



**Факультет Кибернетики и Информационной безопасности
КАФЕДРА КИБЕРНЕТИКИ (№ 22)**

Направление подготовки 09.04.04 Программная инженерия

Пояснительная записка

к научно-исследовательской работе студента на тему:

**Модернизация модели виртуального Агента с целью добавления
эмоционального взаимодействия**

Группа _____ М20-504

Студент _____ Клычков М. Д.

Руководитель _____ Самсонович А. В.

Научный консультант _____

Оценка руководителя _____ Оценка комиссии _____

Члены комиссии

Москва 2022



**Факультет Кибернетики и Информационной безопасности
КАФЕДРА КИБЕРНЕТИКИ (№ 22)**

Задание на НИР

Студенту гр. М20-504 Клычкову Матвею Дмитриевичу

ТЕМА НИР

**Модернизация модели виртуального Агента с целью добавления
эмоционального взаимодействия**

ЗАДАНИЕ

№ п/п	Содержание работы	Форма отчетности	Срок исполнения	Отметка о выполнении Дата, подпись
1.	Аналитическая часть			
1.1.	Изучение и анализ существующих когнитивных архитектур	Пункт ПЗ		
1.2.	Изучить методы машинного обучения	Пункт ПЗ		
1.3.	Изучение и анализ проблемной области распознавания речи	Пункт ПЗ		
1.4.	Классификация и определение эмоций	Пункт ПЗ		
1.5.	Оформление расширенного содержания пояснительной записи (РСПЗ)	Текст РСПЗ	10.04.2022	
2.	Теоретическая часть			
2.1.	Описать модель работы виртуального актора			
2.2.	Разработка алгоритмов распознавания речи	Алгоритмы		
2.3.	Модификация модификация методов машинного обучения применительно к задаче распознавания речи	Алгоритмы		
3.	Инженерная часть			
3.1.	Построить модель тембора голоса	Модели		
3.2.	Построить семантическую модель распознавания голоса	Модели		
3.3.	Проектирование инструментов для анализа текста	Макеты		
3.4.	Проектирование приложения	Макеты		
3.5.	Результаты проектирования оформить с помощью UML диаграмм	UML диаграммы		
4.	Технологическая и практическая часть			
4.1.	Реализация модели виртуального актора	Исполняемые файлы, исходный текст		

4.2.	Реализация программного приложения	Исполняемые файлы, исходный текст		
4.3.	Реализация и дообучение моделей машинного обучения	Исполняемые файлы, исходный текст		
5.	Оформление пояснительной записи (ПЗ) и иллюстративного материала для доклада.	Текст ПЗ, презентация	06.05.2022	

ЛИТЕРАТУРА

- Tikhomirova, D. V., Chubarov, A. A., Samsonovich, A. V. (2019). Empirical and modeling study of emotional state dynamics in social videogame paradigms. Cognitive Systems Research.
- Samsonovich A.V. Comparative Analysis of Implemented Cognitive Architectures// Biologically Inspired Cognitive Architectures. 2011. Vol. 233. P. 469-479. doi: 10.3233/978-1-60750-959-2-469
- Larue, O., West, R., Rosenbloom, P. S., Dancy, C. L., Samsonovich, A. V., Petters, D., Juvina, I. (2018). Emotion in the common model of 74 A.V. Samsonovich / Cognitive Systems Research 60 (2020) 57–76 cognition. Procedia Computer Science, 145, 740–746.
- Madl, T., Franklin, S., Chen, K., Trappl, R. (2018). A computational cognitive framework of spatial memory in brains and robots. Cognitive Systems Research, 47, 147– 172.
- Laird, J. E., Lebiere, C., Rosenbloom, P. S. (2017). A standard model of the mind: Toward a common computational framework across artificial intelligence, cognitive science, neuroscience, and robotics. AI Magazine, 38(4), 13–26./

Дата выдачи задания:
10.10.2021

Руководитель
Студент

Клычков М. Д.
Самсонович А. В.

Реферат

Пояснительная записка содержит 46 страниц, 19 рисунков, – 32 источников литературы.

Ключевые слова: UNITY3D, eBICA, СОЦИАЛЬНО-ЭМОЦИОНАЛЬНЫЙ ИНТЕЛЛЕКТ, ВИРТУАЛЬНЫЙ АКТОР, МАШИННОЕ ОБУЧЕНИЕ, РЕККУРЕНТНЫЕ СЕТИ

Объектом исследования являются экспертные системы.

Предмет исследования - модель когнитивной архитектуры для создания Виртуального Актора.

Целью данной научно-исследовательской работы является создание прототипа Виртуального Актора, модель интеллекта которого основана на когнитивной архитектуре eBICA. Существует два глобальных подхода к созданию социально-эмоционального интеллекта. Один основанный на нейросетях, другой на когнитивных архитектурах. В ходе работы над НИР был разработан и протестирован с участием испытуемых прототип Виртуального Актора, обладающий социально-эмоциональным интеллектом и помещенный в виртуальное окружение, которое создано при помощи графического движка Unity3d, была реализована система для анализа речи с выявлением эмоций соответствующим речевым признакам, так же были в дополнении к вышеуказанной системе, были спроектированы и реализованы методы воздействия на виртуального агента, основывающиеся на семантической составляющей речевого контекста. Данная работа является актуальной поскольку на данный момент эта область находится на начальных этапах развития и активной интеграции в различные индустрии. Созданная и протестированная модель интеллекта затем может быть интегрирована в другие проекты с Виртуальным Актором: виртуальный слушатель, виртуальный клоун, виртуальный танцор.

Содержание

Введение	4
1 Исследование существующих когнитивных архитектур и анализ их недостатков	5
1.1 Изучение и анализ существующих когнитивных архитектур	5
1.2 Изучение и анализ когнитивной архитектуры eBICA	10
1.3 Нейронные сети и их типы	12
1.4 Синтез и распознавание речи	14
1.5 Классификации и определение эмоций	16
1.6 Выводы	17
2 Описание моделей, отвечающих за генерацию поведения виртуального актора	18
2.1 Постановка задачи	18
2.2 Описание работы модели актора	19
2.3 Распознавание речи	22
2.4 Рекуррентные нейронные сети	25
2.5 Выводы	25
3 Проектирование модели поведения виртуального агента	26
3.1 Проектирование модифицированного прототипа Виртуального Актора	26
3.2 Модель тембора	28
3.3 Модель семантики	30
3.4 Инструменты для анализа текста	30
3.5 Построение блок-схем и UML диаграмм	33
3.6 Выводы	36
4 Реализация программного продукта	37
4.1 Текущая версия реализованной модели	37
4.2 Интеграция моделей из python в C#	38
4.3 Реализация и дообучение модели	39
4.4 Выводы	40
Заключение	41

Список литературы	42
Список литературы	43

Введение

В последнее время все большую и большую популярность набирают технологии предоставляющие возможность участвовать человеку в виртуальном мире либо технические средства, позволяющие представление виртуальной реальность в реальном мире. Виртуальные Акторы способны в будущем заменить докладчика на конференциях различного характера и представлениях.

Целью исследования является создание общей вычислительной модели механизмов, лежащих в основе человеческих эмоций. Осуществляется это путем создания Виртуального агента, подкрепленного когнитивной архитектурой и помещенного в виртуальное окружение. В данной парадигме человек (испытуемый, проводящий сеанс игры) может взаимодействовать с Виртуальным Актором, воплощенного в виде аватара в игре с трехмерной графикой, и оценивать его по различным параметрам. Такая модель может интерпретировать человеческое поведение и на основе экспериментов можно делать выводы о ее социальной приемлемости и точности имитирования человеческих эмоций.

В первом разделе обозначены когнитивные архитектуры, выявляются их преимущества и недостатки по сравнению с когнитивной архитектурой eBICA. Определяются типы распознавания и синтеза речи. Ставятся цели и задачи научно-исследовательской работы.

Во втором производится анализ когнитивных архитектур и производится анализ методов распознавания человеческой речи, а так же выявление в ней эмоциональных составляющих. Так же приводятся теоретические выкладки описания модели социально-эмоционального интеллекта.

В третьем разделе приводится подробное описание работы алгоритма, реализующего когнитивную модель социально-эмоционального интеллекта с учетом интегрированной системы распознавания речи. Строятся блок-схемы для кодовой реализации модели интеллекта.

В четвертом разделе приводятся реализация программного продукта представляющего собой когнитивно эмоциональный интеллект, способный взаимодействовать на эмоциональной основе.

1. Исследование существующих когнитивных архитектур и анализ их недостатков

Проводится анализ по выявлению существующих недоработок прототипа. Выявляются недостатки и преимущества по сравнению с другими моделями искусственного интеллекта.

1.1 Изучение и анализ существующих когнитивных архитектур

Одной из наиболее известных когнитивных архитектур является архитектура, составленная Jonathan Gratch и Stacy Marsella, что описано в работе [1]. Цель их исследования - создать общую вычислительную модель механизмов, лежащих в основе человеческих эмоций, которая сможет всецело их описать. Хотя такая модель может давать объяснение человеческого поведения, они рассматривают разработку вычислительных моделей эмоций как ключевой объект исследований для искусственного интеллекта, который будет способствовать развитию большого количества вычислительных систем, которые моделируют, интерпретируют или влияют на человеческое поведение. На рисунке (Рис. 1.1) демонстрируется Когнитивно-мотивационно-эмоциональная система.

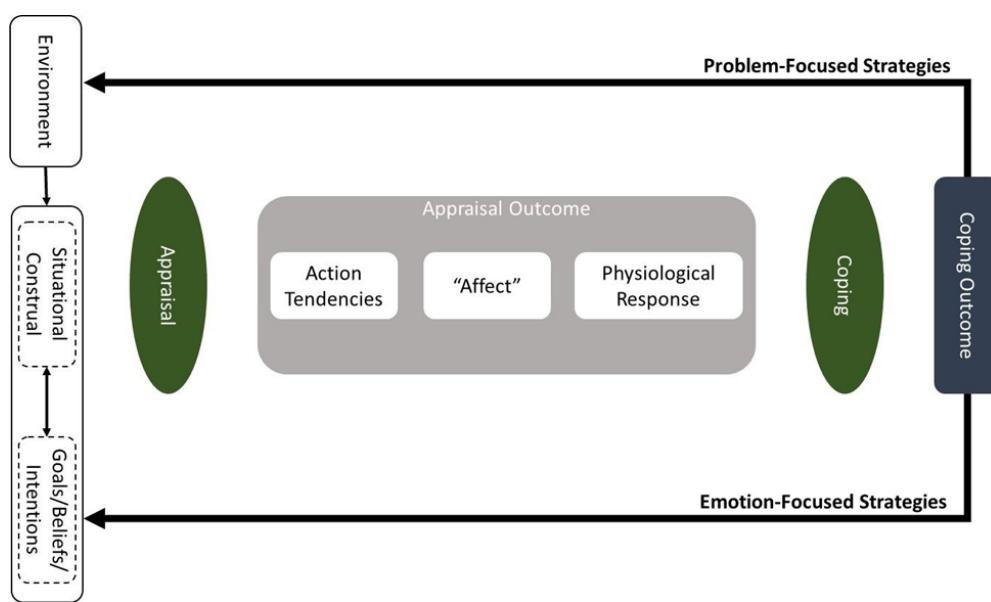


Рис. 1.1 – Когнитивно-мотивационно-эмоциональная система по материалам Smith and Lazarus.

Теория оценки служит концептуальной основой их работы, но эта психологическая теория недостаточно точна, чтобы служить спецификацией вычислительной модели. Для этого они

переделывают теорию с точки зрения методов и представлений искусственного интеллекта. Когнитивно-мотивационно-эмоциональная система Craig Smith и Richard Lazarus, показанная на рисунке 1, является представителем современных теорий оценки. Эмоция концептуализируется как двухступенчатая система контроля. Оценка характеризует отношения между человеком и его физическим и социальным окружением, называемые отношениями человека и окружающей среды, копирование поведения для восстановления или поддержания этих отношений. Поведение возникает в результате тесной связи познания, эмоций и реакций совпадения: когнитивные процессы служат для построения индивидуальной интерпретации того, как внешние события соотносятся с его целями и желаниями (отношения человека и окружающей среды). Система использует эти характеристики для изменения отношений между человеком и окружающей средой, мотивируя действия, которые изменяют среду (копирование, ориентированное на проблему), или мотивируя изменения в интерпретации этих отношений (копирование, ориентированное на эмоции).

Модель PAD была разработана Albert Mehrabian и James A. Russell в 1974 году для описания и измерения эмоциональных состояний, как говорится в Работе [2]. В данной модели используются три числовых измерения для представления всех эмоций:

- A – arousal (возбуждение);
- P – pleasure (удовольствие);
- D – dominance (доминирование).

