Федеральное государственное бюджетное образовательное учреждение высшего образования «Комсомольский-на-Амуре государственный университет»

На правах рукописи

Алена Анатольевна Животова

МАТЕМАТИЧЕСКАЯ МОДЕЛЬ, АЛГОРИТМЫ И ПРОГРАММНЫЙ КОМПЛЕКС ДЛЯ ПОВЫШЕНИЯ КАЧЕСТВА МАШИННОГО ПЕРЕВОДА УЗКОСПЕЦИАЛЬНЫХ ТЕХНИЧЕСКИХ ТЕКСТОВ НА АНГЛИЙСКИЙ ЯЗЫК

Специальность 1.2.2. – Математическое моделирование, численные методы и комплексы программ

Диссертация на соискание учёной степени кандидата технических наук

Научный руководитель: кандидат технических наук, доцент, Бердоносов Виктор Дмитриевич

ОГЛАВЛЕНИЕ

ВВЕДЕНИЕ	. 5
ГЛАВА 1 ПРЕДРЕДАКТИРОВАНИЕ КАК ПЕРСПЕКТИВНОЕ НАПРАВЛЕНИ	Œ
РАЗВИТИЯ СИСТЕМ МАШИННОГО ПЕРЕВОДА	14
1.1 Развитие систем машинного перевода	14
1.1.1 Анализ на основе ТРИЗ-эволюционного подхода	15
1.1.2 ТРИЗ-эволюционный анализ систем МП	18
1.2 Интерактивный перевод и предредактирование	33
1.3 Методы предредактирования для повышения качества машинного перевода.	36
1.3.1 Правила для авторов по написанию исходных текстов	36
1.3.2 Концепция контролируемого языка	36
1.3.3 Предредактирование, основанное на правилах	37
1.3.4 Решение задачи перефразирования как способ предредактирования 3	38
ГЛАВА 2 МАТЕМАТИЧЕСКАЯ МОДЕЛЬ ОПТИМИЗАЦИОННОГО	
ПРЕДРЕДАКТИРОВАНИЯ	40
2.1 Обобщенная математическая модель процесса перевода	40
2.1.1 Термины и определения	41
2.1.2 Оценка задачи перевода	48
2.1.3 Предредактирование исходного текста	50
2.1.4 Перевод	52
2.1.5 Оценка качества перевода	53
2.1.6 Постредактирование переведенного текста	56
2.1.7 Обучение переводчика	57
2.1.8 Обобщенная модель перевода	59
2.2 Математическая модель поиска весов значимости признаков текста	50
2.3 Математическая постановка задачи машинного перевода	51
2.4 Математическая постановка задачи оптимизационного предредактирования.	54
2.5 Метод градиентного спуска для решения задачи автоматического	
оптимизационного предредактирования	67

ГЛАВА 3	методика реализации оптимизационного	
	ПРЕДРЕДАКТИРОВАНИЯ	69
3.1 Модел	ть оптимизационного предредактирования	69
3.2 Обуче	ение модели оптимизационного предредактирования	. 70
3.3 Обуче	ение модели оценки сложности задачи перевода	. 72
ГЛАВА 4	ПРОГРАММНЫЙ КОМПЛЕКС ОПТИМИЗАЦИОННОГО	
	ПРЕДРЕДАКТИРОВАНИЯ УЗКОСПЕЦИАЛЬНЫХ	
	РУССКОЯЗЫЧНЫХ ТЕКСТОВ ДЛЯ ИХ ПЕРЕВОДА НА	
	АНГЛИЙСКИЙ ЯЗЫК	. 76
4.1 Архит	гектура программного комплекса и его подсистем	. 76
4.1.1 N	Модуль автоматической очистки сырых данных из памятей переводов	
	САТ для тренировки языковой модели	. 78
4.1.2 N	Модули машинного перевода	. 78
4.1.3 N	Модуль оценки качества машинного перевода	. 79
4.1.4 N	Модуль препроцессинга текстовых данных для взвешенной оценки	
	параметров русскоязычного текста	. 79
4.1.5 N	Модуль вероятностной оценки сложности переводческой задачи для	
	систем машинного перевода	. 82
4.1.6 I	Предредактор русскоязычных узкоспециальных текстов для систем	
	машинного перевода	. 86
4.2 Данны	ые для обучения и тестирования программного комплекса	. 87
4.2.1 H	база данных структурного анализа предложений технических	
	русскоязычных текстов	. 87
4.2.2 I	Корпус параллельных двуязычных текстов нефтегазовой тематики для	
	тренировки языковых моделей в задачах перефразирования	
	узкоспециальных технических русскоязычных текстов и повышения	
	качества их перевода на английский язык	. 89
4.3 Тести	рование программного комплекса	90
4.3.1 I	Тостановка задачи тестирования	90
4.3.2 I	Начальные условия и границы проведения тестирования	. 90

