случае скрипт становится слишком сложным, в случае же отвлечения их должно быть много (исключением может быть лишь отвлечение, которое встречается крайне редко).

- В качестве другого подхода можно выделить циклы и сказать, что вершина следующая в диалоге за вершиной повторения с некоторой вероятностью будет являться началом ответа.

Здесь используется наблюдения, о том что после отвлечения на странный вопрос, оператор зачастую повторяет ту же или схожую фразу для возвращения в сценарий. Более того, эта идея используется в графе, который реализован для реального окружения. Там диалоговая система так же повторяет фразу, её сокращенную версию или её иную формулировку, которую произносил, перед тем, как перейти в отвлечение.

Поскольку подходы используют разные идеи их так же имеет смысл комбинировать и использовать данные полученные в обоих подходах.

2.3. Подход поиска отвлечений с большим количеством рёбер входящих в одну вершину

Первым этапом данного подхода является выбор вершин: выбираются вершины, у которых много входящих ребер. Для значений размеров кластеров в интервале от двенадцати до пятнадцати был выбран параметр три. При увеличении количества кластеров соответственно должен увеличиваться и порог.

В некоторых случаях отвлечения не ограничиваются одной вершиной. В связи с этим появляется необходимость определить, где заканчивается то или иное отвлечение, и где, следовательно, оператору необходимо вернуться в вершину с которой он в это отвлечение ушёл. Для этого предположим следующую гипотезу – если в вершину можно попасть только пройдя через отвлечение, то это значит, что фраза является частью соответствующего отвлечения.

Для того, чтобы найти соответствующий хвост для отвлечения будет использован следующий алгоритм: заблокируем вершину и рассмотрим все вершины, достижимые из вершины старта. Важно так же помнить о возможности того, что в кластерах могут быть допущены небольшие ошибки, поэтому если почти весь кластер недостижим после блокировки, то он так же входит в отвлечение.

Ниже представлен псевдокод функции поиска таких хвостов для отвлечений:

```

function DIGRESSION_FINDER(graph, node, dialogs)
  for dialog  $\in$  dialogs do
    flag  $\leftarrow$  True
    for msg  $\in$  dialog do
      msgs[msg.cluster]  $\leftarrow$  msgs[msg.cluster] + 1
      if (msg.cluster == node.cluster)  $\wedge$  flag then
        achieved.add(msg.cluster)
        achieved_msgs[msg.cluster].add(msg)
      else
        flag  $\leftarrow$  False
      end if
    end for
  end for
  for cluster  $\in$  graph.clusters do
    if clusternotinachieved then
      digression_clusters.add(cluster)
    else
      if  $\frac{\textit{achieved\_msgs}[\textit{cluster}]}{\textit{msgs}[\textit{cluster}]} > 0.9$  then
        digression_clusters.add(cluster)
      end if
    end if
  end for
  return digression_clusters
end function

```

Во время тестирования на реальных диалогах человека с человеком была выявлена важная особенность: на работу алгоритма очень сильно влияет качество кластеризации. Поскольку фразы в голосовых диалогах не всегда верно

переводятся в текст, сами фразы сравнительно короткие (а в случае людей-операторов ещё и очень вариативные), поэтому кластеризация оказалась очень некачественной.

2.4. Подход поиска отвлечений с поиском циклов

Предполагается, что оператор всегда действует согласно скрипту, а человек может отвлекаться, поэтому в подавляющем большинстве случаев после ответа на сторонний вопрос, оператор задаёт свой вопрос заново. В связи с данной спецификой, можно выделить несколько полезных особенностей:

- Подавляющее большинство поддиалогов отвлечений достаточно короткие.
- Вопрос повторяет оператор, поэтому циклы можно искать с повторением только операторской вершины.
- Возможно появление отвлечения внутри отвлечения, поэтому необходимо найти оба.

Для решения особенности первого пункта было добавлено ограничение на расстояние между повторяющимися кластерами в диалоге. Оно должно быть не слишком велико, поскольку отвлечения длины больше 2 операторских фраз встречаются в диалогах очень редко, при этом попадают длинные циклы, которые не являются отвлечениями и при их удалении, удаляется большая содержательная часть диалога.

Проблема последнего пункта решается тем, что циклы будут удаляться по мере нахождения. Таким образом внутренний цикл будет удалён раньше внешнего.

В этом случае в качестве потенциальных отвлечений выберем вершины человеческих фраз которые являются первыми после начала цикла. После перевода графа в продуктовый режим, соответствующие вершины станут рёбрами и таким образом в графе появятся рёбра-триггеры, которые будут перенаправлять ход диалога в вершину-отвлечение.

Ниже находится псевдокод для поиска одного цикла:

```
function REMOVE_CYCLES(dialog)
  for msg ∈ dialog do
    msg_num ← msg_num + 1
    if (last_msg[msg.cluster] − msg_num) ≤ 5 then
```

```

    cycle_end ← last_msg[msg.cluster]
    dialog.remove_cluster(cycle_start, msg_num)
  end if
end for
end function

```

2.5. Выборка хороших кластеров

В качестве решения проблемы некачественной кластеризации было предложено использовать данные из фраз человека. На предыдущем этапе фразы человека разделялись на группы по принципу совпадения вершины оператора до и после. Данные о соседних вершинах больше никак не использовались.

Было решено кластеризировать тексты пользователей. Но поскольку, как описывалось выше, кластеризация недостаточно хорошая, то необходимо было отсеять плохие кластеры во избежание каскадных ошибок.

Для этого использовалось попарное сравнение фраз внутри каждого кластера. Для этого сравнивались наборы слов внутри фразы с весами. Если точность превышала порог равный 0.5, то такая пара считалась хорошей.

Константа 0.5 была выведена эмпирически. Для большего значения в кластере начинало содержаться большое количество фраз разных по смыслу.

Проверка всех пар фраз имеет асимптотику по времени $O(n^2)$, где n – количество фраз. Так как на предыдущем этапе восстановления все операции имели асимптотику не более чем $O(n \log n)$, то получилось так, что эта часть занимала значительно больше половины времени от всего восстановления графа.

Было решено для каждого кластера брать случайную выборку, равная утроенному размеру кластера, асимптотически это занимало уже $O(n)$. В силу достаточно больших размеров кластеров (размер их для основной массы данных составляет несколько сотен фраз), статистически, это показывало те же результаты, что и перебор всех пар.

Утверждение 1. При размере диалога от 40 фраз необходимое количество экспериментов для попадания в доверительный интервал $Q = 0.95$ и доверительной вероятности¹ $\varepsilon = 0.1$ достаточно $3n$ экспериментов.

¹Определение доверительного интервала можно изучить здесь.[3]

Доказательство. Значение $p = 0,5$ – наихудшее, в том смысле, что для него вероятность порогового отклонения превысит выбранное значение ε , при наибольшем количестве экспериментов. Таким образом достаточно доказать утверждение для него. Для такого случая результат равен 96, причём вне зависимости от размера возможных исходов, т.е. вне зависимости от размера кластера. Эксперименты рассчитанные для разных значений приведены в учебном пособии.[1]

Таким образом, для более маленьких кластеров можно рассчитать эти значения перебрав все пары полностью, а для больших кластеров, $3n$ экспериментов экспериментов будет достаточно.

Ниже приведён псевдокод для функции, считающей метрику для вершин с большим количеством фраз:

```
function Is_GOOD_CLUSTER(msgs)
  num_of_comparings  $\leftarrow 3 \cdot \text{msgs.size}()$ 
  for (msg1, msg2)  $\in \text{msgs.getPairs}(\text{num\_of\_comparings})$  do
    good_pairs + = msg_compare(msg1, msg2)
  end for
  cluster_threshold  $\leftarrow 0.5$ 
  return  $\frac{\text{good\_pairs}}{\text{num\_of\_comparings}} \geq \text{cluster\_threshold}$ 
end function
```

Функция *msg_compare*(*msg1*, *msg2*) отвечает за сравнение двух сообщений.

2.6. Функция сравнения сообщений

Для алгоритма выше необходима функция сравнения двух сообщений. Для неё должны быть выполнены следующие условия:

- Функция должна быть бинарной. *True* – в случае схожести сообщений, *False* – иначе.
- Функция должна работать за $O(1)$, при условии, что длина сообщений так же принята за $O(1)$.

Второе условие необходимо, поскольку в противном случае при среднем размере групп сообщений от нескольких сотен разница во времени будет в разы увеличиваться.

В качестве алгоритма для сравнения сообщений был выбран следующий: берутся наборы слов в виде множеств и пересекаются с учётом весов слов. Это значение записываем в числитель. В качестве знаменателя берём объединение наборов слов с теми же коэффициентами слов и соответствующими количественными коэффициентами и записываем в знаменатель.

Полученная дробь должна превышать некоторый порог, который ищется вручную.