Модель PAD первоначально использовалась в теории психологии окружающей среды, а основное идеей модели было предположение о том, что физическая среда влияет на людей через их эмоциональное воздействие. На основе данной модели были построены физиологическая теория эмоций и теория эмоциональных эпизодов. Также модель использовалась для изучения неверbalного общения, в потребительском маркетинге и при создании анимированных персонажей, которые выражают эмоции.

В модели PAD используются трехмерные шкалы, которые в теории могут иметь любые числовые значения:

- шкала удовольствия-неудовольствия показывает, насколько приятно или, наоборот, неприятно человек себя чувствует по отношению к чему-то. Например, радость это – приятная эмоция; гнев и страх – неприятные эмоции;
- шкала возбуждения-неактивности измеряет, насколько человек чувствует возбуждение или его отсутствие. В данном случае оценивается именно возбуждение, а не интенсивность эмоций. Например, горе или депрессия характеризуются слабым возбуждением, но

сильной интенсивностью; а гнев или ярость имеют и высокую интенсивность, и высокое состояние возбуждения;

- шкала доминирования-покорности описывает чувство контроля и доминирования по сравнению со смирением и подчиненностью. Например, гнев — это доминирующая эмоция, а страх
- эмоция покорности, хотя обе они имеют неприятный характер.
- эмоция покорности, хотя обе они имеют неприятный характер.

Еще одна интересная когнитивная архитектура описана в статье трех научных деятелей Ron Sun, Nick Wilson, Michael Lynch. Статья имеет название: “Emotion: A Unified Mechanistic Interpretation from a Cognitive Architecture”. В этой статье рассматривается проект, который пытается интерпретировать эмоции - сложное и многогранное явление с механистической точки зрения, чему способствует существующая комплексная вычислительная когнитивная архитектура - CLARION. Эта когнитивная архитектура состоит из ряда подсистем: подсистем, ориентированных на действие, не ориентированных на действие, мотивационной и метакогнитивной подсистем. С этой точки зрения эмоции в первую очередь основаны на мотивации. Основываясь на этих функциональных возможностях, мы механистически (вычислительно) соединяем части вместе в рамках CLARION и фиксируем множество важных аспектов эмоций, как описано в литературе. На (Рис. 1.2) демонстрируются подсистемы когнитивной архитектуры CLARION.

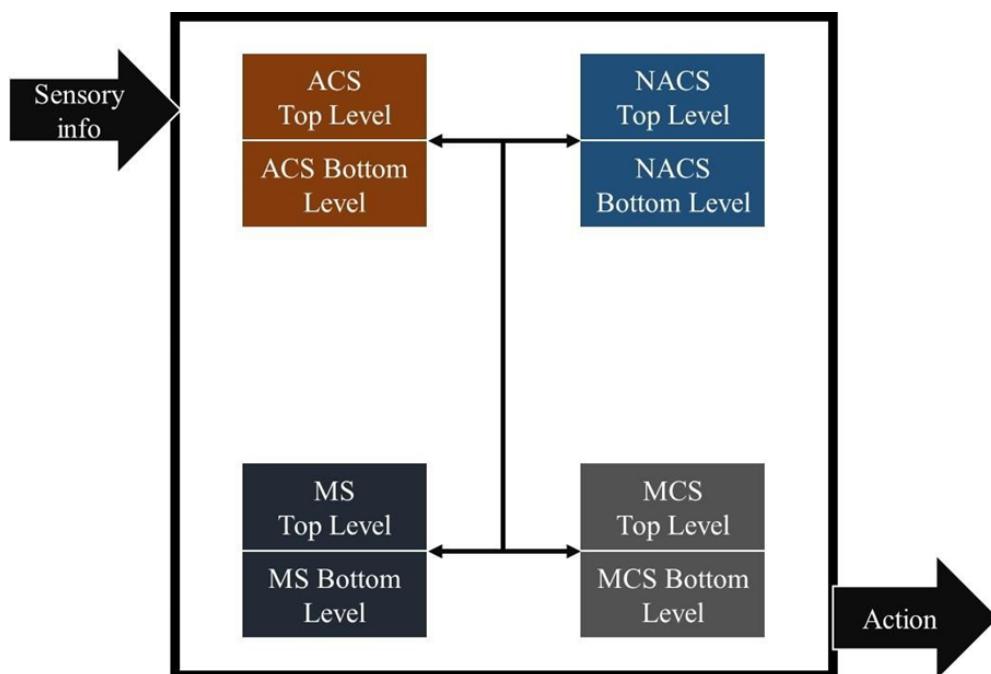


Рис. 1.2 – Подсистемы когнитивной архитектуры CLARION

Основные информационные потоки показаны стрелками. ACS означает подсистему, ори-

ентированную на действия. NACS означает подсистему, не ориентированную на действия. MC — это мотивационная подсистема. MCS означает метакогнитивную подсистему.

Получившая наибольшее распространение из всех формальных моделей представления эмоций является модель OCC (Ortony, Clore, & Collins), которая упоминается в работе [3], предложенная в 1988 году учеными Кембриджского университета. Иерархия содержит три ветви, а именно: эмоции, касающиеся последствий событий (например, радость и жалость), действия агентов (например, гордость и упрек) и аспекты объектов (например, любовь и ненависть). Кроме того, некоторые ветви объединяются в группу сложных эмоций, а именно эмоций относительно последствий событий, вызванных действиями агентов (например, благодарность и гнев). На рисунке (Рис. 1.3) демонстрируется оригинальная модель OCC.

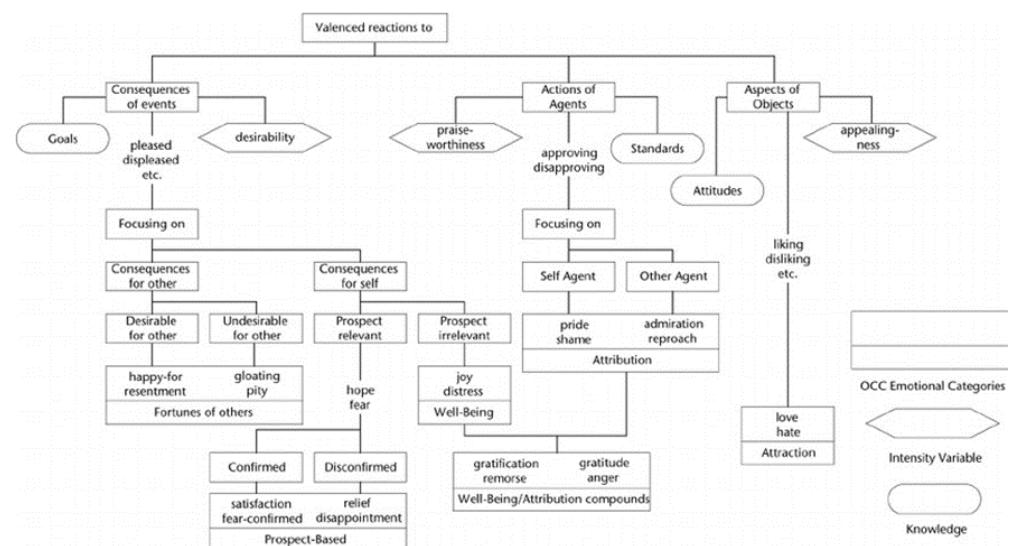


Рис. 1.3 – Оригинальная модель OCC

В основе правил динамики данной модели лежит реакция валентности (Valenced reaction). Под «валентностью» в психологии понимают внутреннюю привлекательность – «хорошую» (положительную валентность) или отвратительность – «плохую» (отрицательную валентность) события, объекта или ситуации. Эмоции формируются под воздействием трех основных факторов — последствий событий (Consequences of events), действий движущих сил (Actions of agents) и аспектов событий.

Рассмотрим левую «ветку» эмоциональной реакции. Последствия событий могут быть приносящими удовольствие (pleased) или доставляющими неудовольствие (displeased). Проведя предварительную оценку, человек фокусируется (Focusing on) на разделении последствий событий для себя (Consequences for self) и для других (Consequences for other), которые, в свою очередь могут оказаться для последних желательными (Desirable for other) или нежелательными (Undesirable for other). По поводу судеб других (Fortunes of others) в зависимости от личного

отношения — положительного или отрицательного — человек может испытывать следующие эмоции: радость за другого (Happy for), обида (Resentment), злорадство (Gloating) или жалость (Pity).

У этой модели есть свои ограничения, заключающиеся как в ее требовании упрощения человеческих эмоций, так и в ее сложном подходе к тому, как надлежит выводить эмоциональные состояния конечных пользователей посредством интерпретации поведения человека через знаки и сигналы, транслируемые людьми. Использование этой модели в ее оригинальном описании затруднено отсутствием математического аппарата, в следствии чего многие исследователи в своих Виртуальных Акторах используют упрощённые версии данной модели.

Также большой интерес представляет когнитивная архитектура, реализованная в физическом роботе, под названием - интегрированное когнитивное универсальное тело (iCub). Это когнитивная архитектура, дизайн которой основан на существующих знаниях в области робототехники, вычислений, нейробиологии и психологии, целью которой является копирование некоторых когнитивных процессов человека для их включения в человекоподобных роботов.

Эта архитектура реализована в человекоподобном роботе. Он был разработан для исследования сообществом когнитивных систем. Кроме того, он имеет лицензию «Стандартная общественная лицензия GNU (GPL)», так что любой человек может свободно использовать все наработки по данному проекту. Данная архитектура реализована в человекоподобном роботе, который имеет 53 степени свободы. По размеру он похож на ребенка трех-четырех лет и ребенка в возрасте 2,5 лет по когнитивным способностям. Кроме того, он может ползать и сидеть. Некоторые особенности, которые описаны в работе [4]

- Не хватает семантической памяти, чтобы помочь ему обобщать события;
- Невозможно сформировать привычки;
- Он учится путем подражания, проб и ошибок;
- Обнаруживает, распознает и отслеживает человеческое лицо, наблюдая за его действиями;
- Действия основаны на жестах рук, таких как встряхивание и манипулирование объектами, например, толкание, подъем и опускание.
- Действия, наблюдаемые роботом, изучаются и сохраняются в базе данных в процессе обучения.

На рисунке (Рис. 1.4) представлена схема работы iCub.

Вспомогательным инструментом при создании актора, наделенного социально-эмоциональным интеллектом, может являться - имитация моторного обучения (IML). IML начинает

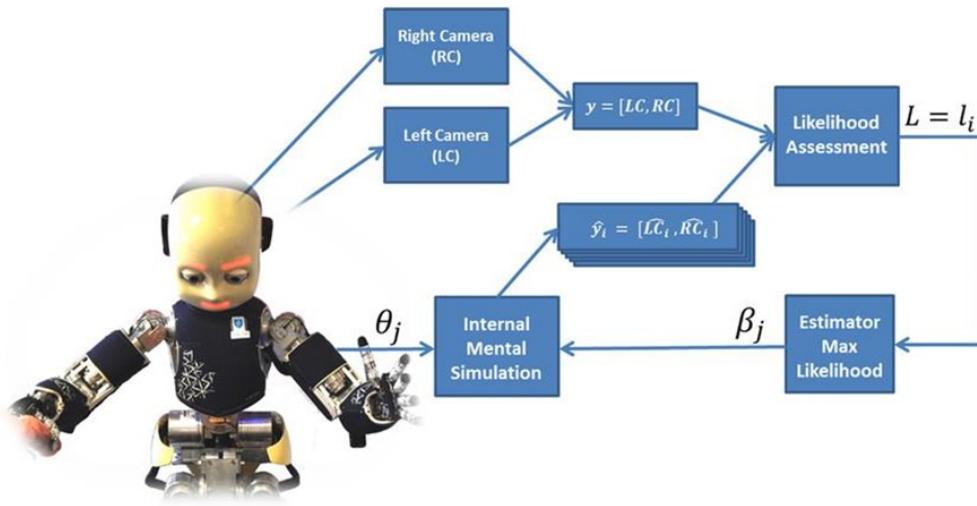


Рис. 1.4 – Схема работы iCub

наблюдать за другим актором, осуществляющим некоторую цепочку действий, затем категоризирует действия (определяет какую цель преследуют данные действия) одновременно отслеживая изменения точки обзора, окружающей среды, положения и типов объектов. Другими словами, когда Виртуальный агент неоднократно наблюдает за определенной новой последовательностью действий, каждый из знакомых элементов действия активирует соответствующее моторное представление через существующие ассоциации. Данное наблюдение формирует связи между элементарными моторными представлениями. Эта связь представлений составляет моторное обучение и улучшает имитационное движение. Способность моторной системы интегрировать разные части организма позволила бы создать обширный репертуар моторного поведения путем смешивания выходных сигналов разных частей организма, чтобы конечный результат отражал относительный и взвешенный вклад каждого в достижении цельной имитации движения. Поскольку невозможно воспроизвести функционирование мозга, были созданы модели, которые пытаются имитировать различные функции и поведение.