•	4.3.3 Методология и план тестирования	91
	4.3.4 Результаты тестирования	92
4.4	Внедрение программного комплекса в контур автоматизации процессов	
	переводческой деятельности	96
СПІ	ИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ1	03
ПРІ	ИЛОЖЕНИЕ А Охранные документы на результаты интеллектуальной	
	деятельности1	115
ПРІ	ИЛОЖЕНИЕ Б Акт о внедрении (использовании) результатов кандидатской	
	диссертационной работы1	25

ВВЕДЕНИЕ

Перевод – рутинная необходимость во многих отраслях, включая науку, производство, медицину и т.д., и с ростом количества информации и скорости ее генерирования растет и потребность в повышении качества перевода наряду с сокращением затрат на него. Современные системы машинного перевода (МП) показывают высокие показатели качества, кардинально изменив к 2023 г. структуру рынка лингвистических услуг, вытесняя переводчиков в пользу пост-редакторов и корректоров МП. Интерес исследователей к теме машинного перевода также значительно возрос. Так, согласно данным базы Science Direct количество статей по ключевым словам «машинный перевод» (Machine Translation) и «качество машинного перевода» (Масhine Translation Quality) в 2023 году выросло на 117% и 146% соответственно по сравнению с 2017 годом [1].

В переводе специфика предметной области имеет ключевое значение, ведь МП тем эффективнее, чем больше обучающих данных (корпусов) загружено в систему, однако для некоторых предметных областей собрать достаточный объем параллельных тренировочных данных не всегда возможно. Так, например, нефтегазовый сектор — один из ключевых для экономики нашей страны с большой долей участия иностранных компаний в проектах освоения месторождений и нефтегазопереработки. Качество перевода в данной области имеет критически важное значение для коммуникации и обмена технологиями. Для подобных предметных узкоспециальных областей, связанных с объектами критической инфраструктурой, собрать достаточный объем двуязычных тренировочных корпусов проблематично ввиду ограничений конфиденциальности данных и секретности разработок.

Несмотря на выдающиеся прорывы нейросетевых, гибридных и больших языковых моделей МП в области семантической точности и гладкости перевода, вопрос качества перевода системами МП нельзя назвать решенным. Результат работы МП — черновик, который пользователь должен оценить и доработать самостоятельно. При этом пользователь без знания языка перевода не имеет инструментов для того, чтобы влиять на результат или хотя бы оценить качество полученного перевода. Эту

проблему активно освещают зарубежные исследователи А. Lear, С. Quinci, С. Canfora, А. Ottman, D. Kenny, Р. Sanchez-Gijon. Предоставляя пользователю средства обработки текста на языке, носителем которого он является, на любом из этапов перевода, можно повысить его качество. Зная ключевые параметры текста и их связь с предполагаемой оценкой качества, становится возможным предложить алгоритмы и инструменты автоматического и/или полуавтоматического редактирования текста с целью его оптимизации для повышения качества перевода на требуемый язык.

Значительный вклад в разработку теоретических и практических основ в области подготовки исходных текстов к переводу, предварительного редактирования и упрощения естественных языков для систем автоматической обработки текстов, в частности систем МП, внесли зарубежные авторы: V. Kumar, F. Azadi, M. Federico, V. Alabau – в области интерактивного перевода; V. Sereton, P. Bouillon, J. Gerlach, A. Taufik, Y. Liang, W. Han, A. G. Arenas, C. Shei, Y. Hiraoka, M. Yamada, R. Miyata, A. Fujita – в области разработки подходов к предредактированию; L. O'Brien, D. Folaron, W. Aziz, M. Toledo – в области контролируемых и упрощенных языков. Среди российских авторов и для русского языка данная тема освещена незначительно, однако известны работы А.В. Ниценко, И. В. Оборневой, А.Д. Дмитриевой, А. Н. Лапошиной и др. в области оценки восприятия текста и упрощения русскоязычных текстов в соответствии с квалификацией реципиента.

Теоретическая актуальность и значимость темы определяется недостаточным уровнем исследований, касающихся алгоритмов предредактирования русскоязычных текстов для повышения качества их машинного перевода на другие языки, в частности на английский язык. Современные системы не анализируют исходный текст с целью оценки сложности задачи перевода и оптимизации результата МП, который должен быть проверен и при необходимости доработан пользователем, не всегда обладающим достаточной для этого компетенцией.