В качестве альтернативных вариантов функций могут быть использованы следующие:

- Сравнение фраз на основе семантической близости. Об этом алгоритме можно прочитать в следующей статье.[2]
- Сравнение фраз с помощью эмбедингов и на основе косинусного расстояния, которое было использовано в этой работе[5].

Все перечисленные и описанные функции подходят под перечень изначальных условий, поскольку функцию возвращающую нецелые значения в диапазоне от 0 до 1 можно превратить в ступенчатую, сделав её бинарной. Все они работают за асимптотически необходимое время.

2.7. Слияние кластеров операторских сообщений

Изначальный алгоритм, который создавал вершины операторов предполагал точный выбор количества кластеров. Новое решение предполагает избыточное разбиение на кластеры. Поскольку количество кластеров на которые разделяются операторские сообщения можно контролировать вручную, это регулируется достаточно просто.

После такого разбиения выделяются кластеры, фразы в которых максимально похожи между собой, в таком разбиении будет меньше ошибок связанных с кардинальными смысловыми различиями между фразами из одного кластера.

Далее, после выбора хороших человеческих вершин, для каждой операторской вершины рассматриваются все соседние вершины фраз человека. Было опробовано несколько метрик для сравнения соседних вершин.

В первой версии рассматривались все соседские вершины. Они не делились на группы по положению до или после, и, соответственно, это показывало худшие результаты по сравнению со второй версией.

Во второй версии выбирались все вершины до, исходя из предположения, что они будут больше коррелировать с содержанием операторских вершин и не будет проблемы с тем, что не учтён порядок.

2.8. Оценка для сравнения графов

В качестве верной модели использовалась модель написанная вручную. Модель хранится в формате JSON и в ней фразы человека соответственно являются условиями перехода.

Для того, чтобы можно было сравнивать текущую модель и изначальную, необходимо было перевести вторую модель к формату первой.

Здесь стоит напомнить, что в восстанавливаемой модели кластеры человеческих фраз так же являлись вершинами. Таким образом для того, чтобы была возможность перевести их в рёбра и соответственно сопоставить изначальному графу, для каждой вершины с человеческими кластерами должны были быть верны следующие условия:

- У вершины должна быть ровно одна предыдущая, причём эта вершина должна быть вершиной операторских фраз.
- Последующая вершина должна быть так же операторской вершиной, причём она должна являться единственной последующей вершиной для данной.
- Ни одна из набора фраз данного кластера не должна совпадать с фразами из кластеров, в которые ведут рёбра из предыдущей операторской вершины.

Реализация была сделана при помощи классов графа, новых вершин, рёбер. После проверки на то, что нет коллизий описанных выше можно провести инициализацию с помощью псевдокода приведенного ниже:

```
function CREATE_NODES(dialogs, digression)
  for dialog ∈ dialogs do
    for msg ∈ dialog do
      if msg.is_incoming then
        if msg.next then
          transactions[msg.prev.cluster].add(vertex[msg.next.cluster])
        else
          end_vertexes.add(msg.prev.cluster)
```

```

        end if
    else
        vertex_msgs[msg.cluster].add(msg)
        vertex_clusters.add(msg.cluster)
    end if
end for
end for
for cluster ∈ vertex_clusters do
    vertexes.add(Vertex(transactions[cluster], vertex_msgs[cluster],
        cluster ∈ end_vertexes, cluster ∈ digressions))
end for
return vertexes
end function

```

Проверив, что все необходимые условия соблюдены, можно преобразовывать граф, выделяя отвлечения следующим образом:

- Каждая вершина новой модели сопоставлялась одной изначальной. Сопоставление проводилось путём сравнения произносимой фразы.
- Вершины фраз человека были выделены в группы, каждая группа относилась к одному ребру и являлась будущим набором исходящих рёбер.
- Технические вершины в изначальном графе игнорировались, но учитывалось их место в последовательности модели.
- Сравнивались наборы вершин отвлечения в изначальном графе и в полученном.
- При создании метрики для графа была использована статья [11]

ГЛАВА 3. РЕЗУЛЬТАТЫ ЭКСПЕРИМЕНТОВ

3.1. Используемые данные

В качестве данных использовалось два принципиально разных типа наборов диалогов. Первые проводились уже с существующим скриптом. Второго типа, это данные диалогов человека с человеком.

1.Особенность первого заключается в том, что там легко кластеризовать фразы оператора, так как они произносятся всегда одинаково, за исключением вариаций для коротких вариантов.

Но необходимо понимать, что случаи повторения фразы имеют другие вершины соседей поскольку рассчитаны на специфичные случаи. Такие объединения были полезны для алгоритма поиска циклов, но находить их было сложно в силу того, что они явным образом не попадали под условия. Так же можно было выявить слишком заниженные коэффициенты, при слиянии фраз, которые не являлись одинаковыми.

Важным вариантом использования таких данных являлась возможность перевести полученный граф в тот же формат и сравнить полученный результат с оригиналом.

2.Особенности второго соответственно следующие: во первых там говорят разные операторы и у них разный стиль подачи одних и тех же данных. Во вторых диалоги более сложные и зачастую отходят от скрипта. В третьих люди решают более сложные вопросы и умеют давать ответы на не предусмотренные скриптом вопросы. На этот вариант данных стоит ориентироваться, но в связи с отсутствием оригинального скрипта в удобном формате напрямую сравнить полученный результат представляется возможным только вручную.

В обоих случаях есть сложности с переводом речи людей в текст, поэтому иногда даже рассматривая текст диалога вручную нельзя понять что человек имел ввиду.

Для всех датасетов первого типа, существуют файлы, содержащие в себе оригинальный граф, который принимается за верный. Все данные сериализованы в формате JSON.

На Рисунке 3 представлен восстановленный граф, получившийся после одного из запусков. С этой моделью происходила непосредственная работа:

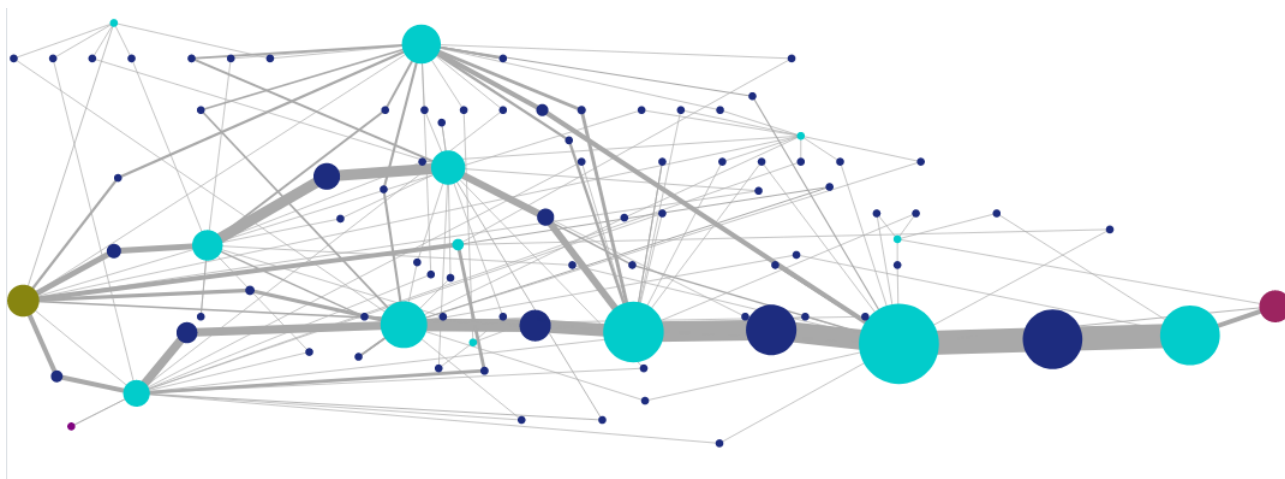


Рисунок 3 – Восстановленный граф

В этой модели голубым цветом обозначены операторские вершины, синим соответственно человеческие. Кроме того, для удобства анализа выделены две фиктивные вершины: стартовая и вершина окочания диалогов. В дальнейшем для приведённых данных они в расчёт не брались.