1.2 Изучение и анализ когнитивной архитектуры eBICA

Архитектура состоит из семи компонентов: интерфейсный буфер, рабочая, процедурная, семантическая и эпизодическая системы памяти, система ценностей и система когнитивных карт. Три основных строительных блока для этих компонентов – это ментальные состояния, схемы и семантические карты. Семантическая память – это коллекция определений схем. Буфер интерфейса заполняется схемами. Рабочая память включает активные психические состояния. Эпизодическая память хранит неактивные психические состояния, сгруппированные в эпизоды - предыдущее содержимое рабочей памяти. Следовательно, эпизодическая память состоит из структур, аналогичных тем, которые обнаруживаются в рабочей памяти, но которые

«заморожены» в долговременной памяти [5]. Процедурная память включает в себя примитивы. Система ценностей включает в себя шкалы, представляющие основные значения. Система когнитивных карт включает, в частности, семантические карты эмоциональных ценностей. Семантическая карта использует абстрактное метрическое пространство (семантическое пространство) для представления семантических отношений между ментальными состояниями, схемами и их 13 экземплярами, а также для присвоения значений их оценкам. На (Рис.1.5) демонстрируется семантическая карта [5].



Рис. 1.5 – Семантическая карта

Для когнитивного семантического отображения может использоваться слабое когнитивное семантическое картирование. Идея заключается в том, чтобы расположить представления на основе очень немногих основных семантических измерениях. Эти измерения могут возникать автоматически, если стратегия состоит в том, чтобы объединить синонимы и антонимы друг от друга. Карта, часть которой показана на рисунке 6 является результатом этого процесса. Эта карта не очень хорошо отделяет различные значения друг от друга: например, основные и сложные чувства. Однако она классифицирует значения в соответствии с их семантикой. Рисунок (Рис. 1.6) демонстрирует примеры простейших эмоциональных элементов в рамках eBICA [6].

(А) Схема имеет оценку в качестве своего атрибута. Это также атрибут головного узла. Значение этого атрибута - «доминантный», что означает, что действие воспринимается как проявление доминирования, или агент воспринимается как «доминантный по отношению ко мне» и т. Д. (В) Психическое состояние имеет оценку атрибут, который представляет собой эмоциональное состояние и самооценку агента в данный момент в данной ментальной перспективе.

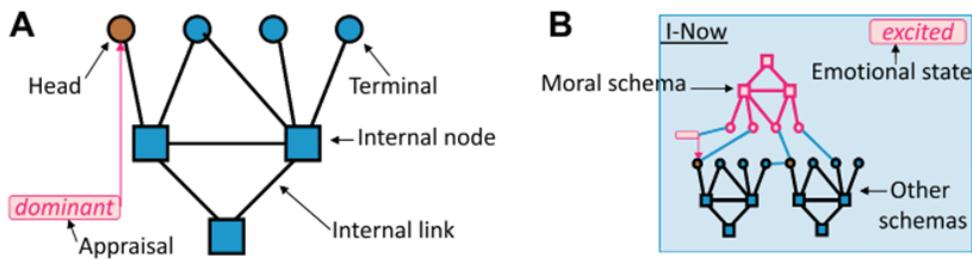


Рис. 1.6 – примеры эмоциональных элементов в рамках eBICA

Показанная ценность этой оценки «взволнована», что означает, что агент находится в возбужденном эмоциональном состоянии. Моральная схема, показанная в В, связывается с частью содержания психического состояния (включая определенный образец оценок) и представляет оценку выбранного образца, например, образец взаимодействий и взаимных оценок двух агентов, упомянутых в ментальном состоянии.

1.3 Нейронные сети и их типы

Нейронная сеть (также искусственная нейронная сеть, ИНС) – математическая модель, а также её программное или аппаратное воплощение, построенная по принципу организации и функционирования биологических нейронных сетей – сетей нервных клеток живого организма. Это понятие возникло при изучении процессов, протекающих в мозге, и при попытке смоделировать эти процессы. Первой такой попыткой были нейронные сети У. Маккалока и У. Питтса. После разработки алгоритмов обучения получаемые модели стали использовать в практических целях: в задачах прогнозирования, для распознавания образов, в задачах управления и др.

Разделяют несколько основных разновидностей Нейронных сетей, согласно работе [7], а именно:

- Нейронные сети прямого распространения
- Сети радиально-базисных функций
- Нейронная сеть Хопфилда (Hopfield network, HN)
- Цепи Маркова (Markov chains, MC или discrete time Markov Chains, DTMC)
- Машина Больцмана (Boltzmann machine, BM)
- Ограниченнная машина Больцмана (restricted Boltzmann machine, RBM)
- Автокодировщик (autoencoder, AE)
- Разреженный автокодировщик (sparse autoencoder, SAE)
- Вариационные автокодировщики (variational autoencoder, VAE)
- Шумоподавляющие автокодировщики (denoising autoencoder, DAE)

- Сеть типа «deep belief» (deep belief networks, DBN)
- Свёрточные нейронные сети (convolutional neural networks, CNN)
- Развёртывающие нейронные сети (deconvolutional networks, DN)

С точки зрения машинного обучения, нейронная сеть представляет собой частный случай методов распознавания образов, дискриминантного анализа. Рекуррентные нейронные сети (РНС, англ. Recurrent neural network; RNN) – вид нейронных сетей, где связи между элементами образуют направленную последовательность [8]. Благодаря этому появляется возможность обрабатывать серии событий во времени или последовательные пространственные цепочки. В отличие от многослойных перцептронов, рекуррентные сети могут использовать свою внутреннюю память для обработки последовательностей произвольной длины. Поэтому сети RNN применимы в таких задачах, где нечто целостное разбито на части, например: распознавание рукописного текста или распознавание речи. Было предложено много различных архитектурных решений для рекуррентных сетей от простых до сложных. В последнее время наибольшее распространение получили сеть с долговременной и кратковременной памятью (LSTM) и управляемый рекуррентный блок (GRU).

В последнее время наибольшую популярность для решения задач тематической классификации, применяемой для выделения семантического смысла текста, приобрели глубокие нейронные сети, так как они позволяют достичь наивысшей точности среди всех известных моделей машинного обучения. В частности, сверточные нейронные сети совершили прорыв в классификации изображений. В настоящее время они успешно справляются и с некоторыми задачами автоматической обработки текстов. Более того, как утверждается в некоторых исследованиях сверточные сети подходят для этого даже лучше рекуррентных нейронных сетей, которые чаще всего используются для анализа текстовых последовательностей. С другой стороны, использование сверточных сетей для классификации текстов мало исследовано. Поэтому исследование применения сверточных нейронных сетей для задачи классификации текстов в качестве альтернативы рекуррентным нейронным сетям представляет практический интерес, что описано в [9].

Для решения поставленной задачи требуется получить способ представления данных в виде, пригодном для обработки сверточной нейронной сетью. Например, в виде матрицы вещественных чисел. Наиболее распространенным является способ отображения каждого слова в многомерное векторное пространство. В рамках данной работы векторные представления слов строились на основе модели word2vec [10]. sa

1.4 Синтез и распознавание речи

Синтез речи - широком смысле – восстановление формы речевого сигнала по его параметрам; в узком смысле – формирование речевого сигнала по печатному тексту. Часть искусственного интеллекта. Синтезом речи – прежде всего называется всё, что связано с искусственным производством человеческой речи.

Первой технологией воспроизводящей синтез речи служила механическая говорящая машина Фабера (1845). Воздушный мех, приводимый в движение ножной педалью служил «лёгкими», а вытесняемый из межа воздух при помощи ряда клавиш направлялся в различные по объёму трубы - разные положения "голосовой щели" и "полости рта" (Рис. 1.7).

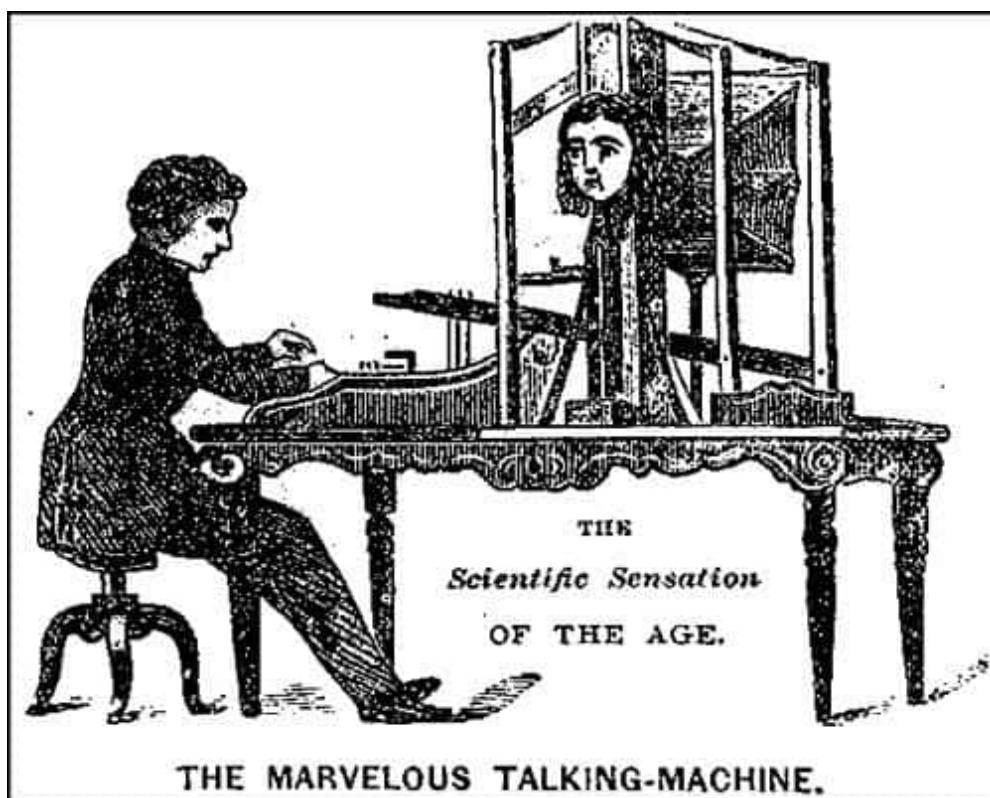


Рис. 1.7 – Машина Фабера 1845

Еще одной технологией служил так называемый "вокодер". Вокодер – синтезатор речи на основе произвольного сигнала с богатым спектром. Изначально вокодеры были разработаны в целях экономии частотных ресурсов радиолинии системы связи при передаче речевых сообщений. Вместо собственно речевого сигнала передают только значения его определённых параметров, которые на приёмной стороне управляют синтезатором речи (Рис. 1.8).

Основу синтезатора речи составляют три элемента:

- Генератор тонального сигнала для формирования гласных звуков;
- Генератор шума для формирования согласных;

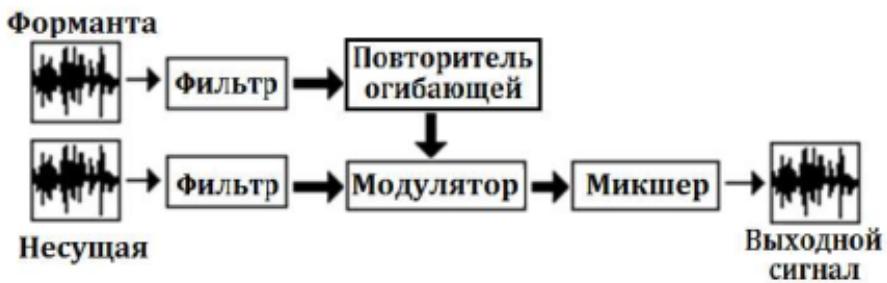


Рис. 1.8 – Вокодер

- Система формантных фильтров для воссоздания индивидуальных особенностей голоса.

Устная речь является собой практически непрерывный звуковой поток, представляющий определенные сложности для сегментирования. В большинстве случаев для выделения фонемы необходимо знание языка. Выделяемость фонемы некоторым образом сопряжена со смыслом и со значением, хотя сама по себе она не является значащей единицей.

На текущий момент синтез речи не представляет собой сложную задачу, так как в связи с темпами развития технологий, компьютеры позволяют хранить требуемую для синтеза информацию и соответственно синтезировать передаваемую фразу по посыпаемой машиной строкой текста, а также передавать нужную эмоциональную окраску в зависимости от семантики синтезируемого речевого фрагмента. Однако с темпами развития технологий, стали появляться задачи, когда машина реагирует не на скриптованное выполнение команд, а должна произвести речевой анализ получаемого текста и в ответ синтезировать ответ, для такого рода задач применяется распознавание речи.

Распознавание речи – автоматический процесс преобразования речевого сигнала в цифровую информацию. Обратной задачей является синтез речи. Распознавание речи – одна из самых интересных и сложных задач искусственного интеллекта. Здесь задействованы достижения весьма различных областей: от компьютерной лингвистики до цифровой обработки сигналов.