Практическая актуальность и значимость темы исследования объясняется тем, что в условиях развития экономики контента, когда цикл генерирования и обновления информации сократился с месяцев до дней, и с учетом необходимости ее локализации в режиме реального времени, требуется оптимизация временных и

материальных затрат на непрерывный перевод больших массивов текстовых данных с сохранением качества перевода, особенно в узкоспециальных технических областях, для который в открытых источниках не достаточно тренировочных данных. Использование вероятностной оценки сложности задачи перевода и алгоритмов оптимизационного предредактирования исходных текстов позволяет снизить зависимость качества МП от человеческого фактора, а также предоставляет необходимые критерии для разработки стратегии управления рисками, связанными с компетенцией исполнителей и пользователей систем МП, при реализации крупных переводческих проектов.

Основная идея диссертации в том, чтобы, используя особенности работы алгоритмов систем МП и основы теории перевода, автоматизировать предварительное редактирование исходных текстов с тем, чтобы оптимизировать их структуру, благодаря чему системы МП будут эффективнее переводить их на требуемый язык и допускать меньше стилистических ошибок, для распознания которых требуется более высокая компетенция пользователя в области языка перевода.

Объектом исследования выступает процесс перевода текстов, **предметом** исследования — методы повышения качества перевода при работе с исходным текстом.

Целью работы является разработка моделей и алгоритмов и их реализация для повышения качества машинного перевода узкоспециальных технических текстов путем автоматического оптимизационного предредактирования.

Задачи исследования:

- Выполнить анализ существующих систем машинного перевода, направлений их совершенствования и способов реализации автоматического оптимизационного предредактирования.
 - Разработать математическую модель процесса перевода.
- Разработать методику и алгоритм вероятностной оценки сложности задачи перевода.
- Разработать методику оптимизационного предредактирования исходных текстов.

- Реализовать разработанные алгоритмы в программном комплексе для оценки сложности задачи перевода русскоязычных текстов на английский язык и оптимизационного предредактирования с целью повышения качества перевода на английский язык.
- Проверить адекватность разработанных алгоритмов на корпусе узкоспециальных технических текстов.

Научная новизна:

- 1. Предложена новая методика для повышения качества машинного перевода текстов с русского языка на английский язык, отличающаяся от существующих применением обратного перевода для сбора тренировочных данных и оптимизационного предредактирования на основе вероятностной оценки сложности задачи перевода.
- 2. Впервые предложена методика оценки сложности переводческой задачи для переводчика на основе его компетенции и специализации и параметров исходного текста, которая позволяет прогнозировать риски некачественного и/или несвоевременного решения задачи перевода.
- 3. Предложен новый алгоритм, позволяющий расширить область применения оптимизационного метода градиентного спуска путём использования элементов нечеткой логики в выражении функции правдоподобия через функцию принадлежности полученного текста низкой сложности задачи перевода для выбранной системы МП.
- 4. Предложен новый алгоритм, позволяющий расширить область применения метода наименьших квадратов для поиска весов значимости параметров исходного текста для вероятностной оценки ожидаемого качества его перевода на целевой язык.
- 5. Предложена новая архитектура и реализован программный комплекс для повышения качества машинного перевода текстов с русского языка на английский язык, отличающийся от существующих применением ансамбля моделей для оптимизационного предредактирования на основе вероятностной оценки сложности задачи перевода с целью повышения качества машинного перевода текстов с русского языка на английский язык.

Основные положения, выносимые на защиту:

- 1. Математическая модель процесса перевода, позволяющая определить понятия компетенции и специализации переводчика через отношения множеств.
- 2. Математическая постановка задач машинного перевода и оптимизационного редактирования при помощи функции правдоподобия и элементов нечеткой логики.
- 3. Решение задачи максимизации правдоподобия оптимизационного предредактирования численным методом градиентного спуска (подъема).
- 4. Методика и алгоритм вероятностной оценки сложности задачи перевода на основе регрессионного анализа зависимости ожидаемого качества машинного перевода от признаков русскоязычного текста с оптимизацией функции потерь методом наименьших квадратов.
- 5. Методика и алгоритмы оптимизационного предредактирования русскоязычных текстов, позволяющие сократить объем необходимых данных для тренировки модели путем применения концепции обратного перевода и время тренировки модели за счет предварительной обработки тренировочных данных.
- 6. Программный комплекс для реализации системы анализа и предредактирования русскоязычных текстов для повышения качества машинного перевода на английский язык.