На Рисунке 4 представлена модель адаптированная для текста работы. На ней будут показаны результаты работы алгоритмов:

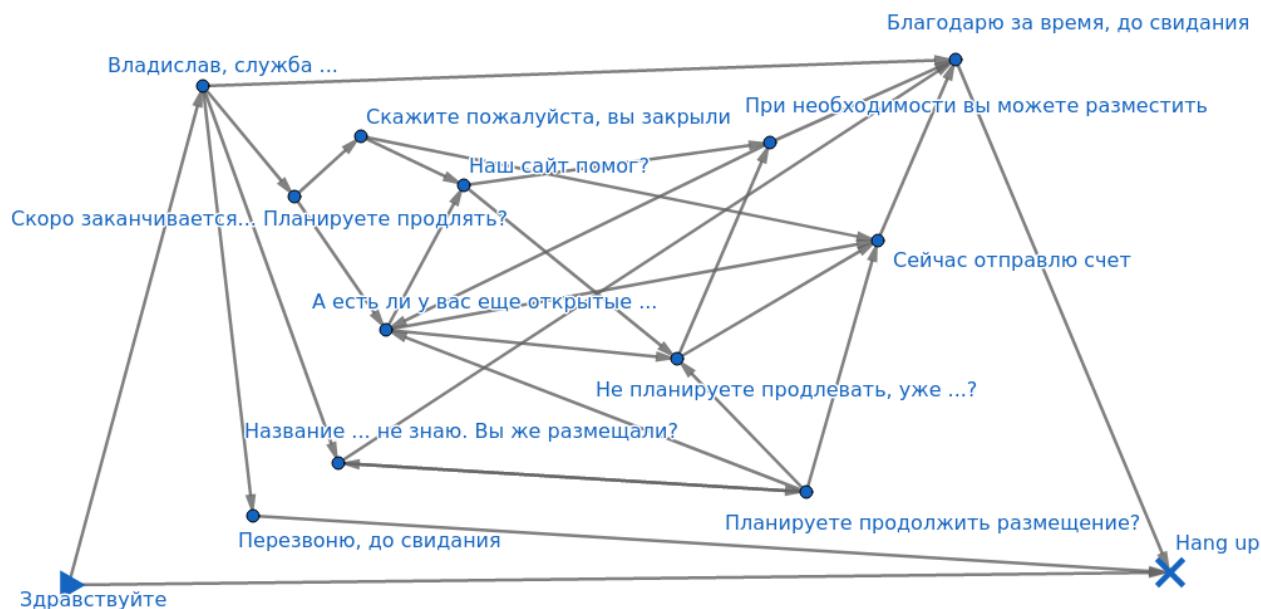


Рисунок 4 – Адаптированный граф

NB: В данном примере вершины человеческих кластеров и вершины операторских кластеров с фразами «Алло», «Я говорю» и «Повторите пожалуйста» были намеренно убраны, а так же склеены вершины разных приветствий, поскольку их наличие делало восприятие графа более затруднительным.

Про эти вершины и особенности будут написаны отдельные заметки везде, где граф будет использоваться в качестве примера.

По тем же причинам были удалены большинство рёбер ведущих в конец диалога — человек может бросить трубку в любой момент, на алгоритмы эти рёбра не влияют.

Изначально при создании графа вершины человеческих фраз просто группируются по тому, какие 2 кластера находятся до и после, поэтому становится возможным заменить их на рёбра, не потеряв ценной информации.

3.2. Результаты улучшения кластеризации

Алгоритм вносит небольшие изменения и способен объединять некоторые одинаковые кластеры.

Основная сложность состояла в анализе, так как для данных из разговоров с искусственным интеллектом кластеризация фраз оператора изначально была очень хорошей в силу того, что фразы повторялись. Для случаев разговоров человека с человеком анализ результатов приходилось проводить вручную в силу отсутствия разметки данных.

Были проведены исследования по влиянию параметров алгоритма на результат. В результате чего были выявлены следующие закономерности:

- Пороговое значение схожести сообщений для каждого кластера должно варьироваться от 0.2 до 0.5. В этом промежутке в качестве пар схожих сообщений выбираются схожие пары сообщений. Для меньшего значения соответственно практически полностью одинаковые. Для большего у них появляется вариативность.
- Пороговое значение количества хороших пар в кластере должно варьироваться так от 0.25 до 0.5. Для больших значений в одном кластере начинают появляться наборы фраз несущие разный смысл.

Ниже, на Рисунке 5 представлены примеры различных значений промежутков порога для количества хороших пар в кластере:

нет	до свидания	але
нет нет	спасибо до свидания	на какую тему подскажите
нет	спасибо всего доброго до свидания	даша
нет александр александрович	всё спасибо вам за напоминание всего доброго	сколько минут
		интернет
		я

(a) 0.3 (b) 0.5 (c) 0.6

Рисунок 5 – Различные значения пороговой функции количества хороших пар

Для графа с Рисунка 4, в качестве примера работы данного алгоритма, рассмотрим пример объединения 2 вершин приветствия¹. Набор входящих и исходящих рёбер для вершин приветствия можно увидеть на Рисунке 6:

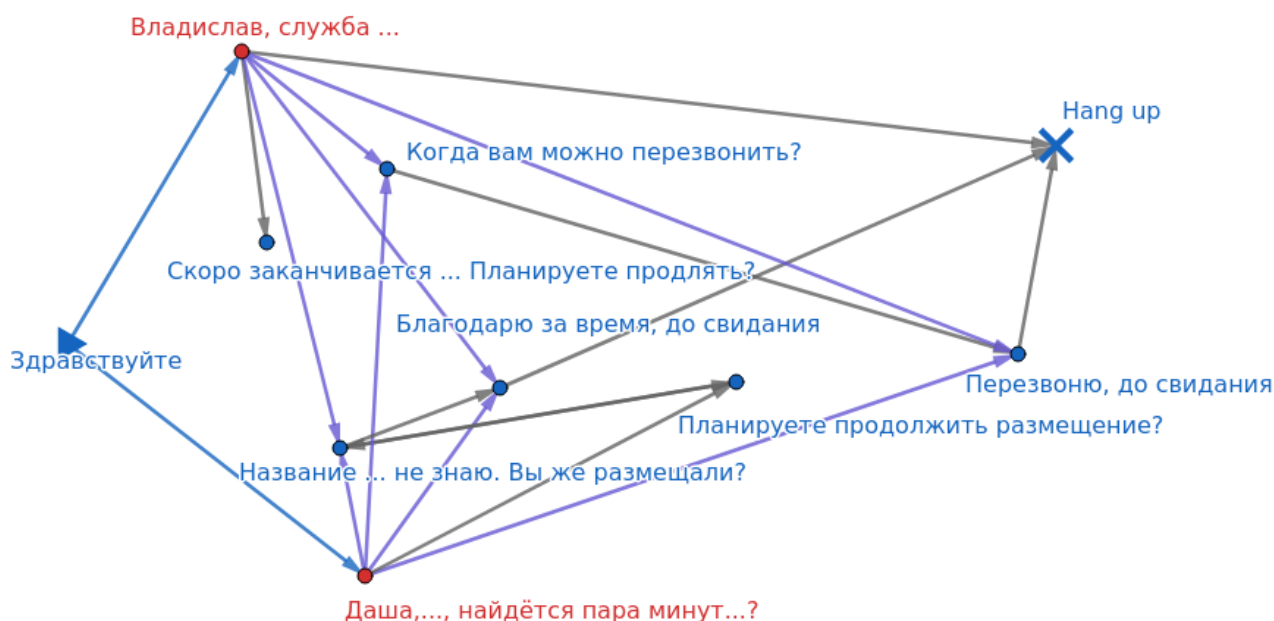


Рисунок 6 – Восстановленный граф

Среди данных вершин были найдены две вершины с приветствием из разных запусков. Важно отметить, что при уменьшении количества кластеров, на которые будут поделены вершины будут сливаться такие пары как «Скоро заканчивается ..., планируете продлевать?» и «Не планируете продлевать, уже ...», в связи с тем, что семантика этих фраз похожа, а приветствия сформулированы по-разному. Таким образом уменьшение количества кластеров проблему не решает.

¹на Рисунке 4 для удобства, они были уже объединены

Пары соответствующих рёбер, ведущих к одинаковым вершинам (рёбра помечены фиолетовым и синим) имеют схожие фразы, которые объединяются при рекластеризации (кластеры : («слушаю, да»), («нет вакансий»), («здравствуйте») и т.д.). Таким образом у вершин появляются множества кластеров до и после с большим пересечением и они объединяются.

Подход с пересечением групп операторских кластеров до и после, не может быть использован, в силу того, что в этом случае фразы «Алло» и «Повторите пожалуйста» будут объединяться в силу того, что в них идут рёбра и из них идут рёбра практически во все вершины.

3.3. Результаты работы алгоритма поиска отвлечений по рёбрам

Для данного алгоритма критично качество кластеризации и наличие небольшого количества выбросов.

В случае с данными из реальных диалогов зачастую не было видно правильно выделенных отвлечений из-за смещения кластеров и операторских и человеческих вершин.

Для данных из диалогов с существующей диалоговой системой алгоритм находил все отвлечения и так же отмечал некоторое количество дополнительных вершины, которые отвлечениями не являются. Основная причина столь большого различия в том, что операторские фразы во втором случае кластеризуются практически идеально.

На данном этапе этот алгоритм корректно использовать в качестве отсеивания вершин, которые не подходят в качестве отвлечений, но были ошибочно выбраны другими алгоритмами. Также было сделано следующее наблюдение: алгоритм хорошо подсвечивает последовательности фраз, которые часто произносятся друг за другом. Например, после фразы «Когда вам перезвонить» – фраза «Перезвоню, до свидания» встечается **74.8%** случаев².