Основными задачами распознавания речи являются:

1. Системы Interactive Voice Response (IVR) в колл-центрах:

- Биометрическая идентификация
- Автоматическая маршрутизация звонка
- Аналитика речи, поиск пауз, тормозов в интерфейсе, аналитика эмоционального состояния
- Сбор оценок операторов колл-центра

2. Персональные помощники:

- Распознавание поисковых запросов
- Распознавание команд управления

Существует ряд технологий предоставляющих возможность осуществлять распознавание речи:

- Алиса (Yandex)
- iOS Siri (Apple/Nuance)
- Google Assistant
- Amazon Alexa

В данный момент для распознавания речи используется так называемая акустическая модель. По сути это функция, принимающая на вход небольшой участок акустического сигнала (фрейм) и выдающая распределение вероятностей различных фонем на этом фрейме. Таким образом, акустическая модель дает возможность по звуку восстановить, что было произнесено — с той или иной степенью уверенности.

Еще один важный аспект акустики — вероятность перехода между различными фонемами. Из опыта известно, что одни сочетания фонем произносятся легко и встречаются часто, другие сложнее для произношения и на практике используются реже. Можно обобщить эту информацию и учитывать ее при оценке «правдоподобности» той или иной последовательности фонем.

Однако акустическая модель — это всего лишь одна из составляющих современных систем. Словари на текущий момент состоят из миллионов, а то и триллионов слов, многие из них совпадают по своему звучанию и могут даже совпадать. Вместе с тем, при наличии контекста роль акустики падает: невнятно произнесенные, зашумленные или неоднозначные слова можно восстановить «по смыслу». Для учета контекста используются так же вероятностные модели, такой тип языковых моделей называется n-gram language models.

1.5 Классификации и определение эмоций

Многообразие эмоций, их качественных и количественных проявлений исключают возможность простой и единой классификации. Каждая из характеристик эмоций может выступать в качестве самостоятельного критерия, основания для их классификации (таб. 1.1).

Таблица 1.1 – характеристики эмоции как основания для их классификации

Знак	Положительные, отрицательные, амбивалентные
Модальность	Радость, гнев, страх и др.
Влияние на поведение и деятельность	Осознаваемые, неосознаваемые
Предметность	Предметные, беспредметные
Степень произвольности	Произвольные, непроизвольные
Происхождение и развлечения	Врожденные, приобретенные, первичные, вторичные
Уровень	Высшие, низшие
Длительность	Кратковременные, длительные
Интенсивность	Слабые, сильные

По знаку эмоциональные переживания можно разделить:

1. на положительные
2. отрицательные
3. амбивалентные

Основной функцией положительных эмоций является поддержание контакта с позитивным событием, поэтому им присуща реакция приближения к полезному, необходимому стимулу. Кроме того, по мнению П.В. Симонова, они побуждают нарушать достигнутое равновесие с окружающей средой и искать новую стимуляцию.

Для отрицательных эмоций характерной является реакция удаления, прерывания контакта с вредным или опасным стимулом. Считается, что они играют более важную биологическую роль, поскольку обеспечивают выживание индивида.

Амбивалентными эмоциями являются противоречивые эмоциональные переживания, связанные с двойственным отношением к чему-либо или кому-либо (одновременное принятие и отвержение).

1.6 Выводы

Были изучены и проанализированы основные когнитивные архитектуры, особое внимание уделялось когнитивной архитектуре eBICA. Была рассмотрена проблема синтеза и распознавание речи. Были изучены материалы описывающие классификацию и определение эмоций.

2. Описание моделей, отвечающих за генерацию поведения виртуального актора

В данном разделе приводится теоретическое описание модели.

2.1 Постановка задачи

В рамках научно-исследовательской работы был расширен подход решения поставленной задачи задачи, который выражается в использовании машинного обучения. Данный подход используется в совокупности когнитивной архитектурой eBICA. eBICA – “emotional biologically inspired cognitive architecture” – “эмоциональная биологически вдохновленная когнитивная архитектура”.

В этой архитектуре эмоциональные элементы добавлены практически ко всем процессам за счет модификации основных строительных блоков архитектуры. Ключевым моментом этой когнитивной архитектуры являются оценки, которые связаны со схемами и психическими состояниями как их атрибуты, моральные схемы, которые контролируют модели оценок и представляют социальные эмоции, а также семантические пространства, которые дают значения этих оценок.

Как видно из (Рис. 2.1), архитектура представляет собой конгломерат компонентов: интерфейсный буфер, рабочая, процедурная, семантическая и эпизодическая системы памяти, система ценностей и система когнитивных карт [6]. Три основных строительных блока для этих компонентов - это ментальные состояния, схемы и семантические карты. Семантическая память - это коллекция определений схем. Буфер интерфейса заполняется схемами.

Рабочая память включает активные психические состояния. Эпизодическая память хранит неактивные психические состояния, сгруппированные в эпизоды - предыдущее содержимое рабочей памяти. Следовательно, эпизодическая память состоит из структур, аналогичных тем, которые обнаруживаются в рабочей памяти, но которые «заморожены» в долговременной памяти. Процедурная память включает в себя примитивы. Система ценностей включает в себя шкалы, представляющие основные значения. Система когнитивных карт включает, в частности, семантические карты эмоциональных ценностей. Семантическая карта использует абстрактное метрическое пространство (семантическое пространство) для представления семантических отношений между ментальными состояниями, схемами и их экземплярами, а также для присвоения значений их оценкам.

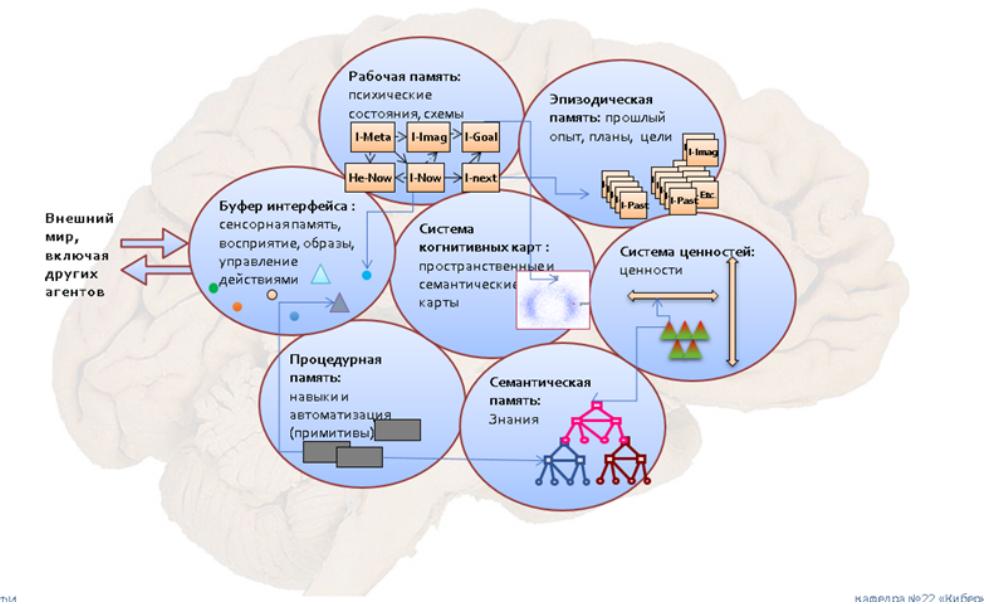


Рис. 2.1 – Структура когнитивной архитектуры eBICA

2.2 Описание работы модели актора

На момент начала выполнения работы уже была реализована система работы виртуальных агентов, и она состоит в том, что сперва считываются действия и объекты с заданными для них параметрами и значениями из Excel файла, а также инициализируются значения. Затем выбор действий происходит в следующем порядке: в радиусе вокруг оценок виртуального актера выбираются действия, которые попадают в этот радиус, а также проверяются различные условия, необходимые для выполнения действия. Если условия не выполняются, то действие не может быть выбрано. Затем, после того как в список добавлены все действия, которые могут быть выполнены, рассчитываются вероятности на основе оценок и рассчитанных констант. Также, если действие повторное, его вероятность несколько занижается. После расчета вероятностей выбирается действие, которое влияет на оценки виртуального актера.

Помимо этого происходит замена состояний объектов и виртуальных актеров. После этого происходит пересчет оценок Appraisals и Feelings [6].

Основная модель eBICA определяет поведение виртуального актера исходя из следующих факторов:

- соматический;
- рациональный;
- когнитивный.

Нравственный фактор регулирует отношения первого актера со вторым на основе системы ценностей (представленной семантической картой) и моральных схем. Под когнитивным фак-

тором понимается учет соображений нравственности, этики и морали, общей системы ценностей, понятий о добре и зле, о собственном достоинстве, эмпатии, соображений эстетики, стремлений к простоте и элегантности, и т.д. Учет этих соображений возможен на основе когнитивных оценок (appraisals) всех релевантных агентов, событий, их возможных действий и последствий этих действий, фактов, свойств, отношений, и т.д. Возможен вариант модели, в которой ответное действие может выбираться лишь из двух вариантов: положительная реакция на действие человека и отрицательная. Данная версия модели весьма неплохо работает даже с таким ограничением. Но невозможно придерживаться данной парадигмы при увеличении количества возможных вариантов для взаимодействия между акторами. В данной модели необходимо учесть пересчет оценок Appraisals и Feelings. Для пересчета оценок Appraisals используется следующая формула 2.1:

$$Appraisals = (1 - r) * Appraisals + r * Action \quad (2.1)$$

где Appraisals - оценка, r - эмпирически вычисленная константа экспоненциального затухания, Action - оценка совершающего действия на семантической карте.

Одновременно с Appraisals пересчитываются так называемые "чувства" Feelings согласно режиму работы моральной схемы. Аффективное пространство VAD – это трехмерное векторное пространство, точки которого соответствуют определенным эмоциональным состояниям, или аффектам, представленным триплетами значений (Valence, Arousal, Dominance). Существуют и сходные модели: PAD (Pleasure, Arousal, Dominance), EPA (Evaluation, Potency, Arousal) и другие. Здесь мы используем модель VAD. Соответственно, под «семантической картой» здесь часто понимается ее конкретная разновидность: аффективная карта (или когнитивная семантическая карта).

Шкалы имеют следующие значения:

- dominance – варьируется при значении от 0 (покорность) до +1 (доминантность) и описывает соответствующие чувства;
- valence – при значениях от -1 до 0 показывает уровень негатива или радости соответственно;
- arousal – значения от -1 до 1 показывают уровень возбуждения (заинтересованности), к примеру, гнев по уровню возбуждения сильнее раздражительности, но слабее ярости.

Оценки представлены в виде векторов на трехмерной семантической карте [5], 1.5. Моральная схема определяет общую установку на оценку поведения акторов, согласно их ролям и типу ситуации. Ее целью (как агента) является достижение и поддержание «нормального» по-

ложении дел, определенного набором Feelings. Вообще говоря, моральная схема состоит из двух частей: части, распознающей тип ситуации и осуществляющей привязку (binding), и части, реализующей динамику схемы. В случае парадигмы актора можно считать, что моральная схема одна, уже привязана, и потому первая часть ее не актуальна.

Субъективные оценки (Feelings) генерируются по определенным правилам на основании истории объективных оценок и состояний системы. Грубо говоря, Feelings – это субъективное представление о том, каким оцениваемый актор является «на самом деле», и, следовательно, какого поведения от него нужно ожидать и на какое место его нужно ставить своим поведением. Следовательно, выбор поведения актора должен осуществляться так, чтобы приблизить Appraisals к Feelings.

Значение Feelings определяет моральная схема, которая может работать в одном из трех режимов. Первый режим основывается на формуле 2.2:

$$Feelings = \beta * Appraisals \quad (2.2)$$

где β – эмпирически вычисленная константа.

В данном режиме схема говорит, что если актор ведет себя хорошо, то к нему нужно относиться как к хорошему, и т.д.

Цель данного процесса – распознать и классифицировать актора, выработать отношение к нему и приписать ему определенную роль во взаимоотношениях.

В данном режиме моральная схема работает пока разница между квадратами норм Feeling и Appraisals не станет меньше некоторого значения.

Суть второго режима заключается в том, что значение Feeling фиксировано и экстремально по абсолютной величине, т.е. находится на сфере, ограничивающей семантическую карту (предположим, что есть такая сфера). Направленность вектора Feeling может быть либо произвольной, определенной предысторией, либо дискретной – вдоль одной из осей.

Третий режим состоит в том, что значения Feelings меняются 2.3, подстраиваясь под текущие значения Appraisals (здесь r_1 может быть отличным от r):

$$Feelings = (1 - r_1) * Feelings + r_1 * (Appraisal - Feelings) \quad (2.3)$$

Соответственно значения Appraisals и Feelings как говорится в работе [Samsonovich05] пересчитываются после каждого действия первого актора, направленного на второго актора. Так же пересчет оценок происходит после определения и совершения одним из акторов ответного или самостоятельного действия. В данном контексте под термином “самостоятельное действие”

ствие” имеется в виду действие, основанное лишь на текущем состоянии мира и значений векторов Appraisals и Feelings акторов, отобранные на (Рис. 2.2).