Практическая значимость работы обусловлена возможностью интегрирования программного обеспечения, реализующего вышеперечисленные алгоритмы, основанные на вероятностной оценке сложности задачи перевода и алгоритмах оптимизационного редактирования, в системы управления и автоматизации переводческой деятельности. Разработанные алгоритмы позволят повысить качество перевода русскоязычных узкоспециальных технических текстов в условиях ограниченного объема эталонных двуязычных корпусов для обучения нейросетевых моделей, оптимизировать затраты на перевод, повысить надежность существующих систем, снизить зависимость качества перевода от человеческого фактора.

Разработанные в ходе исследования алгоритмы и программные комплексы реализованы и внедрены в практическую деятельность ведущего предприятия

лингвистической отрасли в г. Комсомольске-на-Амуре — ООО «Агентство переводов «ФИАС-Амур» (акт внедрения результатов диссертации на соискание ученой степени кандидата технических наук № 6/23/1 от 10.06.2023); получены свидетельства о регистрации программ ЭВМ и БД для 6 программных модулей, 2 программных комплекса, 2 баз данных.

Достоверность результатов исследования определяется применением апробированных математических методов, включая теорию множеств, численных методов оптимизации, таких как метод наименьших квадратов и метод градиентного спуска, статистических методов, а именно метода максимального правдоподобия, а также использованием современных комплексов программ анализа данных и экспериментально.

Личный вклад автора. Все результаты, представленные в работе, получены автором самостоятельно. Из совместных работ в работу включены только результаты, полученные лично автором. Соавторы публикаций по теме диссертации участвовали в обсуждении постановочной части решаемых задач и результатов, полученных по разработанным автором методам и алгоритмам.

Соответствие паспорту специальности. Диссертационная работа со-ответствует области исследования специальности 1.2.2 «Математическое моделирование, численные методы и комплексы программ» по п. 2 «Разработка, обоснование и тестирование эффективных вычислительных методов с применением современных компьютерных технологий» (п. 1,2 научной новизны), п. 3 «Реализация эффективных численных методов и алгоритмов в виде комплексов проблемно-ориентированных программ для проведения вычислительного эксперимента» (п. 5 научной новизны), п. 4 «Разработка новых математических методов и алгоритмов интерпретации натурного эксперимента на основе его математической модели» (п. 3,4 научной новизны), п. 8 «Комплексные исследования научных и технических проблем с применением современной технологии математического моделирования и вычислительного эксперимента» (п. 1-4 научной новизны).

Апробация результатов исследования. Основные результаты работы докладывались и обсуждались на следующих научных конференциях:

- краевой конкурс молодых ученых Хабаровского края «XXVI краевой конкурс молодых ученых в сфере научных исследований», I место в секции «Физикоматематические науки и информационные технологии» (г. Хабаровск, 2024 г.);
- 29-ая международная конференция по компьютерной лингвистике и интеллектуальным технологиям «Диалог» (г. Москва, 2023 г.);
- международная научно-практическая конференция «Информационные технологии и интеллектуальные системы принятия решений» (ITIDMS 2022) (г. Москва, 2023 г.);
- VII международная научно-практическая конференция «Информационные технологии и высокопроизводительные вычисления» (ITHPC-2023), диплом III степени за лучший оклад среди молодых ученых (г. Хабаровск, 2023 г.);
- краевой конкурс молодых ученых Хабаровского края «XXV краевой конкурс молодых ученых в сфере научных исследований» (г. Хабаровск, 2023 г.);
- VI всероссийская национальная научная конференция молодых учёных «Молодёжь и наука: актуальные проблемы фундаментальных и прикладных исследований» (г. Комсомольск-на-Амуре, 2023 г.);
- международный конкурс «2022 International Activity of Innovation Entrepreneurship Creation», организован представительством МА ТРИЗ в Китае при поддержке Китайской ассоциации по науке и технологиям (Китай, 2022 г.);
- II международная научно-практическая конференция молодых ученых «Наука, инновации и технологии: от идей к внедрению» (г. Комсомольск-на-Амуре, 2022 г.);
- международная мультидисциплинарная конференция по промышленному инжинирингу и современным технологиям «Far East Con-2020» Дальневосточного федерального университета (г. Владивосток, 2020 г.).

Публикации. Основные теоретические и практические результаты диссертационного исследования опубликованы в 12 научных работах, в том числе в 3 работах в издании, рекомендованном ВАК, в 3 работах в изданиях, индексируемых в международной базе Scopus.

Объём и структура работы. Диссертация включает в себя введение, четыре основные главы, заключение, список используемой литературы и 2 приложения, изложена на 125 страницах. Текст работы содержит 13 таблиц и 22 рисунка и список литературы из 124 наименований. В приложениях содержатся копии свидетельств о государственной регистрации программ для ЭВМ и баз данных, копия акта внед-рения результатов диссертации.