3.4. Результаты алгоритма поиск циклов

После того, как циклы найдены, их необходимо удалить. Таким образом основной граф уменьшается в размерах и его становится легче анализировать. Из особенностей реализации важно отметить, что удаляются не все циклы, а только размера не более чем 5. Рассмотрим пример следующего диалога:

Оператор: Здравствуйте

²Среди диалогов, которые не обрываются после первой фразы.

Человек: Ало

*Оператор: Меня зовут Даша, вы ***³?*

*Человек: *шуриание*⁴*

Оператор: Извините, я вас не расслышал, повторите пожалуйста

Человек: Я говорю давайте задавайте

*Оператор: Вы сейчас ***?*

Человек: сейчас нет

*Оператор: Вы нашли ***?*

Человек: да

*Оператор: Убрать *** в архив?*

*Человек: *неразборчивый текст**

Оператор: Извините, я вас не расслышал, повторите пожалуйста

*Человек: ненадо пока убирать ****

Оператор: Извините за беспокойство, всего доброго, до свидания

Если удалить все фразы между двумя просьбами, то состояния которые не являются отвлечениями будут утеряны. В данном случае они составляют основную часть диалога.

Ниже, на Рисунке 7 приведен график, в котором сравниваются количества вершин, до и после применения алгоритма удаления циклов. Для каждого графа так же указано, датасет какого размера был использован при создании:

³Здесь и далее *** будут использоваться для обезличивания частей диалога, которые не содержат важных для алгоритма данных, но несут в себе персональную информацию или информацию о компаниях.

⁴Здесь и далее *описание* для пропуска не содержательных частей с коротким описанием.

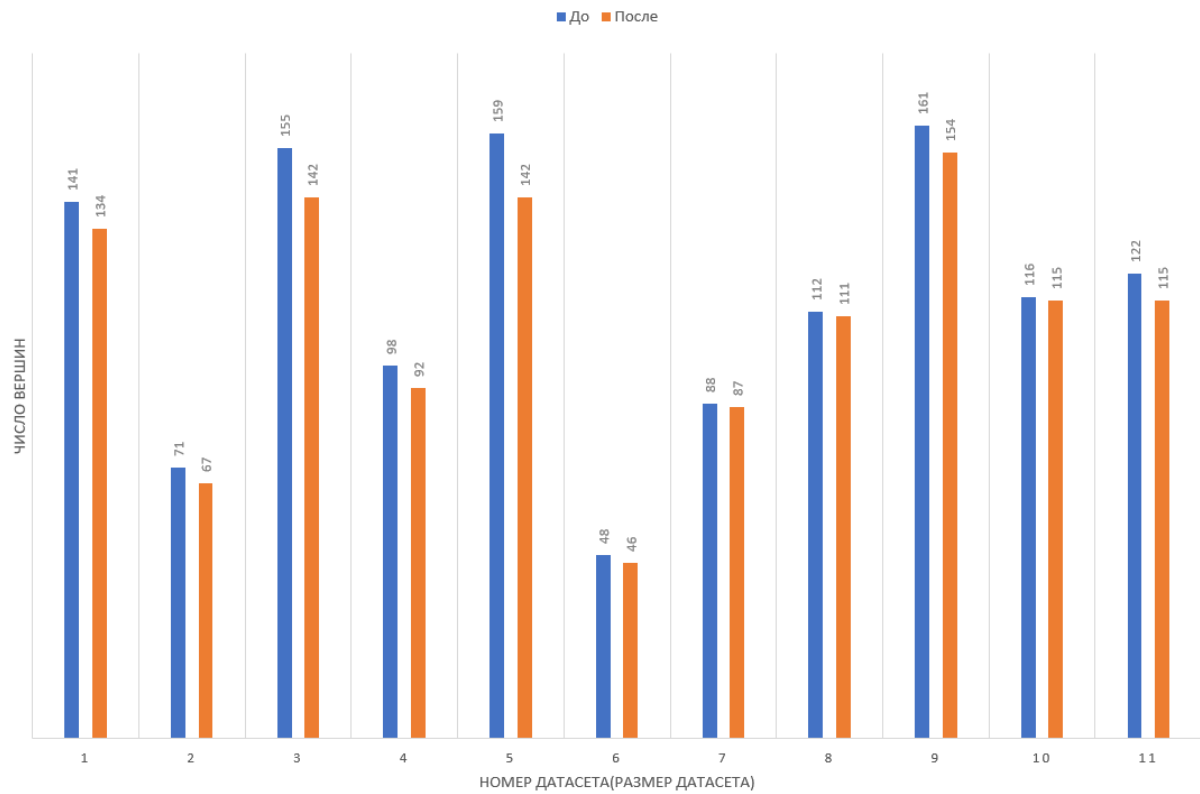


Рисунок 7 – Изменение количества вершин

Ниже, на Рисунке 8 приведен график, в котором сравниваются количества рёбер, до и после применения алгоритма удаления циклов.

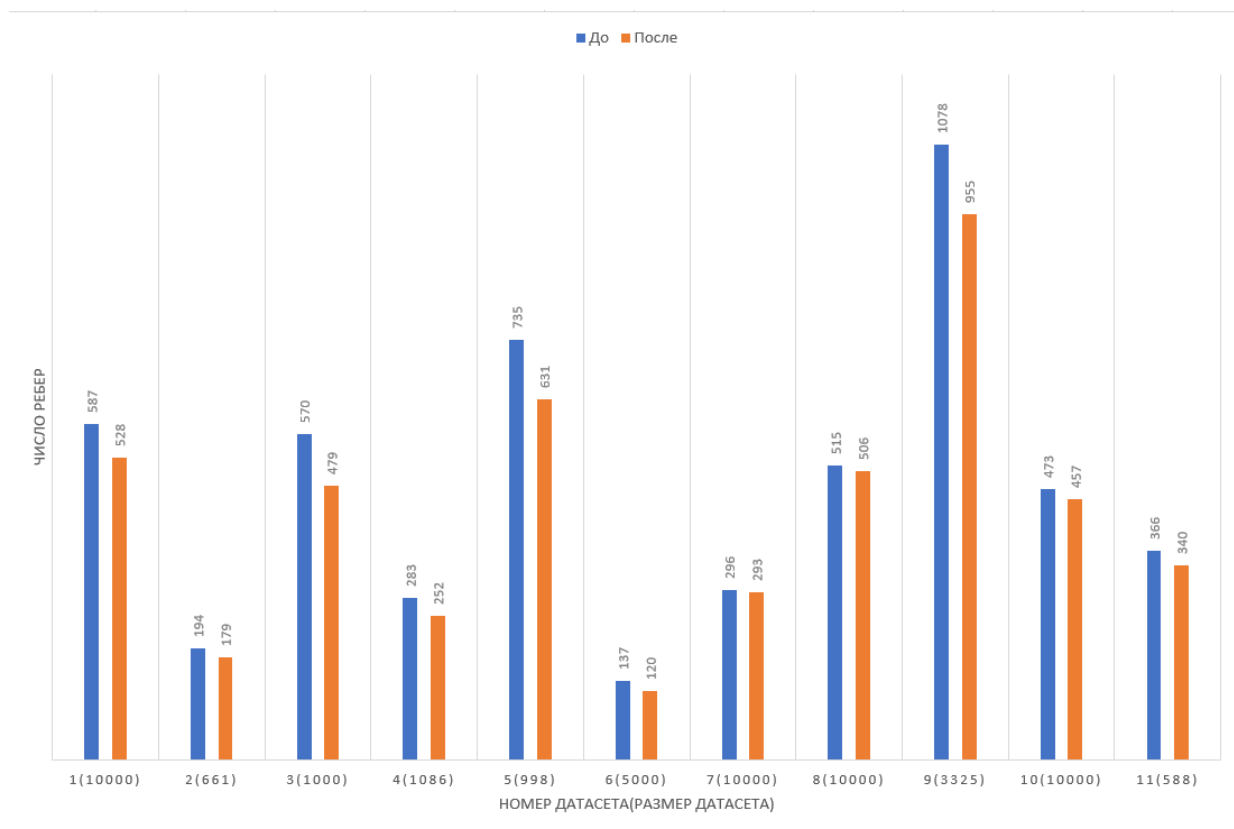


Рисунок 8 – Изменение количества вершин

Среди этих датасетов есть 2 принципиально разных типа:

Первый тип – **оператор + человек**, эти датасеты были предоставили компании заказчиками, и в них человек придерживается некоторого скрипта, но заметно от него отклоняется, поэтому при меньших размерах порядок вершин в них такой же. К этому типу относятся соответственно датасеты 2, 4, 5, 9, 11.

Второй тип – **искусственный интеллект + человек**, эти датасеты были собраны соответственно компанией при звонках со скриптами написанными вручную. К этому типу соответственно относятся 1, 3, 6, 7, 8, 10.