Рис. 2.2 – Корреляция значений Appraisals (оранжевая) и Feelings (синяя) для показателя доминантности на протяжении времени/действий с шагом в 5 секунд

Согласно (Рис. 2.2) мы видим, что работа модели сводится к выбору действия, которое будет максимально приближать Appraisals к Feelings и вектор соматического состояния к начальному положению.

2.3 Распознавание речи

Как было упомянуто ранее, для распознавания речи используются акустические модели, однако этого может быть недостаточно, так как словари на текущий момент состоят из миллионов, а то и триллионов слов, многие из них совпадают по своему звучанию и могут даже совпадать и в написании и в зучании.

За основной параметр в использовании распознавания речи при использовании классической акустической модели - берется фонема, однако в слове она может находиться в трех различных состояниях: начало, середина и конец. Фонемы являются позиционно-зависимыми и контекстно- зависимыми: формально «одна и та же» фонема звучит существенно по-разному в зависимости от того, в какой части слова она находится и с какими фонемами соседствует. Вместе с тем, простое перечисление всех возможных вариантов контекстно- зависимых фонем вернет очень большое число сочетаний, многие из которых никогда не встречаются в реальной жизни; чтобы сделать количество рассматриваемых акустических событий разумным, близкие контекстно- зависимые фонемы объединяются на ранних этапах тренировки и рассматриваются вместе.

Тренировка акустической модели — сложный и многоэтапный процесс. Для тренировки используются алгоритмы семейства Expectation-Maximization, такие, как алгоритм Баума-Велша. Суть алгоритмов такого рода — в чередовании двух шагов: на шаге Expectation имеющаяся модель используется для вычисления матожидания функции правдоподобия, на шаге Maximization параметры модели изменяются таким образом, чтобы максимизировать эту оценку (Рис. 2.3).

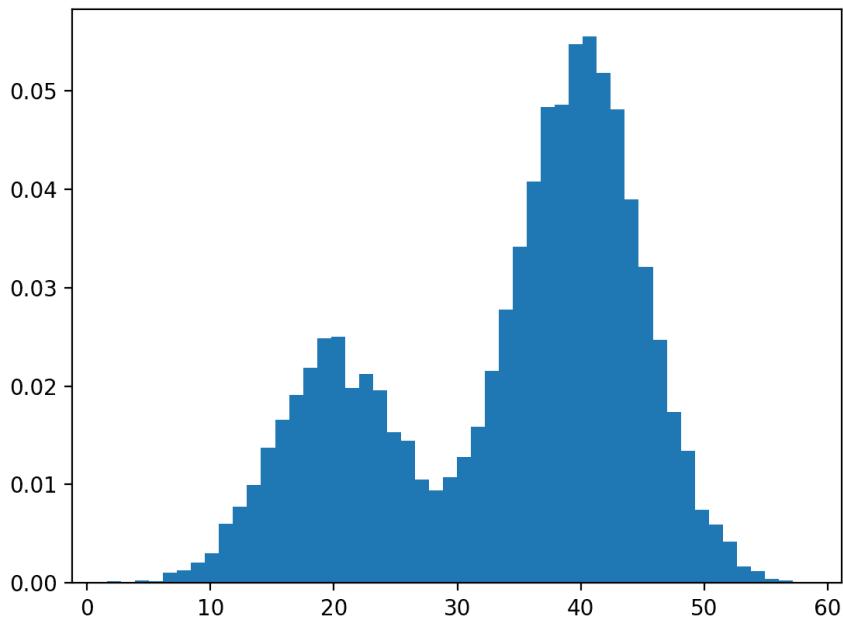


Рис. 2.3 – Гистограмма набора данных, построенная из двух разных гауссовых процессов

На ранних этапах тренировки используются простые акустические модели: на вход даются простые MFCC features (Рис. 2.4), фонемы рассматриваются вне контекстной зависимости, для моделирования вероятности эмиссии в HMM используется смесь гауссиан с диагональными матрицами ковариаций (Diagonal GMMs — Gaussian Mixture Models).

Результаты каждой предыдущей акустической модели являются стартовой точкой для тренировки более сложной модели, с более сложным входом, выходом или функцией распределения вероятности эмиссии. Существует множество способов улучшения акустической модели, однако наиболее значительный эффект имеет переход от GMM-модели к DNN (Deep Neural Network), что повышает качество распознавания практически в два раза. Нейронные сети лишены многих ограничений, характерных для гауссовых смесей, и обладают лучшей обобщающей способностью. Кроме того, акустические модели на нейронных сетях более устойчивы к шуму и обладают лучшим быстродействием.

Нейронная сеть для акустического моделирования тренируется в несколько этапов. Для

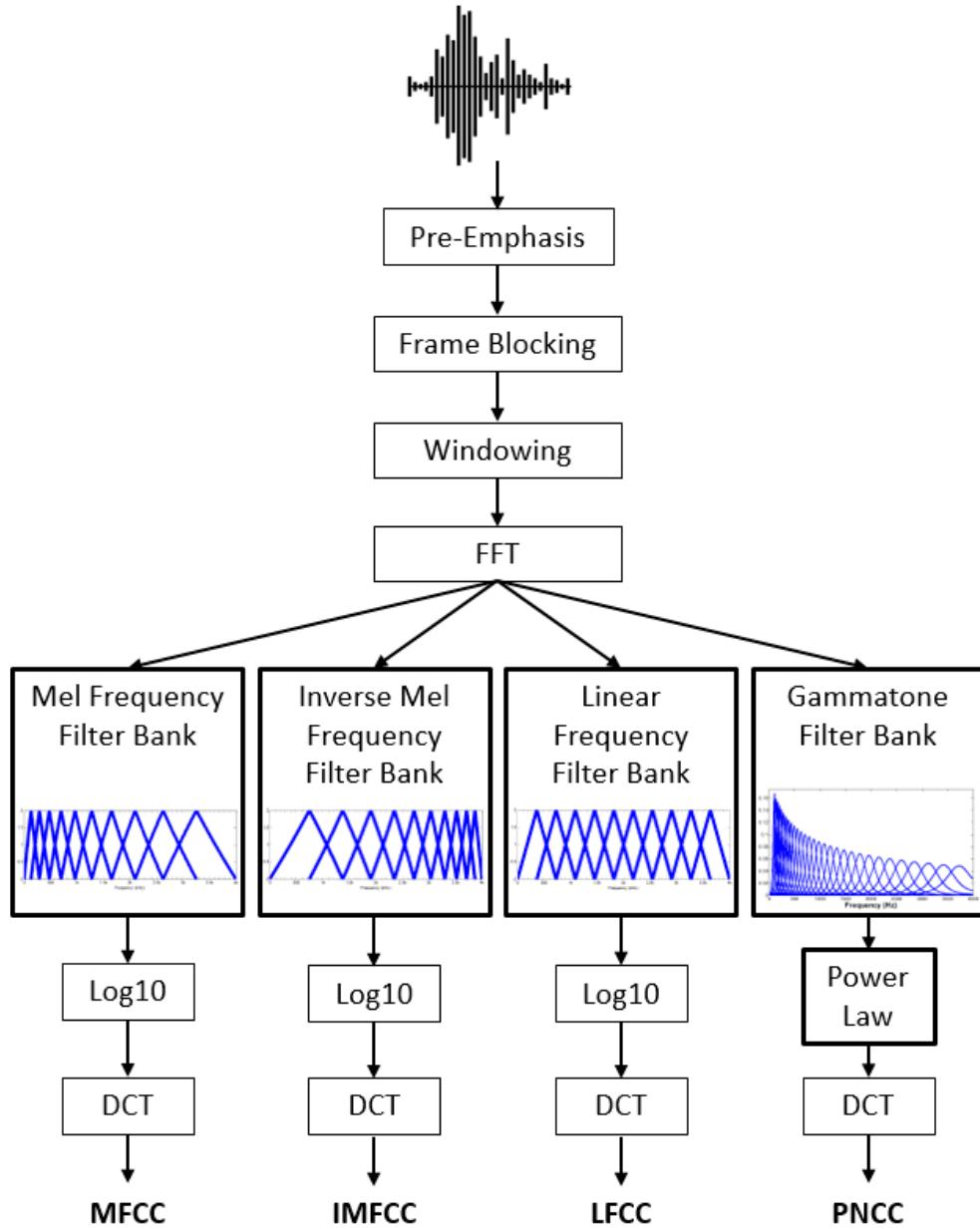


Рис. 2.4 – Четыре метода выделения признаков MFCC, IMFCC, LFCC и PNCC

инициализации нейросети используется стек из ограниченных машин Больцмана (Restricted Boltzmann Machines, RBM). RBM – это стохастическая нейросеть, которая тренируется без учителя. Хотя выученные ей веса нельзя напрямую использовать для различия между классами акустических событий, они детально отражают структуру речи. Можно относиться к RBM как к механизму извлечения признаков (feature extractor) – полученная генеративная модель оказывается отличной стартовой точкой для построения дискриминативной модели. Дискриминативная модель тренируется с использованием классического алгоритма обратного распространения ошибки, при этом применяется ряд технических приемов, улучшающих сходимость и предотвращающих переобучение (overfitting). В итоге на выходе нейросети – несколько фреймов MFCC-features (центральный фрейм подлежит классификации, остальные образуют

контекст), на выходе — около 4000 нейронов, соответствующих различным сенонам. Эта нейросеть используется как акустическая модель в production-системе.

2.4 Рекуррентные нейронные сети

LSTM — это класс возвратных нейронных сетей. Поэтому, прежде чем мы сможем перейти к LSTM, важно понять нейронные сети и рекуррентные нейронные сети [8].

Нейронные сети - Искусственная нейронная сеть представляет собой слоистую структуру из связанных нейронов, вдохновленную биологическими нейронными сетями. Это не один алгоритм, а комбинация различных алгоритмов, которая позволяет нам выполнять сложные операции с данными. Рекуррентные нейронные сети - это класс нейронных сетей, предназначенный для работы с временными данными. Нейроны RNN имеют состояние / память ячейки, и ввод обрабатывается в соответствии с этим внутренним состоянием, которое достигается с помощью петель в нейронной сети. В RNN существуют повторяющиеся модули «tanh» слоев, которые позволяют им сохранять информацию. Однако недолго, поэтому нам нужны модели LSTM.

LSTM - Это особый вид рекуррентной нейронной сети, способной изучать долгосрочные зависимости в данных. Это достигается за счет того, что повторяющийся модуль модели имеет комбинацию четырех слоев, взаимодействующих друг с другом [8].

На (Рис. 2.5) отображены структуры вышеупомянутые виды рекуррентных сетей:

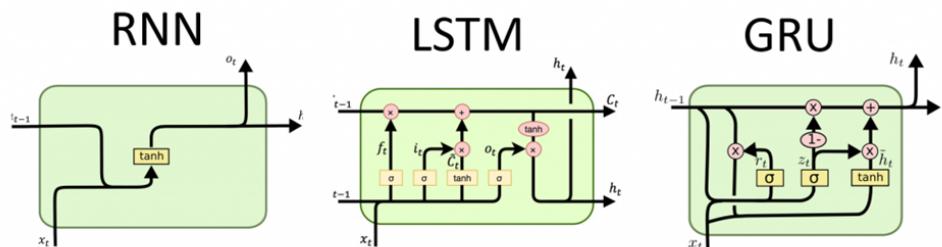


Рис. 2.5 – Структуры рекуррентных нейронных сетей

2.5 Выводы

Было изучено описание работы модели виртуального актора. Изучены подходы и методы применительно к задаче распознавания речи. Были рассмотрены ключевые подходы машинного обучения, которые будут использоваться в работе, а именно рекуррентные нейронные сети

3. Проектирование модели поведения виртуального агента

В этом разделе описывается и обосновывается выбор инструментария для проектирования и программного воплощения модели поведения актора в заданной парадигме. Описываются ключевые моменты проектирования и программной реализации модели поведения актора.

3.1 Проектирование модифицированного прототипа Виртуального Актора

В данном разделе описывается реализация когнитивной модели в виде функций, интерпретируемых в форме псевдокода. Здесь приведены основные функции, вызываемые при работе алгоритма по выбору действия для Виртуального Актера. В начале будет описан алгоритм по выбору ответного действия Виртуального Актера на действие человека.

На вход поступает действие человека, целью которого является влияние на Виртуального Актера. В зависимости от характеристики действия пересчитываются оценки Appraisals человека и Виртуального Актера по формуле 3. Соответственно уже на данном этапе меняется характеристика взаимодействующих акторов, что может повлиять на выбор ответного действия пингвина. Следующим шагом в соответствии с режимом работы моральной схемы производится вычисление новых значений векторов Feelings для взаимодействующих акторов.