В первой главе предложено использовать ТРИЗ-эволюционный подход к выявлению направлений развития гибридных систем нейронного МП, сформулированы положения ТРИЗ-эволюционного анализа, позволяющие систематизировать исследуемую область знаний с высокой степенью детализации. Применение ТРИЗэволюционного анализа систем МП позволило: систематизировать данные об эволюции систем МП; определить ключевые этапы развития систем МП; выделить главные производственные параметры, определяющие направления развития систем нейронного МП. Показано, что перспективным направлением является развитие методов и алгоритмов автоматизированной предобработки исходных текстов для нейронного МП, при этом исследование в этой области позволит добиться повышения качества МП. ТРИЗ-эволюционная карта позволила определить проблемы повышения качества МП, которые могут быть решены путем оптимизационного предредактирования исходных текстов. Кроме того, в главе рассмотрены основные способы предобработки исходных текстов для систем МП, таких как использование контролируемого языка, правил предредактирования и решение задачи перефразирования текста в контексте решаемой задачи.

Во второй главе предложено решение задачи оптимизационного предредактирования методом градиентного спуска (подъема); приводится обобщенная модель процесса перевода, разработанная на основе теории множеств, описывающая систему понятий прикладной лингвистики в области машинного перевода, включая формализацию понятий опыта, специализации и компетенции переводчика; описаны математическая постановка задач перевода и оптимизационного предредактирования через функцию правдоподобия и с использованием элементов нечеткой логики. В ходе моделирования впервые обоснована целесообразность и разработана

методология оценки сложности переводческой задачи. Результаты выполненного моделирования показывают, что уже на этапе оценки исходного текста, возможно предсказать ожидаемое качество перевода на основе параметров исходного текста и компетенции и специализации переводчика.

В третьей главе рассматриваются методика осуществления автоматического оптимизационного предредактирования текста с целью повышения качества машинного перевода относительно формализованных требований и методика расчёта сложности задачи перевода заданного текста для заданного переводчика в соответствии с формализованными требованиями к переводу.

В четвёртой главе описана реализация методов автоматического оптимизационного предредактирования узкоспециального технического русскоязычного текста с целью повышения качества его МП на английский язык и оценки сложности задачи перевода для системы МП; описывается программный комплекс, реализованный в соответствии с этими методиками. Описывается архитектура программного комплекса и основных его подсистем: генератора МП текстов с русского языка на английский язык; препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста; вероятностной оценки сложности переводческой задачи для систем МП; автоматической очистки сырых данных из памятей переводов САТ для тренировки языковой модели; тренировочного модуля языковой модели для перефразирования русскоязычных технических текстов, предредактора русскоязычных узкоспециальных текстов для систем МП. Описаны полученные в ходе реализации описанных методов и алгоритмов массивы данных: база данных показателей структурного анализа предложений технических русскоязычных текстов, корпус параллельных двуязычных текстов нефтегазовой тематики для тренировки языковых моделей в задачах перефразирования узкоспециальных технических русскоязычных текстов и повышения качества их перевода на английский язык. Приводятся результаты тестирования программного корпуса на реальных данных. Описаны возможности интеграции программного комплекса в контур автоматизации процессов переводческой деятельности.

ГЛАВА 1 ПРЕДРЕДАКТИРОВАНИЕ КАК ПЕРСПЕКТИВНОЕ НАПРАВЛЕНИЕ РАЗВИТИЯ СИСТЕМ МАШИННОГО ПЕРЕВОДА

1.1 Развитие систем машинного перевода

Одним из способов решения задачи повышении качества перевода наряду с сокращением затрат на него является автоматизация процессов перевода. Первые попытки такой автоматизации берут свое начало еще в 1950х годах. С тех пор технологии МП прошли большой путь, но только в 2014-2017 гг. [2, 3, 4] произошел качественный скачок роста идеальности данных систем, который привел рынок лингвистических услуг к пониманию того, что внедрение и развитие данных технологий является одним из наиболее перспективных направлений развития отрасли [5, 6]. Современные системы МП показывают высокие показатели качества, кардинально изменив к 2023 г. структуру рынка лингвистических услуг России [7], вытесняя переводчиков в пользу пост-редакторов и корректоров машинного перевода.