В обоих типах при увеличении размера датасета увеличивается количество вершин. Важно отметить, что большинство из вершин, это ответы человека, количество же операторских вершин варьируется от 12 до 25. Для случая диалогов человека с человеком размеры графа могут быть изначально больше, это объясняется тем, что в нём большее разнообразие фраз.

Вырезанные ребра и вершины соответственно отделяются и становится проще анализировать граф.

Если рассмотреть какие вершины из операторских попадают в циклы, то хорошо видно тенденцию того, что есть несколько вершин у которых большое

количество ответов попадает в циклы. Некоторые такие вершины, как можно заметить исчезают вообще. Для остальных можно подобрать пороговое значение при котором можно будет считать, что вершина является триггером отвлечения.

Для примера графа указанного выше отметятся следующие вершины, содержащие большое количество фраз в циклах (см. Рисунок 9):

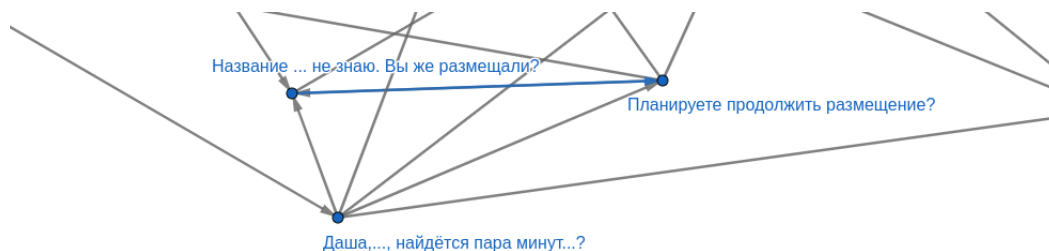


Рисунок 9 – Отмеченный цикл

Интересным будет заметить, что из двух вершин именно «Название ... не знаю, вы же размещали.» попадёт в отвлечения, т.к. внутри цикла находятся именно её сообщения, а сообщения «Планируете продолжить размещение?» находятся на границах цикла.

Более того, самый большой процент (84.7%) попавших в циклы вершин будет иметь фразу «Я говорю», которая была исключена, т.к. сильно усложняла граф. Она является примером **общего** для всех графов отвлечения (ответом на вопрос «повторите» и т.п.). Таким образом, поведение алгоритма выделившего её является корректным.

3.5. Программная реализация

Предложенные алгоритмы были реализованы на языке Python с использованием интегрированной среды разработки Pycharm. Кроме того использовался текстовый редактор Visual Studio Code с внутренним плагином для визуализации графов описанных в формате JSON. Разработка велась в компании Dasha.AI.

Компания использует приватный Git-репозиторий, находящийся на хостинге Gitlab. Для удобства локальной разработки приложение запускалось в Docker-контейнере.

ЗАКЛЮЧЕНИЕ

В ходе выполнения работы была рассмотрена задача выделения отвлечений в графовой модели для голосовой диалоговой системы. Для её решения были разработаны и написаны алгоритмы основанные на особенностях разговоров оператора с человеком. Были учтены особенности существующей модели графовых диалогов и данные полученные при её использовании.

Были написаны алгоритмы выделения циклов и поиска отвлечений по рёбрам. Эти алгоритмы позволяют выделить значительное количество отвлечений и уменьшают размер основного графа. Кроме того в ходе разработки стало ясно, что алгоритмы очень чувствительны к кластеризации и возникла необходимость улучшить её качество. Для чего был разработан алгоритм на основе данных о соседних фразах.

Поскольку работа является частью большего проекта. То интеграция в реальное окружение запланирована по завершению всех частей проекта. Разработка должна будет автоматизировать часть рабочих процессов.

Одно из возможных направлений развития работы, использовать информацию о пользователе в диалоге, сделав его таким образом более естественным. Так же существует возможность выделить некоторые общие отвлечения для всех диалогов, получая таким образом возможность предсказывать отвлечения в виде подсказок.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 *Мухин О.* Моделирование систем: учебник [Электронный ресурс]. — 2014. — URL: <http://stratum.ac.ru/education/textbooks/modelir/lection34.html> (дата обр. 30.04.2020).
- 2 *Нзюк Н. Б., Тузовский А. Ф.* Классификация текстов на основе оценки семантической близости терминов // *Izvestiya Tomskogo Politehnicheskogo Universiteta Inziniring Georesursov*. — 2012. — Т. 320, № 5. — С. 43–48.
- 3 *Чернова Н.* Теория вероятностей: Учеб. пособие/Новосиб. гос. ун-т // Новосибирск. 2007, 160 с. — 2007.
- 4 Dialogue-based Management of user Feedback in an Autonomous Preference Learning System. / J. Lucas-Cuesta [et al.] // Vol. 1. — 09/2010. — P. 330–336.
- 5 *Dönmez İ., Pashaei E., Pashaei E.* Word Vector Space for Text Classification and Prediction According to Author. —
- 6 *Doshi F., Roy N.* Efficient model learning for dialog management // 2007 2nd ACM/IEEE International Conference on Human-Robot Interaction (HRI). — 2007. — P. 65–72.
- 7 *Jokinen K., McTear M.* Spoken Dialogue Systems. — Morgan & Claypool Publishers, 2010. — (Synthesis lectures on human language technologies). — ISBN 9781598295993. — URL: <https://books.google.ru/books?id=uawwulnD020C>.
- 8 *Li H.* Text Clustering // *Encyclopedia of Database Systems* / ed. by L. LIU, M. T. ÖZSU. — Boston, MA : Springer US, 2009. — P. 3044–3046. — ISBN 978-0-387-39940-9. — DOI: 10.1007/978-0-387-39940-9_415. — URL: https://doi.org/10.1007/978-0-387-39940-9_415.
- 9 *Wikipedia contributors.* Dialog manager — Wikipedia, The Free Encyclopedia. — 2020. — URL: https://en.wikipedia.org/w/index.php?title=Dialog_manager&oldid=941891414 (visited on 04/18/2020).
- 10 *Williams J. D., Young S.* Partially observable Markov decision processes for spoken dialog systems // *Computer Speech & Language*. — 2007. — Vol. 21, no. 2. — P. 393–422.

- 11 *Wills P., Meyer F. G.* Metrics for graph comparison: A practitioner's guide
// Plos one. — 2020. — Vol. 15, no. 2. — e0228728.

ПРИЛОЖЕНИЕ А. ИСХОДНЫЙ КОД НА ЯЗЫКЕ PYTHON

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  import json
3  import logging
4  import os
5
6  import numpy as np
7  import requests
8  from flask import Flask, Response, jsonify, request
9  from flask_caching import Cache
10 from sklearn.cluster import AgglomerativeClustering, KMeans
11
12 from reconstruction import algorithm_descriptions, reconstruct_graph
13 from reconstruction.data_structure import DatasetStorage
14 from reconstruction.data_structure.graph import save_graph
15 from reconstruction.data_structure.message import messages_to_json
16 from reconstruction.embedding import EmbeddingStorage
17
18 @app.route('/graph', methods=['POST'], strict_slashes=False)
19 @cache.cached(key_prefix=make_cache_key)
20 def build_graph():
21     dataset_name, dialogs, alg_name, hparams = parse_request()
22     graph = reconstruct_graph(dialogs, alg_name, hparams, embedding_storage)
23     save_graph(graph, dataset_name)
24     return jsonify(graph)

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  from reconstruction.reconstruct_graph import reconstruct_graph,
   ↪ algorithm_descriptions

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  from sklearn.cluster import KMeans
3
4  import reconstruction.algorithm as alg
5  from reconstruction.embedding.mean_vectorizer import MeanEmbeddingVectorizer
6
7  algorithms = {
8      "cluster_all_messages": {
9          },
10     "cluser_incoming_outgoing": {
11         },
12     "cluster_incoming_outgoing_on_each_step": {
13         },
14
15     "cluster_outgoing_and_incoming_after": {
16         },

```

```

17     "cluster_prev_and_next": {
18     },
19     "cluster_prev_and_next_recluster_incoming": {
20         "reconstruction": alg.cluster_prev_and_next_recluster_incoming,
21         "clustering": {
22             "emb_name": "tenth.norm-sz500-w7-cb0-it5-min5.w2v",
23             "vectorizer": MeanEmbeddingVectorizer,
24             "alg": KMeans
25         }
26     }
27 }