Работа моральной схемы, следующая - по умолчанию она выключена и значения Feelings вычисляются по формуле 4. В этот момент времени моральная схема не оказывает никакого воздействия на принятие решения Виртуальным Актером. При отклонении абсолютного значения вектора Feelings от начального положения больше, чем на заранее заданную константу StartMoralSchema - моральная схема начинает свою работу. Ее функционирование происходит теперь в двух режимах:

1. Feelings вычисляется по формуле номер 5 в случае, если (расхождение между Feelings и Appraisals больше заранее определенной константы). Данный режим работы моральной схемы называется конфликтным. В таком ключе продолжается работа пока верна формула, описанная выше.
2. Feelings константна и экстремальна по своим параметрам. Это означает, что Виртуальный Актор понимает каким образом нужно относиться к человеку: как к другу или врагу, как к подчиненному или начальнику и т.д. В данном режиме схема работает до тех пор, пока расхождения между Feelings и Appraisals не станет критическим. В этом случае снова включается режим номер 1.

Возможен также вариант, когда отклонение Feelings от начального положения крайне мало. В таком случае моральная схема снова выключается и Feelings вычисляется по формуле 4.

Первая функция описывает работу моральной схемы. В данной функции определяется режим работы моральной схемы и каким образом будет меняться субъективная оценка человека в отношении Виртуального актора.

На данном этапе получается следующая картина - обработано действия человека, направленное на пингвина, и пересчитаны "оценки" и "чувства" взаимодействующих акторов. Далее на основе соматического критерия и когнитивного фактора вычисляются вероятности для всех возможных действий в данной ситуации. Набор всех действий для Виртуального Актора заранее описан в *.json файле и у каждого действия есть параметр, который определяет в каком контексте оно применимо. Соответственно после определения действий применимых в данной ситуации, для каждого из них вычисляется вероятность на основе когнитивного фактора по формуле 6, на основе соматического по формуле 7 и итоговая вероятность по формуле 8.

Следующим шагом необходимо определить ответное действия для пингвина из возможных с использованием ранее посчитанных для них вероятностей. Все действия сортируются по возрастанию численного значения их вероятностей. При помощи функции рандомной генерации чисел, основанной на равномерном распределении, генерируется дробное число от 0 до 1. Далее по формуле 9 определяется ответное действие.

Где - вероятность i -го действия, k - число, пробегающее от 1 до n , где n - общее количество рассматриваемых действия для Виртуального Агента. Минимальное k , при котором выполняется неравенство, описанное в формуле 9, соответствует номеру действия в списке действий, отсортированном по возрастанию вероятностей. Это действие и будет совершено Виртуальным Актором.

Предпоследним этапом работы алгоритма является пересчет значений Appraisals и Feelings по описанной в начале данного пункта схеме.

В конце действие, которое должен выполнить пингвин возвращается в виде строкового значение в функцию, отвечающую за перемещение и действия пингвина в виртуальном окружении и выбранное действие визуализируется Виртуальным Актором. После этого человек может снова совершить воздействие на пингвина, и вся данная процедура снова повторится.

Далее будет описан алгоритм выбора самостоятельного действия Виртуальным Актором. Под самостоятельным действием понимается такой действия, которое предпринимает пингвин, основываясь на состоянии виртуального окружения и отношении с человеком. Алгоритм выбора самостоятельного действия схож с алгоритмом выбора ответного действия на действие со стороны человека, направленное на Виртуального актора. Исполнение начинается с шага

вычисления вероятностей для возможных в данной ситуации действий на основе когнитивного и соматического фактора. В данной парадигме возможно проявление одного из следующих взаимодействий со стороны пингвина:

- Подойти к корзине со снежками с желанием поиграть
- Подойти к корзине с рыбой и попросить поесть
- Подойти к человеку с целью взаимодействия с ним
- Пойти спать
- Посмотреть в сторону человека
- Поприветствовать человека

Здесь может быть проблема с бесконечной очередью вызовов действий без временного интервала между ними. Тогда взаимодействие с Виртуальным агентом может стать проблематичным. Возможны два пути решения данной проблемы. Первое решение — это создание некоторого минимального промежутка времени после выполнения самостоятельного действия пингвина, в течение которого, программно будет запрещено совершение самостоятельного действия для Виртуального Актора. Такой подход в определенной мере решает проблему бесконечной непрерывной очереди действий пингвина, однако он довольно искусственный и слабо согласуется с моделью социально-эмоционального интеллекта на основе когнитивной архитектуры eBICA. Второй подход заключается в создании действий “заглушек”, которые практически не влияют на соматические и когнитивные оценки акторов. Данными действиями могут быть, например:

- Пингвин стоит на месте
- Пингвин перемещается по виртуальному окружению в какую-либо его точку

При совершении этих действий оценки Виртуального Актора и человека практически не будут меняться и будет создан временной интервал между социальными действиями пингвина.

3.2 Модель тембора

В связи с задачей эмоционального взаимодействия ?пользователя с виртуальным актором ставится задача распознавания эмоциональной составляющей из человеческой речи. А также классификация эмоционального воздействия.

Так как в данной работе для анализа человеческой речи используется нейросеть, а это значит, что требуется дата-сет для обучения данной нейросети. В большинстве работ по классификации эмоций присутствующих в речи человека используются открытые данные RAVDESS. Этот дата-сет включает в себя 7356 аудио записей с эмоциональным наполнением.

Область распознавания речи и ее эмоциональной составляющей довольно обширна. И решение данной задачи с нуля является слишком трудозатратной. Поэтому в работе будет использована предобученная модель. Данная модель позволит преобразовывать аудиоинформацию в векторное пространство.

Вектор полученный из записи человеческой речи будет использован для обучения нейронной сети, которая по признакам вектора определит эмоциональную составляющую.

Для того, чтобы получить решение данной задачи будут использованы такие модели (Рис. 3.1):

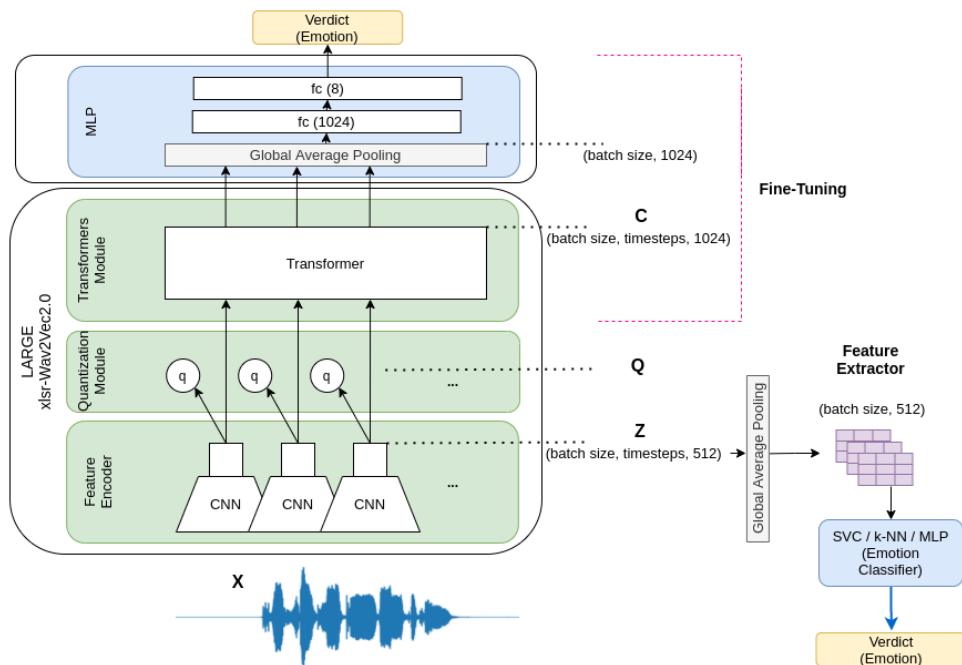


Рис. 3.1 – Предлагаемые конвейеры для распознавания речевых эмоций.

SVM с ядром «RBF» - является одной из наиболее популярных методологий обучения по precedентам, предложенной В. Н. Вапником и известной в англоязычной литературе под называнием SVM (Support Vector Machine).

k-ближайших соседей (kNN) с большинством голосов для выбора класса - расшифровывается как k Nearest Neighbor или k Ближайших Соседей – это один из самых простых алгоритмов классификации, также иногда используемый в задачах регрессии.

Многослойный персептрон (MLP) - это класс искусственных нейронных сетей прямого распространения, состоящих как минимум из трех слоёв: входного, скрытого и выходного. За исключением входных, все нейроны используют нелинейную функцию активации.

3.3 Модель семантики

Текст человеческой речи позволяет передать не только непосредственно смысл, который закладывает в нее говорящий, но и эмоциональную составляющую. Конечно текстовый формат представления не позволяет передать интонацию, но даже по тексту можно выявить эмоцию говорящего.

В данной работе требуется в дополнение к эмоциональной оценке по аудиозаписи получить эмоциональную оценку по тексту. Что позволит наиболее точно установить эмоциональное взаимодействие с виртуальным актором. Так же такой анализ позволит выявить случаи сарказма в речи ?пользователя.

Для упрощения будем классифицировать текста шкале: негативная оценка - нейтральная оценка - позитивная оценка. Сегодня уже существуют достаточно точные решения, определяющие эмоциональную составляющую текста. Но большинство таких решений являются нейросетями обученными на полноценных текстах редко включающих в себя диалоговую составляющую. Тогда как в работе требуется только анализ диалогов. Поэтому будет взята модель анализа тональности текста и дообучена на диалоговых данных. Данные будут собраны посредством фильтрации диалогов из публично доступных данных с сохранением разметки.

3.4 Инструменты для анализа текста

После того как речь была распознана и конвертирована в текстовый формат ставится задача определить семантический смысл предложения.

Возможность идентификации семантической близости между словами сделала модель word2vec широко используемой в NLP-задачах, которые подробно описываются в [13]. Идея word2vec основана на контекстной близости слов. Каждое слово может быть представлено в виде вектора, близкие координаты векторов могут быть интерпретированы как близкие по смыслу слова [14].

Таким образом, извлечение семантических отношений (отношение синонимии, родственные отношения и другие) может быть автоматизировано. Установление семантических отношений вручную считается трудоемкой и необъективной задачей, требующей большого количества времени и привлечения экспертов. Но среди ассоциативных слов, сформированных с использованием модели word2vec, встречаются слова, не представляющие никаких отношений с главным словом, для которого был представлен ассоциативный ряд [15].

В работе рассматриваются дополнительные критерии, которые могут быть применимы для решения данной проблемы. Наблюдения и проведенные эксперименты с общеизвестными характеристиками, такими как частота слов, позиция в ассоциативном ряду, могут быть использованы для улучшения результатов при работе с векторным представлением слов в части

определения семантических отношений для русского языка.

Представление слов в виде векторов позволяет применять математические операции. В большинстве примеров можно встретить вычитание векторов, когда результат вычисления $\text{vec}(\text{'Madrid'}) - \text{vec}(\text{'Spain'}) + \text{vec}(\text{'France'})$ будет ближе к $\text{vec}(\text{'Paris'})$, чем к другим векторам из распределения. Таким образом, разница векторов может быть использована для поиска семантических отношений между словами [16].

Word2vec не возвращает напрямую семантические отношения между словами. В ассоциативном ряду, который может быть возвращен в качестве близких слов к запрашиваемому (главному) слову, отражаются слова, которые часто употребляются рядом в контексте. Бессспорно, в ассоциативном ряду встречаются синонимы, антонимы, гипонимы, гиперонимы, холонимы, меронимы, ассоциации и другие типы, которые могут быть определены как семантические отношения.

Для реализации используются такие пакеты языка программирования python как: ufal.udpipe и wget. Как было сказано ранее для нахождения близости слов мы используем готовую модель, которую взяли с ресурса <https://rusvectores.org/static/models> с помощью пакета wget, принцип работы описан в [22]. Данный пакет позволяет скачивать данные с веб ресурсов посредством GET запроса.