В литературе уже описаны многие аспекты истории развития систем и технологий МП. Автором одних из наиболее подробных обзоров является John Hutchins [8, 9, 10], в которых подробно описаны ключевые этапы развития систем на момент издания, их достоинства и недостатки и т.д. Многие исследователи описывают историю развития МП в виде хронологического обзора с целью дать представление об истории развития предметной области [11, 12, 13] описывают в большей степени социальный аспект развития технологий МП. Также в литературе приведены различные варианты классификации технологий МП по используемым методам, по структуре используемых для обучения данных и тп. Так, например, Lane Schwartz [14] очень подробно описывает историю развития МП, основные парадигмы и классификации технологий со ссылками на авторов и разработчиков, а также приводит обзор литературы по данной теме. Из всех работ по классификации технологий МП и категоризации статей по данной теме отдельно следует отметить архив статей по МП [15] — картотека материалов по МП за период с 1950 по 2017 год. В ней не только

собрана ключевая информация по МП, но и проведена подробнейшая категоризация по технологиям и разделам данной предметной области.

Несмотря на широкую освещенность темы в литературе, многие из существующих статей, опубликованных в рецензируемых журналах, устарели. Технологии за последние годы совершили качественный скачек, но рецензируемых материалов, описывающих новые технологии в разрезе эволюции очень мало, а описание перспектив и направлений развития ограничено и/или отсутствует их обоснование. Систематизация информацию о технологиях МП необходима выявления ключевых направлений и перспектив развития и исследования. Предлагается в анализ и систематизацию включить: краткое изложение ключевых аспектов технологии МП; оценку развития главных параметров при переходе от одной технологии к другой; определение ключевых проблем, ограничивающих применение в каждой технологии; компактное, но емкое описание решений, за счет которых произошло преодоление выявленных ограничений.

Такой анализ позволит получить представление о развитии систем МП, составить карту проблем, требующих решения, что сократит время на выявление актуальных проблем МП за счёт систематизации, визуализации и структурирования ключевых данных, а определение релевантных путей совершенствования систем МП облегчит и ускорит постановку задач исследований и формулирования гипотез.

1.1.1 Анализ на основе ТРИЗ-эволюционного подхода

Развернувшаяся в последние десятилетия информационная революция обострила и потребности создания и развития новых методов извлечения и систематизации знаний [16], которые обеспечивали бы возможность изучения и структурирования огромного объёма информации за ограниченное время. Наметить пути разрешения данного противоречия позволяет ТРИЗ-эволюционный подход. ТРИЗ (теория решения изобретательских задач) — область знаний, исследующая механизмы развития искусственных систем с целью создания практических методов решения инновационных задач [17]. Благодаря своей высокой эффективности и универсальности, ТРИЗ получила международное признание, успешно применяется и развивается во

многих областях деятельности человека, прежде всего, в промышленном производстве, науке и образовании [18, 19].

Согласно ТРИЗ, в развитии искусственных систем происходит чередование этапов количественного роста и качественных скачков. В процессе количественного роста в результате неравномерного развития характеристик искусственной системы проявляются противоречия, которые препятствуют повышению идеальности системы. Идеальность, в ТРИЗовском понимании, это отношение суммы параметров, характеризующих пользу к сумме параметров, характеризующих затраты. Как правило, в качестве оценки идеальности выступает главный производственный параметр. Противоречие — это проявление несоответствия между требованиями, предъявляемыми к системе, и ограничениями, налагаемыми на нее. Выявление и анализ противоречий лежат в основе прогнозирования развития систем. Разрешение противоречий неизменно ведет к повышению идеальности системы.

ТРИЗ-эволюционный подход [20] предполагает выстраивать эволюцию исследуемых систем по мере увеличения их идеальности. Развитие системы происходит при разрешении противоречий, за счёт использования инструментов ТРИЗ [21] (приёмы разрешения противоречий, вепольный анализ, законы и тренды развития систем и так далее). Таким образом, выстраивается эволюция систем от противоречия к противоречию. Такое выстраивание, позволяет не только систематизировать знания по соответствующим системам, но и предлагать новые высокоэффективные решения. Анализ при помощи ТРИЗ-эволюционного подхода позволяет систематизировать знания по соответствующим системам; определить главные производственные параметры, выявить закономерности и перспективы развития с учетом повышения идеальности; предлагать новые высокоэффективные решения путём разрешения выявленных противоречий. Данный подход был многократно описан в литературе на примере таких предметных областей как объектно-ориентированное программирование [22], автоматизированные системы управления [23], установки коксования [24] и др.

Согласно ТРИЗ-эволюционному подходу [25] исследование развития некоторой системы включает три этапа: анализ развития системы (I), построение ТРИЗ-

эволюционной карты (II) и анализ ТРИЗ-эволюционной карты (III). Этап I включает следующие шаги:

Шаг 1. Описание исходного объекта ТРИЗ-эволюции.

Шаг 2. Выявление противоречий у выбранного объекта. Противоречия записываются по формуле:

$$Con_i:MP_i,MP_k,$$
 (1)

где MP_j – улучшаемый параметр, MP_k – ухудшаемый параметр.