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  from collections import Counter, defaultdict
3
4  import numpy as np
5  from sklearn.metrics.pairwise import cosine_similarity
6
7  from reconstruction.data_structure.graph import Node
8  from reconstruction.converter.converter import Graph
9  from reconstruction.digressions.digressions import cycle_finder
10
11
12  def calc_centroid(messages):
13
14
15  def create_nodes(messages, cluster_labels, remove_msgs=[]):
16      #secsitive code - create nodes using information about msgs from cycle,
17      ↪ remove msgs keep in dialogs for analyse
18
19  def add_start_node(graph, dialogs, skip):
20      graph["nodes"].append({
21          "id": -1,
22          "group": -1,
23          "centroid": 0,
24          "texts": [],
25          "seq_num": -1,
26          "seq_num_std": 0,
27          "incoming": False,
28          "size_coef": 0.4,
29          "size": len(dialogs),
30          "x": 0,
31          "y": 0
32      })
33      start_cluster = -1
34      graph["start"] = start_cluster
35      weights = defaultdict(int)
36      for dialog in dialogs:

```

```

37         if (len(dialog) > 0):
38             msg = dialog.log[0]
39             cluster = msg.cluster
40             if cluster and cluster not in skip:
41                 weights[cluster] += 1
42     for cluster in weights:
43         graph["links"].append({
44             "source": int(start_cluster),
45             "target": int(cluster),
46             "weight": int(weights[cluster])
47         })
48
49
50 def add_end_node(graph, dialogs, skip):
51
52
53 def nodes_to_graph(nodes, dialogs, node_size_limit=0):
54     #secsitive code - create nodes, count size, linkes with weights, work with
55     ↪ digressions
56
57 def num_of_vertexes(msgs, remove_msgs=[]):
58     vertexs = defaultdict(set)
59     none_counter = 0
60     for msg in msgs:
61         if msg in remove_msgs:
62             continue
63         if msg.next_msg is not None and msg.next_msg not in remove_msgs:
64             vertexs[msg.cluster].add(msg.next_msg.cluster)
65     edges_num = 0
66     for _, set_list in vertexs.items():
67         edges_num += len(set_list)
68     print("vertexs:", len(vertexs), " edges: ",
69         edges_num, " nones: ", none_counter)
70
71
72 def create_graph(messages, cluster_labels, dialogs, node_size_limit=0,
73     ↪ call_remove_cycles=False):
74     #sensitive code - create nodes, graph, call cycles finder, count statistics

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  from reconstruction.config import algorithms
3  from reconstruction.clustering import MessageClusterizer
4  from reconstruction.create_graph import create_graph
5  from reconstruction.coords import graphviz_coords
6  import logging
7
8
9  def get_vectorizer(embedding_storage, clustering_config):

```

```

10
11
12 def get_clustertering_alg(clustering_config, **kwargs):
13
14
15 def preprocess_dialogs(vectorizer, dialogs):
16
17
18 def reconstruct_graph(dialogs, alg_name, hparams, embedding_storage):
19     alg_config = algorithms[alg_name]
20     reconstruction_alg = alg_config["reconstruction"]
21
22     clustertering_alg = get_clustertering_alg(alg_config["clustering"],
23                                               **{"n_clusters":
24                                                  ↪ int(hparams["n_clusters"])}))
25
26     vectorizer = get_vectorizer(embedding_storage, alg_config["clustering"])
27     preprocess_dialogs(vectorizer, dialogs)
28
29     clustered_messages, cluster_labels = reconstruction_alg(
30         dialogs, clustertering_alg, **hparams)
31
32     min_node_size = int(hparams["min_node_size"])
33
34     remove_cycles = hparams["remove_cycles"]
35     graph = create_graph(clustered_messages, cluster_labels, dialogs,
36                           node_size_limit=min_node_size,
37                           ↪ call_remove_cycles=remove_cycles)
38
39     digressions_by_different_income(graph, dialogs)
40
41     graphviz_coords(graph)
42     return graph

```

```

1 # Copyright (C) 2020 Dasha.AI Inc.
2 from reconstruction.algorithm.cluser_incoming_outgoing import \
3     cluser_incoming_outgoing
4 from reconstruction.algorithm.cluster_all_messages import cluster_all_messages
5 from reconstruction.algorithm.cluster_incoming_outgoing_on_each_step import \
6     cluster_incoming_outgoing_on_each_step
7 from reconstruction.algorithm.cluster_outgoing_and_incoming_after import \
8     cluster_outgoing_and_incoming_after
9 from reconstruction.algorithm.cluster_prev_and_next import \
10    cluster_prev_and_next
11 from reconstruction.algorithm.cluster_prev_and_next_recluster_incoming import \
12    cluster_prev_and_next_recluster_incoming

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  import itertools
3  import logging
4  import multiprocessing
5  import random
6  import sys
7  import time
8  from collections import Counter, defaultdict
9  from functools import partial
10 from typing import Dict, List
11
12 import numpy as np
13 from sklearn.cluster import KMeans
14
15 from reconstruction.data_structure import Dialog
16 from reconstruction.data_structure.dialog import get_dialog_messages
17 from reconstruction.clustering import MessageClusterizer
18 from reconstruction.algorithm import cluster_prev_and_next
19 from reconstruction.algorithm.utils import split_incoming_outgoing
20
21
22 def vocab_distribution(words):
23     total = len(words)
24     vocab_dist = dict(Counter(words))
25     for word in vocab_dist.keys():
26         vocab_dist[word] /= total
27     return vocab_dist
28
29
30 def distribution_text_similarity(msg1, msg2, vocab_dist):
31     words1 = msg1.text.split()
32     words2 = msg2.text.split()
33     counter1 = dict(Counter(words1))
34     counter2 = dict(Counter(words2))
35     unique_words1 = set(words1)
36     unique_words2 = set(words2)
37
38     cross = sum(
39         min(counter1[word], counter2[word]) * vocab_dist[word]
40         for word in unique_words1 & unique_words2
41     )
42
43     union = sum(
44         max(counter1.get(word, 0), counter2.get(
45             word, 0)) * vocab_dist[word]
46         for word in unique_words1 | unique_words2
47     )
48
49     similarity = cross / union

```



```

50     return similarity
51
52
53 def remove_invalid_clusters(msgs, clusters, similarity_threshold=0.4,
    ↪ cluster_threshold=0.4):
54     start = time.time()
55     words = [w for t in msgs for w in t.text.split()]
56     vocab_dist = vocab_distribution(words)
57     print("vocab_distribution:", time.time() - start, file=sys.stderr)
58
59     start = time.time()
60     cluster_messages = defaultdict(list)
61     for cl_n, msg in zip(clusters, msgs):
62         cluster_messages[cl_n].append(msg)
63
64     print("cluster_messages:", time.time() - start, file=sys.stderr)
65
66     start = time.time()
67     msgs_in_cluster = dict(Counter(clusters))
68     good_connects = defaultdict(int)
69     msgs_in_stat = defaultdict(int)
70     pool = multiprocessing.Pool(5)
71
72     for cl_n, msg_l in cluster_messages.items():
73         if cl_n in msgs_in_cluster:
74             for i, msg in enumerate(msg_l):
75                 for j in range(i):
76                     # check ~3 * cl_num random pairs
77                     if (random.randint(0, msgs_in_cluster[cl_n] // 3) == 0):
78                         msgs_in_stat[cl_n] += 1
79                         sim = distribution_text_similarity(
80                             msg, msg_l[j], vocab_dist)
81                         if sim > similarity_threshold:
82                             good_connects[cl_n] += 1
83     bad_clusters = {
84         cl_n for cl_n, correct in good_connects.items()
85         if (correct / msgs_in_stat[cl_n] < cluster_threshold)
86     }
87
88     print(msgs_in_cluster)
89     print("msg list size: ", len(cluster_messages),
90         " keys: ", cluster_messages.keys)
91     print([
92         cl_n for cl_n, correct in good_connects.items()
93         if (correct / msgs_in_stat[cl_n] > cluster_threshold)
94     ], " good one")
95     print([
96         cl_n for cl_n, correct in good_connects.items()
97         if (correct / msgs_in_stat[cl_n] < cluster_threshold)
98     ], "bad one")

```

```

99
100     start = time.time()
101     new_msgs, new_clusters = [], []
102     for msg, cl_n in zip(msgs, clusters):
103         if cl_n not in bad_clusters:
104             new_msgs.append(msg)
105             new_clusters.append(cl_n)
106     print("filter:", time.time() - start, file=sys.stderr)
107
108     return new_msgs, new_clusters
109
110
111 def recluster_incoming(clustered_messages, cluster_labels, dialogs, cluster,
112 ↪ **hparams):
113     """Делит входящие сообщения иначе, с учётом того, как они кластеризовались
114     ↪ независимо от сообщений оператора
115     1. Кластеризуем входящие сообщения независимо
116     2. Разделяем ноды из входящих, если там много сообщений каких-то типов
117     3. Объединяем ноды с полностью одинаковыми кластерами
118     """
119     print(hparams, file=sys.stderr)
120
121     proportion_threshold = float(hparams.get("proportion_threshold", 0.25))
122     cluster_threshold = float(hparams.get("cluster_threshold", 0.5))
123     similarity_threshold = float(hparams.get("similarity_threshold", 0.5))
124
125     start = time.time()
126
127     messages = get_dialog_messages(dialogs)
128     incoming, _ = split_incoming_outgoing(messages)
129     print(" get_dialog_messages split_incoming_outgoing:",
130           time.time() - start, file=sys.stderr)
131     start = time.time()
132     incoming_cluster_labels = cluster.fit_predict(incoming)
133     print("fit_predict:", time.time() - start, file=sys.stderr)
134
135     start = time.time()
136     incoming, incoming_cluster_labels = remove_unvalid_clusters(
137         incoming, incoming_cluster_labels,
138         similarity_threshold=similarity_threshold,
139         cluster_threshold=cluster_threshold)
140
141     print("remove_unvalid_clusters:", time.time() - start, file=sys.stderr)
142
143     start = time.time()
144     pair_clusters = defaultdict(int) # new clusters, incoming only
145     for msg, cluster in zip(incoming, incoming_cluster_labels):
146         if msg in clustered_messages: # cluster_messages doesn't include
147             ↪ starts of dialogs
148             pair_clusters[msg] = cluster

```