Для того чтобы модель word2vec, описанная в [23] могла определить расстояние между словами, ей следует передать слово, ставится задача определить его часть речи, для того, чтобы автоматизировать такой процесс по определению части речи слова из текста, мы должны с помощью системы word2vec понять по слову какими категориями частей речи оно обладает, самые распространённые варианты – это существительные и глагол, далее после определения, мы передаем их на анализ дистанции. Принцип работы модели представлен на (Рис. 3.2):

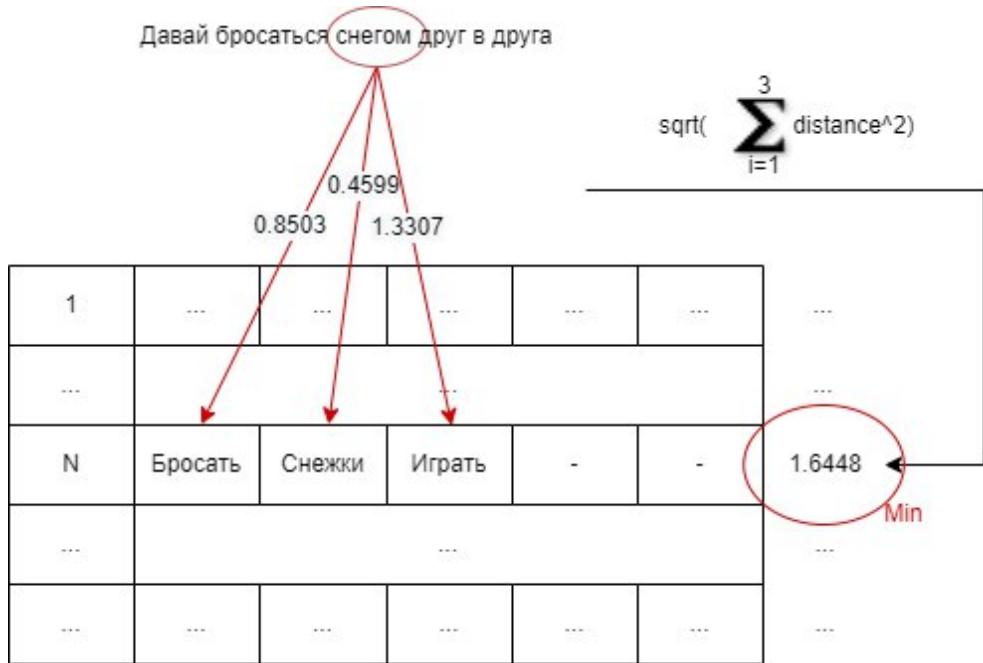


Рис. 3.2 – принцип работы word2vec

В рамках данной работы используется модель ruwikiruscorpora_tokens_elmo_1024_2019 - модель UDPipe для нахождения синонимов. Данная модель позволяет узнать семантическую близость слов.

Чтобы начать применять указанную выше нейронную сеть непосредственно, следует подготовить действия виртуальных акторов так, чтобы ими могла оперировать частично или полностью сама нейронная сеть. Для этого возьмем разобъем описание каждого действия на ключевые слова так, что каждое действие будет ассоциировано с набором слов. При этом создан класс отражающий действие актора как сущность - VAAction(Virtual Actor Action), принцип которой описан в [24].

Данный класс содержит в себе описание действия и ключевые слова, описывающие данное действие. Среди таких слов отсутствуют предлоги, а сами слова представлены в нормальной форме (для существительных - именительный падеж единственное число).

Для построения сценариев берутся тексты, полученные в разделе выше. Для каждого текста выделяются ключевые слова, которые наиболее близки к ключевым словам, описывающим действия акторов. Что происходит в процессе итерации через все слова текста так, что для каждого слова применяется значение близости слова к действию виртуального актора. Для работы с текстом создан класс WText (Web Text), который является ответственным за итерацию через все слова текста. С помощью методов класса задается функция, которая будет применяться к каждому слову. В качестве такой функции берется функция, которая находит близость слова с ключевыми словами экземпляра класса VAAction.

Также класс WText формирует набор определяющих текст слов, эти слова выбираются так, что значение семантической близости больше, чем заданный заранее порог, данной работе порог равен 0.08. После того как все тексты были переведены в экземпляры класса WText, была произведена оценка кол-ва текстов с одинаковым кол-вом определяющих слов.

3.5 Построение блок-схем и UML диаграмм

В данном разделе строятся блок-схемы реализованного алгоритма (Построение такой диаграммы рекомендуется в работе [17]). На (Рис. 3.3) представлена упрощенная блок-схема работы алгоритма по выбору ответного действия для актора. Данный алгоритм срабатывает каждый раз после совершения действия человеком, которое направлено на пингвина. Также данный алгоритм может приостанавливать любое продолжительное действие самого пингвина.

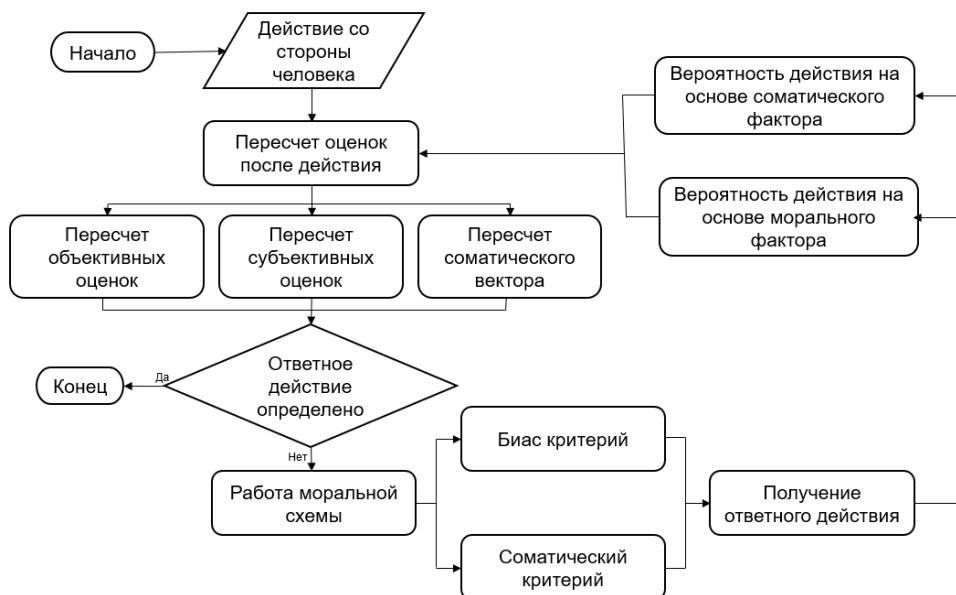


Рис. 3.3 – Блок-схема общей работы алгоритма

На вход поступает действие человека, направленное на пингвина. На основе оценок данного действия пересчитывается Appraisals и Feelings человека и пингвина, а также вектор соматического состояния. Затем последовательно срабатывают две функции: Критерий Биас и Соматический критерий. В них определяется вероятность для каждого действия на основе когнитивного и соматического состояния соответственно. Далее на основе данных вероятностей определяется итоговое действие, которое будет выполнено пингвином. Данное действие определяется наибольшей вероятностью.

рассмотрим более подробно алгоритм пересчета Feelings (субъективных оценок) для актеров. На (Рис. 3.4) представлена блок-схема данного алгоритма. На вход поступают Appraisals и Feelings актера. Затем, если до этого никогда данному актору не присваивалась константная

оценка Feelings, то Feelings присваивается значение согласно уравнению, когда моральная схема выключена и не оказывает никакого воздействия на вектор Feelings. Если Feelings актора уже принимал заданное константное экстремальное значение, то Feelings присваивается значение в зависимости от разницы норм Feelings и Appraisals.

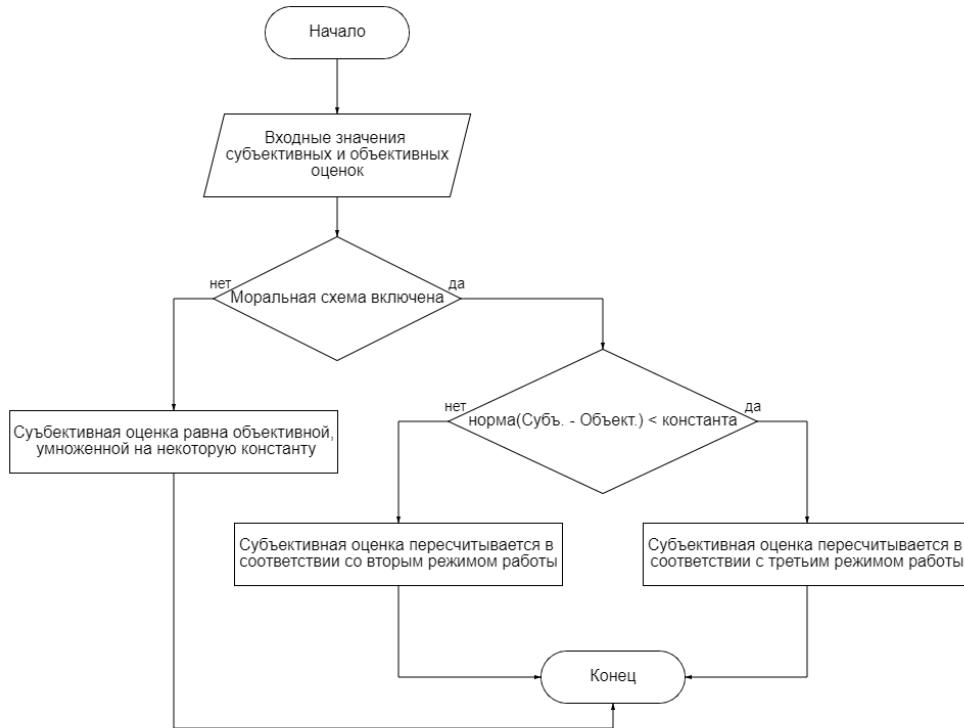


Рис. 3.4 – Блок-схема алгоритма пересчета значений Feelings акторов

На (Рис. 3.5) показана диаграмма классов, описывающая процедуру взаимодействия Виртуального Агента и человека. Диаграмма наглядно демонстрирует взаимосвязь человека и Виртуального Актора с методами логирования и выбора действий для пингвина с пересчетом “оценок” и “чувств”. При использовании данной схемы алгоритм поведения Виртуального Актора, основанного на когнитивной архитектуре eBICA, можно легко перенести и адаптировать под другую парадигму и виртуальное окружение.

Класс Human описывает характеристики человека, проводящего игровую сессию с пингвиным. Класс Penguin представляет собой Виртуального Агента, реализованного в виде пингвина в виртуальном окружении. Human и Penguin наследуются от общего класса Actor и имеют характеристики: Appraisals - “оценки”, Feelings - “чувства”, CoordinateX, CoordinateY, Azimuth - координаты местонахождения в виртуальном окружении и угол поворота относительно севера (север представляет собой сильно удаленную точку от площадки взаимодействия человека и пингвина и определен заранее). Также у классов Penguin и Human есть уникальные поля. У Penguin поле Somatic, которое представляет собой вектор соматического состояния, представ-

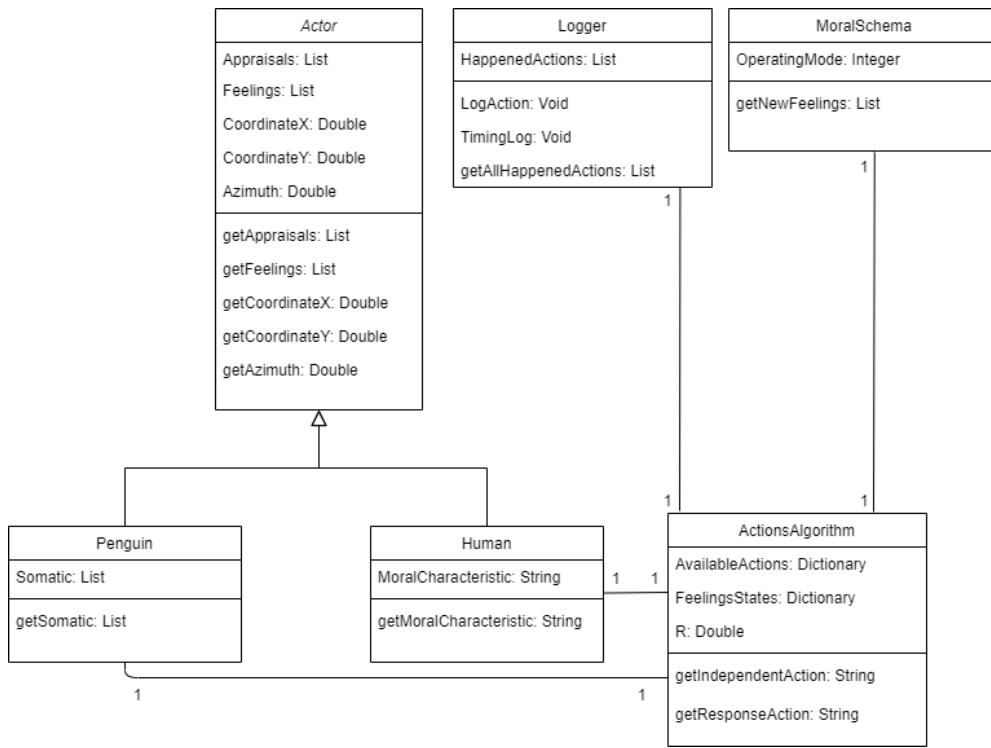


Рис. 3.5 – Диаграмма классов, описывающая взаимодействие Виртуального агента и человека

ленного структурой List. У Human поле MoralCharacteristic - представляет собой строковое значение, которое показывает, как пингвин относится к человеку с точки зрения моральной схемы. Penguin связан ассоциацией с классом ActionsAlgorithm. Данный класс отвечает за выбор действий для Виртуального Актора и обрабатывает действия человека. ActionsAlgorithm связан с классом MoralSchema, который представляет собой моральную схему. Принцип ее работы был описан в разделе 3.1. MoralSchema связана ассоциацией с Logger. Этот класс отвечает за логирование всех действий со стороны человека и пингвина. Также срабатывает функция TimingLog каждые 4 секунды для сохранения в файл текущего состояния виртуального окружения.