Шаг 3. Определение инструментов ТРИЗ, позволяющих разрешить выявленные противоречия.

Шаг 4. Описание итераций ТРИЗ-эволюции, т.е. переходов к последующим объектам, в которых разрешены отдельные противоречия. Описание содержит краткое текстовое описание разрешение противоречия(ий) и символьное описание итерации по формуле:

$$\{Con_{i},...,Con_{n}\} \xrightarrow{\left(\{IP_{j},...,IP_{m}\}\right)} Sol_{x}:\{MP_{k},...,MP_{v}\},\tag{2}$$

где $Con_i,...,Con_n$ — противоречия, разрешенные в рамках итерации ТРИЗ-эволюции; $IP_j,...,IP_m$ — инструменты ТРИЗ, позволяющие сформировать и/или описать решение при переходе от одного объекта ТРИЗ-эволюции к другому; Sol_x — решение, согласно которому разрешены перечисленные противоречия; $MP_k,...,MP_y$ — изменяемые в результате решения параметры.

Далее *Шаги 1-4* повторяются для всех наиболее значимых объектов исследуемой области.

На этапе II происходит графическое представление итераций ТРИЗ-эволюции выявленных и описанных на этапе I. Графическая форма итерации представлена на рисунке 1.1.

На этапе III ТРИЗ-эволюционная карта анализируется с целью определить неразрешенные противоречия и противоречия, разрешение которых привело к повышению идеальности не за счет устранения мешающего параметра. Данные противоречия позволят сформулировать постановку задачи перспективных направлений развития системы. Кроме того, по ТРИЗ-эволюционной карте можно отследить качественные скачки роста идеальности и ключевые тенденции развития системы.

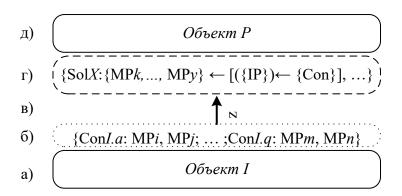


Рисунок 1.1 — Графическое представление итерации ТРИЗ-эволюционной карты: (а) — исходный объект итерации; (б) — множество противоречий исходного объекта итерации; (в) — переход к конечному объекту итерации, где Z — порядковый номер итерации на ТРИЗ-эволюционной карте; (г) — множество разрешений противоречий исходного объекта итерации, реализованных в конечном объекте итерации; (д) — конечный объект итерации

1.1.2 ТРИЗ-эволюционный анализ систем МП

Отправной точкой автоматизации МП можно считать 1933 год [26]. Именно в этом году впервые были запатентованы независимо друг от друга две машины для осуществления перевода: «Механический мозг» Жоржа Арцруни и «Машина для подбора и печатания слов при переводе с одного языка на другой или на несколько других одновременно» П. П. Троянского (патент СССР № 40995 от 5 сентября 1933 г.) [27]. Развиваться же направление МП стало гораздо позднее, во времена холодной войны.

Первым подходом к автоматизации перевода считается пословный перевод [28, 29, 30]. Ключевая идея: разделяем предложение по словам, находим перевод по двуязычному словарю, подставляем перевод на место каждого слова. Недостатки подхода:

- чем сложнее предложение, тем хуже качество перевода;
- требует долгого и кропотливого формирования словарей, что трудозатратно и ограничивает количество языковых пар;
- не учитывает разницу грамматический строй, морфологию, согласование падежей;
 - не учитывает контекстно-зависимые значения слов.

Исходя из описанных недостатков опишем противоречия данного подхода, используя формулу (1) (см. Таблицу 1.1).

Таблица 1.1 – Спецификация противоречий пословного подхода к переводу

Описание противоречия	Формула
	противоречия
Противоречие 1.1: при увеличении количества возможных	Con1.1: MP2↑, MP1.1↓
смысловых значений одной лексической единицы недопустимо	
снижается качество перевода	
Противоречие 1.2: при повышении качества перевода	Con1.2: MP1↑, MP3↓
работоспособность системы недопустимо снижается	
Противоречие 1.3: при увеличении количества языковых пар	Con1.3: MP4↑, MP5↑
недопустимо увеличивается время на составление обучающих	
данных	
Противоречие 1.4: при повышении сложности структуры входных	Con1.4: MP6↑, MP1↓
данных качество перевода недопустимо снижается	