```

146
147 node_clusters = defaultdict(lambda: defaultdict(list))
148 new_clustering = defaultdict(int)
149
150 for msg, cluster in zip(clustered_messages, cluster_labels):
151     if msg in pair_clusters:
152         node_clusters[cluster][pair_clusters[msg]].append(msg)
153
154 new_cluster_num = cluster_labels.max()
155 alloc_clusters = defaultdict()
156 changing_msgs = 0
157
158 for _, clusters in node_clusters.items():
159     node_size = sum(map(len, clusters.values()))
160     small_part = []
161     for income_cluster_n, msgs in clusters.items():
162         cluster_proportion = (len(msgs) / node_size)
163         if (cluster_proportion > proportion_threshold):
164             if income_cluster_n not in alloc_clusters:
165                 new_cluster_num += 1
166                 alloc_clusters[income_cluster_n] = new_cluster_num
167
168                 cluster_num = alloc_clusters[income_cluster_n]
169                 for msg in msgs:
170                     new_clustering[msg] = cluster_num # check
171                     changing_msgs += 1
172             else:
173                 small_part.append(income_cluster_n)
174
175     if len(small_part) == 1:
176         alone_cluster = small_part[0]
177         if alone_cluster not in alloc_clusters:
178             new_cluster_num += 1
179             alloc_clusters[alone_cluster] = new_cluster_num
180             cluster_num = alloc_clusters[alone_cluster]
181             for msg in clusters[alone_cluster]:
182                 new_clustering[msg] = cluster_num # check
183                 changing_msgs += 1
184 print("split clusters:", time.time() - start, file=sys.stderr)
185
186 start = time.time()
187 list_of_changes = []
188 for msg in new_clustering:
189     i = clustered_messages.index(msg)
190     clustered_messages[i].cluster = new_clustering[msg]
191     cluster_labels[i] = new_clustering[msg]
192     list_of_changes.append(i)
193 print("rename clusters:", time.time() - start, file=sys.stderr)
194 return clustered_messages, cluster_labels
195

```

```

196
197 def cluster_prev_and_next_recluster_incoming(dialogs: List[Dialog], cluster:
↳ MessageClusterizer, **hparams):
198     """Кластер входящих определяется как единый кластер, если следующий
↳ исходящий кластер одинаковый для всех сообщений
199
200     1. Строим кластеры для всех исходящий сообщений
201     2. Для каждого кластера находим все последующие входящие сообщения до
↳ следующего исходящего
202     3. Эти входящие сообщения разделяются на кластеры в зависимости от того
↳ какой исходящий кластер следующий
203     4. Номер кластера задается на основе следующего исходящего кластера.
204
205     Таким образом переход от кластера входящих сообщений к кластеру исходящих
↳ становится однозначным.
206
207     Входящие сообщения разделяются, с учётом того, как они кластеризовались
↳ независимо от сообщений оператора
208         1. Кластеризуем входящие сообщения независимо
209         2. Разделяем ноды из входящих, если там много сообщений каких-то типов
210         3. Объединяем ноды с полностью одинаковыми кластерами
211
212     Arguments:
213         dialogs {List[Dialog]} -- Список диалогов
214         cluster {MessageClusterizer} -- Алгоритм кластеризации исходящих
↳ сообщений.
215
216     Returns:
217         messages {List[Message]} -- Список сообщений
218         cluster_labels {np.array} -- Список меток кластера для сообщений
219     """
220
221     start = time.time()
222     clustered_messages, cluster_labels = cluster_prev_and_next(
223         dialogs, cluster)
224     print("Init clustering:", time.time() - start, file=sys.stderr)
225     start = time.time()
226     clustered_messages, cluster_labels = recluster_incoming(
227         clustered_messages, cluster_labels, dialogs, cluster, **hparams)
228     print("Overall recluster:", time.time() - start, file=sys.stderr)
229     return clustered_messages, cluster_labels

```

```

1 # Copyright (C) 2020 Dasha.AI Inc.
2 import json
3 from collections import defaultdict
4
5 class Graph:
6     def __init__(self, messages, digressions=[]):
7         msggs_dict = defaultdict(set)

```

```

8         self.vertexes = []
9         for msg in messages:
10             msgs_dict[msg.cluster].add(msg)
11         for cluster, msgs in msgs_dict.items():
12             tr_dict = defaultdict(set)
13             transactions = []
14             for msg in msgs:
15                 if msg.next_msg:
16                     tr_dict[msg.next_msg.cluster].add(msg)
17             for next_cluster, msgs_tr in tr_dict.items():
18                 transactions.append(Transaction(msgs_tr, int(cluster),
19                     ↪ int(next_cluster)))
20             self.vertexes.append(Vertex(int(cluster), transactions, cluster in
21                 ↪ digressions))
22         self.digressions = digressions
23
24     def __str__(self):
25         str_vert = []
26         for vertex in self.vertexes:
27             str_vert.append(str(vertex))
28         return '{"nodes": {' + ", ".join(str_vert) + '},' +
29             ↪ str(self.get_digr()) + '}'
30
31     def get_digr(self):
32         digr_tr = None
33         return '"digressionNodes": ' + json.dumps(
34             {
35                 "dig:root": {
36                     "Id": "dig:root",
37                     "OnEnter": {"Type": "None"},
38                     "Transitions": [digr_tr]
39                 }
40             })
41
42     def str_transactions_msgs(self):
43         res = ''
44         for vertex in self.vertexes:
45             for transition in vertex.transitions:
46                 res += transition.str_msgs()
47         return res
48
49     class Vertex:
50         def __init__(self, cluster_n, transitions, is_digressions=False):
51             self.cluster_n = cluster_n
52             self.transitions = transitions
53             self.is_digressions = is_digressions
54
55         def __str__(self):
56             str_trans = []

```

```

55     for transition in self.transitions:
56         str_trans.append(str(transition))
57     return '"{}":'.format(self.cluster_n) + json.dumps({
58         "Id": str(self.cluster_n),
59         "OnEnter": {
60             "Type": "Chain",
61             "InnerReactions": [
62                 {
63                     "Type": "Simple",
64                     "Reactions": [
65                         {
66                             "msgId": "RawTextChannelMessage",
67                             "text": 'text{}'.format(self.cluster_n)
68                         }
69                     ]
70                 }
71             ],
72             "Transitions": [
73                 [json.loads(json_tr) for json_tr in str_trans]
74             ]
75         }
76     })
77
78     class Transaction:
79         def __init__(self, msgs_set, from_cluster, to_cluster):
80             self.msgs = msgs_set
81             self.from_cluster = from_cluster
82             self.to_cluster = to_cluster
83             self.tr_name = 'fr_' + str(from_cluster) + '_to_' + str(to_cluster)
84
85         def __str__(self):
86             return json.dumps({
87                 "Condition": {
88                     "Type": "Fact",
89                     "FactName": "input_info",
90                     "InnerCondition": {
91                         "Type": "Field",
92                         "FieldName": "type",
93                         "InnerCondition": { "Type": "ItemEquals", "Item": self.tr_name
94                             ↪      }
95                     },
96                     "Id": str(self.to_cluster),
97                     "Priority": 0
98                 }
99             })
100
101         def compare(self, tr):
102             self.msgs.intersection(tr.msgs)
103
104         def str_msgs(self):

```

```

104     msgs = []
105     for msg in self.msgs:
106         msgs.append('{ "label":' + str(msg.labels or '{}') + ', "msg":' +
            ↪ '"{}"'.format(msg.text) + '}')
107     return '{ "name":' + '"{}"'.format(self.tr_name) + ', "msgs": [' + ",
            ↪ ".join(msgs) + ']'

```