В дополнение к классам приложения до модернизации проектируются классы для работы с человеческой речью (Рис. 3.6):

Для работы с человеческой речью непосредственно используется класс Voice, который отвечает за фильтрацию речи и разбиение аудиодорожки на фреймы. Фреймы передаются в классы EmotionFromVoice и TextFromVoice. Первый извлекает эмоциональную составляющую из фрейма. Второй распознает речь и возвращает текст, который используется в классах EmotionFromText и CommandFromText. EmotionFromText как и EmotionFromVoice извлекает эмоциональную составляющую, но уже из текста. CommandFromText итерируется по полученному тексту и сопоставляет ключевые слова с действиями пингвина по алгоритму из рисунка 3.2 . Полученные команды передаются в ActionAlgorithm для того, чтобы увеличить вероятность выполнения

распознанных действий. Все эмоциональные оценки передаются в Penguin и меняют его состояние соответственно.

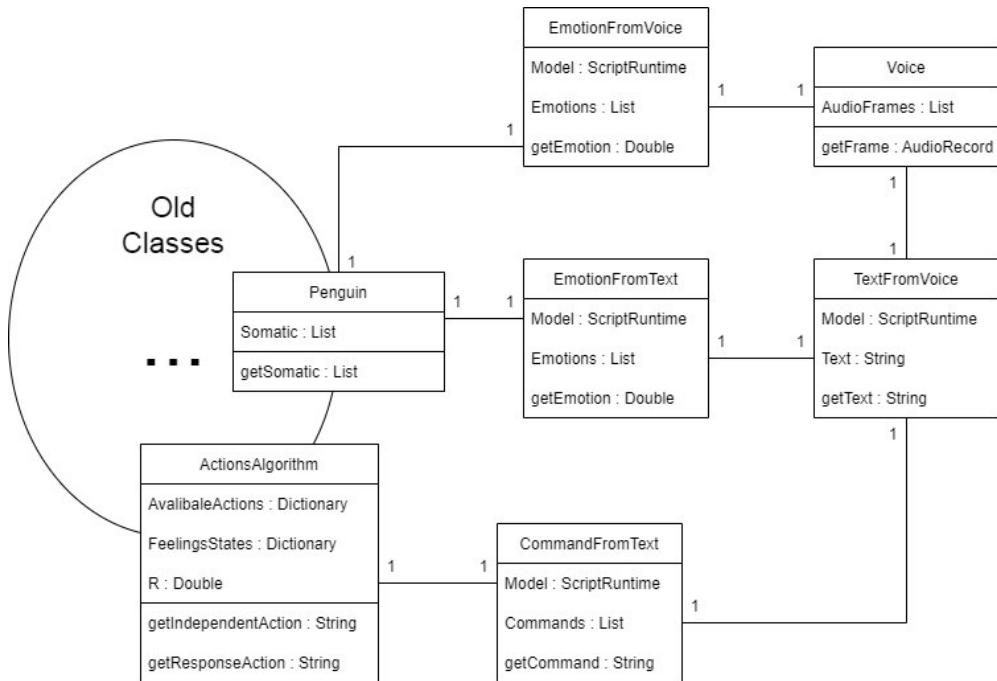


Рис. 3.6 – Расширенная диаграмма классов с участием голосового распознавания

3.6 Выводы

Производилось проектирование модифицированного прототипа виртуального актора с учетом всех ключевых особенностей парадигмы объектно ориентированного программирования. Были разработаны модели семантического, эмоционального анализа речи.

4. Реализация программного продукта

В данном разделе приводятся измененный алгоритм поведения Виртуального Актора.

4.1 Текущая версия реализованной модели

Для решения задачи разработки экспериментальной платформы на базе виртуального окружения для изучения социально – эмоционального поведения акторов целесообразно использовать специальные программные окружения для разработки виртуальных окружений и сопутствующих компонент – «игровые движки». По результатам анализа стало понятно, что такая среда должна прежде всего обладать:

1. Большим сообществом разработчиков
2. Подробной документацией
3. Относительно низкой сложностью
4. Простотой установки
5. Модульной системой

Всеми данными свойствами обладает лишь Unity3D, как самый распространённый, на данный момент инструмент построения виртуальных окружений. На базе Unity3D сделано больше игр, чем на любой другой технологии, поэтому он и обладает наиболее развитым сообществом разработчиков на данный момент. В сети имеется большое множество документации и курсов по данной технологии. Благодаря модульной системе, можно найти специальные программные модули, которые легко встраиваются в разработанный продукт, расширяя возможности всей системы. Преимуществом данной среды также следует считать относительную простоту переноса разработанного виртуального окружения на другие платформы (например, смартфоны, планшеты или же любую из существующих операционных систем). Также стоит отметить, что программная среда имеет бесплатную лицензионную версию для небольших команд разработчиков.

Unity предлагает разработчику возможность использовать три различных сценарных языка: C#, JavaScript (его модификация) и Boo (собственный диалект Python). Для разработки данной экспериментальной платформы был выбран C#, как де facto, стандарт при разработке на Unity3D.

В качестве среды разработки используется Visual Studio 2019. Для контроля версий использовать облачное хранилище Dropbox и GitLab. В качестве средства проектирования использовать

вался стандартный модуль Visual Studio для построения диаграмм классов.



Рис. 4.1 – Скриншот взаимодействия человека с Виртуальным Актором

В работе было реализовано программное приложение с помощью движка юнити результаты работы которого вы можете видеть на (рис. 4.1)

4.2 Интеграция моделей из python в C#

В расширенной диаграмме классов, было выделено какие стадии обработки проходит аудиозапись (Рис. 3.6).

Получение Audio из UNITY3D осуществлялось при помощи класса Microphone и далее помещалось в контейнер для аудио данных AudioClip, который хранит аудиофайл либо сжатым в ogg vorbis, либо без сжатия.

Далее, полученная аудио проходит обработку при помощи моделей реализованных на языке программирования python, при непосредственном использовании IronPython, который представляет из себя реализацию языка программирования Python с открытым исходным кодом, которая тесно интегрирована с .NET Framework. IronPython может использовать библиотеки .NET Framework и Python, а другие языки .NET могут также легко использовать код Python».

Существует несколько методов для работы со скриптами в Ironpython (Рис. 4.2):

Мы использовали метод ExecuteFile(), так как в нашем случае он самый подходящий.

В указанном выше методе происходит следующее:

- В коде C# в метод ExecuteFile(@"/home/...") помещен путь к файлу .py
- Функция из Python определяется в C#,
- Возвращается результат по завершению исполнения кода .py,



Рис. 4.2 – Методы для работы со скриптами в Ironpython

Таким образом проходит интеграция реализации динамического языка программирования IronPython в проект.

4.3 Реализация и дообучение модели

Как говорилось ранее для создания возможности эмоционального взаимодействия пользователя с виртуальным актором - берется предобученная модель "jonatasgrosman/wav2vec2-large-xlsr-53-russian". Данная модель работает в составе библиотеки huggingsound. На выходе модели для отдельного аудиофайла получаются 512-мерные вектора. Далее высчитывается средней вектор, который передается в модели для распознавания эмоций. Используется python3.8 и библиотека машинного обучения pytorch и sklearn.

Для SVC с ядром RBF менялся параметр регуляризации, а именно он принимал значения 1, 10 и 100.

Для k-NN менялось количество соседей от 10 до 40 с шагом в 10.

Для MLP выходной слой не менялся и состоял из 8 нейронов, тогда как количество внутренних слоев менялось от 1 до 3 так, что количество нейронов в них оставалось неизменным и равным 1024. В качестве функции активации была выбрана функция гиперболического тангенса.

В результате было получено (рис. ??):

При обучении датасет разбился на парти в 100 образцов в каждой так, чтобы была возможность обучать модели на GPU. Обучение проводилось до 10 эпох так, что результирующей мо-

деюсь становилась та, что показывала максимальный результат в какой-либо эпохе. Так как ставится задача классификации, то в качестве функции потерь используется функция потерь перекрестной энтропии. В качестве оптимизатора использовался оптимизатор Адам.

4.4 Выводы

Было реализовано приложение позволяющее взаимодействовать с виртуальным окружением, в частности с виртуальным актором. Были дообучены модели машинного обучения с целью построения эмоционального взаимодействия. Для чего в приложение были интегрированы различные технические средства.

Заключение

В ходе работы над НИР был разработан и протестирован с участием испытуемых прототип Виртуального Актора, обладающий социально-эмоциональным интеллектом и помещенный в виртуальное окружение, которое создано при помощи графического движка Unity3d, была реализована система для анализа речи с выявлением эмоций соответствующим речевым признакам, так же были в дополнении к вышеуказанной системе, были спроектированы и реализованы методы воздействия на виртуального агента, основывающиеся на семантической составляющей речевого контекста. Данная работа является актуальной поскольку на данный момент эта область находится на начальных этапах развития и активной интеграции в различные индустрии. Созданная и протестированная модель интеллекта затем может быть интегрирована в другие проекты с Виртуальным Актором: виртуальный слушатель, виртуальный клоун, виртуальный танцор.

1. Были проанализированы различные когнитивные архитектуры.
2. Проанализированы методы распознавания речи.
3. Доработан алгоритм поведения Виртуального Актора.
4. Унифицирован алгоритм поведения Виртуального Актора.
5. Был разработан, алгоритм выявления эмоций из человеческой речи.
6. Добавлено возможность эмоционального взаимодействия с виртуальным актором.
7. Было Реализован визуальный агент и сцена, используя межплатформенную среду разработки компьютерных игр Unity3d.

Список литературы

1. V. S. A. Comparative analysis of implemented cognitive architectures //Biologically Inspired Cognitive Architectures 2011. – IOS Press. – 2011.
2. V. S. A. Emotional biologically inspired cognitive architecture //Biologically Inspired Cognitive Architectures. – 2013.
3. P. I. E. Emotions and feelings. – 2007.
4. V. S. A. Socially emotional brain-inspired cognitive architecture framework for artificial intelligence //Cognitive Systems Research. – 2020.
5. Langley P. Laird J. E. R. S. Cognitive architectures: Research issues and challenges //Cognitive. – 2009.
6. V. S. A. Emotional biologically inspired cognitive architecture //Biologically Inspired Cognitive. – 2013.
7. Bai S. Kolter J.Z. K. V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. – 2018.
8. *ресурс*] W. [Рекуррентная нейронная сеть. – 2020.
9. H. Y. W. K. K. Y. M. S. Comparative study of CNN and RNN for natural language processing. 2017. arXiv preprint.
10. И.А. Батраева А.Д. Нарцев А. Л. ИСПОЛЬЗОВАНИЕ АНАЛИЗА СЕМАНТИЧЕСКОЙ БЛИЗОСТИ СЛОВ ПРИ РЕШЕНИИ ЗАДАЧИ ОПРЕДЕЛЕНИЯ ЖАНРОВОЙ ПРИНАДЛЕЖНОСТИ ТЕКСТОВ МЕТОДАМИ ГЛУБОКОГО ОБУЧЕНИЯ. – 2020.
11. Standard L. HTML Standard - whatwg. –.
12. П. Ф. Машинное обучение. Наука и искусство построения алгоритмов, который извлекают знания из данных. – 2015.
13. Y. K. Convolutional neural networks for sentence classification // Proc. of the 2014 Conf. on Empirical Methods in Natural Language Processing (EMNLP 2014). – 2014.
14. Селезнев К., Владимиров А. Лингвистика и обработка текстов // Открытые системы. – М., 2013.

15. *Mirkin B.* Core Concepts in Data Analysis: Summarization, Correlation and Visualisation, DOI. — 2011.
16. *Manning C., Schuetze H.* Foundations of Statistical Natural Processing. MIT. — M., 1999.
17. *Weisfeld. M.* The Object-Oriented Thought Process. — Fourth Edition. — Addison-Wesley Professional. — 2013.
18. *Константин Симаков И. К.* Особенности очистки адресных данных // Открытые системы. СУБД. — 2013.
19. *Ильинский Д., Черняк Е.* Системы автоматической обработки текстов // Открытые системы. СУБД. — М., 2014. — URL: <https://www.osp.ru/os/2014/01/13039687>.
20. *H. Y. W. K. K. Y. M. S.* Comparative study of CNN and RNN for natural language processing. 2017. arXiv preprint. — 2010.
21. *H. K. В.* Элементы теории ассоциативной семантики // Управление большими системами. — 2012.
22. *М. Л. Ю.* Люди и знаки. — 2010.
23. *MyStem. Я. технология.* MyStem. —. — URL: <https://tech.yandex.ru/mystem>.
24. *Xin. R.* Word2vec parameter learning explained. 2014. arXiv preprint arXiv: 1411.2738. — 2010.