```

1  # Copyright (C) 2020 Dasha.AI Inc.
2  from collections import defaultdict
3  from reconstruction.data_structure.dialog import Dialog
4
5  import numpy as np
6
7
8  def bfs_digressions_finder_0(graph, node, set_of_outgoing_clusters,
    ↪ set_of_incoming_clusters):
9      """Trying to get tail from expected root, and cut vertices
10         which can be potential continue after digressions. As a potential and
11         we mark vertices befor root and next vertexes.
12
13         Have troubles because of cycles, which does not exists in real dialogs"""
14     first_step = set_of_outgoing_clusters[node]
15     queue = set()
16     marked = set()
17     for node in first_step:
18         queue = queue | set_of_outgoing_clusters[node]
19         potential_end = set_of_incoming_clusters[node]
20     for source in set_of_incoming_clusters[node]:
21         potential_end = potential_end | set_of_outgoing_clusters[source]
22     from_q = queue.copy()
23     while queue:
24         cur_node = queue.pop()
25         for next_node in set_of_outgoing_clusters[cur_node]:
26             if (next_node not in potential_end and next_node not in from_q):
27                 from_q.add(next_node)
28                 queue.add(next_node)
29     marked = from_q.copy()
30     # coming to digr part only from dirg + add to marked + add starts
31     # queue = {graph["start"]}
32     # from_q = queue
33     # while queue:
34     #     cur_node = queue.pop()
35     #     print("Queue: " + len(queue))
36     #     for next_node in set_of_outgoing_clusters[cur_node]:
37     #         if (next_node != node and (next_node not in from_q)):
38     #             from_q.add(next_node)
39     #             queue.add(next_node)
40     # marked = marked.difference(from_q)
41     if len(marked) > 0:

```

```

42     graph['digressions'].append({
43         "root": node,
44         "nodes": list(marked)
45     })
46
47
48 def potential_digr_ver(unachive, dialogs_num, dialogs_out_of_dig):
49     possible_list = list() # list with possible errors
50     for cluster in unachive:
51         if (len(dialogs_out_of_dig[cluster]) / dialogs_num[cluster] > 0.8):
52             possible_list.append({
53                 "cluster": cluster,
54                 "ex_list": dialogs_out_of_dig[cluster]
55             })
56     return possible_list
57
58
59 def digressions_finder(graph, node, dialogs):
60     cluster_labels = set()
61     achieve = set()
62     dialogs_num = defaultdict(int)
63     dialogs_out_of_dig = defaultdict(list)
64     for dialog in dialogs:
65         stop = True
66         for msg in dialog.log:
67             dialogs_num[msg.cluster] += 1
68             if (msg.cluster != node) and stop:
69                 achieve.add(msg.cluster)
70                 dialogs_out_of_dig[msg.cluster].append(msg.text)
71             else:
72                 stop = False
73             cluster_labels.add(msg.cluster)
74     unachive = cluster_labels - achieve
75     if len(cluster_labels) / 2 > len(unachive):
76         graph['digressions'].append({
77             "root": node,
78             "nodes": list(unachive),
79             "possible": potential_digr_ver(unachive, dialogs_num,
80                 ↪ dialogs_out_of_dig)
81         })
82
83 def digressions_by_different_income(graph, dialogs):
84     """Find digressions with finding potential starts and cut with some
85     ↪ heuristic"""
86     set_of_incoming_clusters = defaultdict(set)
87     for link in graph["links"]:
88         set_of_incoming_clusters[link["target"]].add(link["source"])
89     num_of_incoming_clusters = dict()
90     for node in graph["nodes"]:

```



```

90     cluster = node["id"]
91     num_of_incoming_clusters[cluster] = len(
92         set_of_incoming_clusters[cluster])
93     if (num_of_incoming_clusters[cluster] > 3): # random const
94         if (not node["incoming"]):
95             digressions_finder(graph, cluster, dialogs)
96
97
98 def find_similar_nodes(nodes, income_shared):
99     nodes_sets = defaultdict(set)
100     potential_merge = [[]]
101     for node in nodes:
102         for transaction in node.transitions.keys():
103             if transaction in income_shared:
104                 nodes_sets[node.cluster_n].add(transaction)
105     for cluster1, set1 in nodes_sets.items():
106         if (len(set1) > 0):
107             nodes = []
108             for cluster2, set2 in nodes_sets.items():
109                 if (len(set1 & set2)) / (len(set1 | set2)) > 0.7 and len(set1)
110                     ↳ * len(set2) > 1 and cluster1 != cluster2:
111                     nodes.append(cluster2)
112             if (len(nodes) > 0):
113                 nodes.append(cluster1)
114             if (len(nodes) > 1):
115                 potential_merge.append(nodes)
116     return potential_merge
117
118 def cut_cycle(dialog, start, end, messages, remove_marks):
119     log = dialog.log
120     if len(log) > end + 1:
121         log[start].next_msg = log[end + 1]
122         log[end + 1].prev_msg = log[start]
123     else:
124         log[start].next_msg = None
125     new_log = log[: start + 1] + log[end + 1:]
126     counter = start
127     for msg in log[start + 1: end]:
128         counter += 1
129         if msg in messages:
130             index = messages.index(msg)
131             remove_marks.append(index)
132     trigger_msg = log[start + 1]
133     end_msg = log[end]
134     while trigger_msg != end_msg and not trigger_msg.incoming:
135         if trigger_msg.next_msg:
136             trigger_msg = trigger_msg.next_msg
137         else:
138             trigger_msg = None

```

```

139     return Dialog(dialog._id, dialog.task, new_log), trigger_msg, log[start +
    ↪ 1: end]

140
141
142 def cycle_finder(dialogs, messages):
143     """Find digressions with finding cycles"""
144     remove_marks = []
145     new_dialogs = []
146     trigger_list = []
147     remove_part = []
148     for dialog in dialogs:
149         new_dialog = dialog
150         have_cycle_flag = True
151         cycles_counter = 0
152         while have_cycle_flag:
153             cycles_counter += 1
154             have_cycle_flag = False
155             nodes_map = {}
156             counter = 0
157             for msg in new_dialog.log:
158                 cluster = msg.cluster
159                 if cluster in nodes_map and counter - nodes_map[cluster] < 5:
160                     #print("Cycle №", cycles_counter)
161                     have_cycle_flag = True
162                     new_dialog, trigger, remove_part_local = cut_cycle(
163                         new_dialog, nodes_map[cluster], counter, messages,
164                         ↪ remove_marks)
165                     remove_part += remove_part_local
166                     if trigger:
167                         trigger_list.append(trigger)
168                     break
169                 if not msg.incoming:
170                     nodes_map[cluster] = counter
171                 counter += 1
172             new_dialogs.append(new_dialog)
173             trigger_clusters = defaultdict(int)
174             for trigger in trigger_list:
175                 trigger_clusters[trigger.cluster] += 1
176             sorted_dict = {k: v for k, v in sorted(
177                 trigger_clusters.items(), key=lambda item: item[1])}
178             print("TRIGGERS STAT:")
179             for cluster, num in sorted_dict.items():
180                 print("        cluster: ", cluster, " num: ", num, " in ",
    ↪ sum(map(lambda x: x.cluster == cluster, messages)))
181     return new_dialogs, remove_marks, remove_part

```

```
1 # Copyright (C) 2020 Dasha.AI Inc.
2 import re
3
4 from reconstruction.embedding.int_to_str import convert_int
5
6
7 def norm_text(text):
8     text = re.sub(r'ё', 'e', text).lower()
9     text = re.sub(r'!"|!|\?|,|\.|%', ' ', text)
10    text = re.sub(r'-' , ' ', text)
11    text = re.sub(r':', ' ', text)
12    text = re.sub(r'::', ' ', text)
13    replace_dict = {
14        'ненадо': 'не надо',
15        "ололо": 'ало',
16        "алло": 'ало',
17        "але": 'ало',
18        "ален": 'ало',
19        "дадада": "да да да",
20        "дадад": "да да да",
21        "додо": "да да",
22        "дада": "да да",
23        "\bпасибо\b": "спасибо",
24        "\bпасиб\b": "спасибо",
25        "пасиба": "спасибо"
26    }
27    for from_w, to_w in replace_dict.items():
28        text = re.sub(from_w, to_w, text)
29
30    text_arr = []
31    for word in text.split():
32        try:
33            converted_text = convert_int(word)
34        except ValueError:
35            converted_text = word
36        text_arr.append(converted_text)
37
38    return " ".join(text_arr)
```

ПРИЛОЖЕНИЕ Б. ДИПЛОМ КОНГРЕССА МОЛОДЫХ УЧЁНЫХ



УНИВЕРСИТЕТ ИТМО

Центр студенческой науки,
конференций и выставок

IX КОНГРЕСС МОЛОДЫХ УЧЕНЫХ

ДИПЛОМПОБЕДИТЕЛЯ КОНКУРСА ДОКЛАДОВ
ДЛЯ ПОСТУПЛЕНИЯ В МАГИСТРАТУРУ

награждается

**ЮЛИЯ
КОНСТАНТИНОВНА
САВОН**Ректор
Университета ИТМО

В. Н. Васильев

Санкт-Петербург, 2020 год