Рассмотрим пример разрешения противоречия Con1.2. Противоречие может быть разрешено при помощи приема «Принцип предварительного действия» (IP10): предварительно сформировать набор правил, в соответствии с которыми будет обрабатываться текст. Данное решение было реализовано в системах дословного перевода. Произошли первая итерация ТРИЗ-эволюции и переход к новой парадигме «Машинный перевод на основе правил». Представим данную итерацию в символьном виде согласно формуле (2):

$$Con1.2: MP1\uparrow, MP3\downarrow \xrightarrow{IP:10} Sol1:MP3\uparrow$$

Графически первая итерация ТРИЗ-эволюции систем МП представлена на рисунке 1.2. Полный перечень итераций ТРИЗ-эволюции от исходного объекта приведен в таблице 1.2. В качестве объектов ТРИЗ-эволюции можно выделить следующие системы МП: дословный МП [31], трансферный МП [32], интерлингвистический МП [33, 34], МП на примерах [35], интерактивный МП [36], статистический МП по словам [37], статистический МП по фразам [38], статистический МП на основе синтаксиса [39], нейронный МП [40], нейронный МП «без учителя» [41], адаптированный нейронный МП [42], гибридный МП [43]. Далее рассмотрим краткую характеристику каждого вида указанных систем МП.

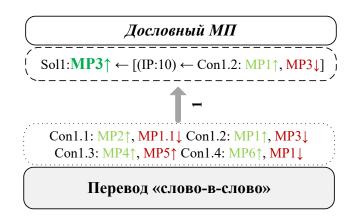


Рисунок 1.2 – Первая итерация ТРИЗ-эволюции систем машинного перевода

Дословный МП: текст разделяется по словам; система ищет перевод каждого слова на основе правил; на основе правил правятся ошибки морфологии, согласования падежей, использования окончаний и т.д.

Трансферный МП: текст на исходном языке анализируется для определения его грамматической структуры; грамматическая структура (члены предложения) перестраивается в структуру, подходящую для воспроизводства текста на языке перевода; выполняется подстановка дословного перевода каждой отдельной части текста в нужное место конечной структуры).

Интерлингвистический МП: исходное предложение при помощи словаря и грамматических правил исходного языка полностью конвертируется в промежуточное представление (совокупность концептов), единое для всех языков мира (interlingua), с последующей конвертацией его в текст на языке перевода при помощи словаря и грамматических правил языка перевода.

МП на примерах: в систему загружаются базы переводов, которыми являются двуязычные корпусы текста (оригинал-перевод); система выполняет поиск ранее переведенных частей, определяет отличные части текста и переводит только отличные от хранящихся в базе переводов части, подставив остальные части из готового перевода.

Uнтерактивный $M\Pi$: на разных стадиях перевода привлекается участие человека. Такое участие может быть выражено в разных формах: с постредактированием, когда человек редактирует уже переведенный машиной текст; с

предредактированием, когда человек редактирует входной текст, приспосабливая его для более легкого понимания машиной; частично автоматизированный перевод, когда человек и машина взаимодействуют в процессе перевода, и участие человека в процессе необходимо для разрешения трудностей.

Статистический МП по словам: система обучается на параллельных корпусах которые разделены на слова и проанализированы по критерию частотности отдельных слов; система определяет порядок слов в переведенных предложениях, в том числе относительный порядок; выполняется перевод отдельных слов путем подстановки статистически более вероятных вариантов на основе обучающих данных; после перевода слова меняются местами согласно статистически более естественному положению; на те места, где предполагается необходимость нового слова, вставляются маркеры (NULL-слова); под каждый маркер выполняется подбор нужного слова (артикля, частицы или глагола).

Статистический МП по фразам: принцип работы совпадает с системами статистического МП по словам с той разницей, для обучения и перевода текст делится не только на слова, но и на n-граммы – пересекающиеся наборы из N слов подряд.

Статистический МП на основе синтаксиса: перед переводом выполняется точный синтаксический разбор предложения, на основании чего строится дерево, используя которое, теоретически, можно обучить систему правильно конвертировать грамматические структуры одного языка в структуры другого, выполняя перевод по словам или фразам. Практически данный подход не был реализован до конца из-за сложности задачи автоматического синтаксического разбора.

Нейронный МП: система состоит из двух нейронных сетей, каждая из который знает только один язык; одна сеть кодирует предложение на языке оригинала в набор цифровых характеристик, а вторая, проанализировав предложение до конца, причем одновременно и слева направо, и справа налево, начинает декодировать числовые характеристики и предсказывать слова перевода, причем каждое предсказанное слово используется для предсказания следующего слова; для более точного выбора слов для перевода используется контекст всего исходного предложения, а также

контекст всех предыдущих предсказанных слов; система обучается на двуязычных корпусах посредством алгоритмов глубокого обучения.

Адаптированный нейронный МП: система нейронного перевода, корпуса для обучения которой составляют тексты некоторой предметной области, исключая попадания в обучающий корпус текстов других предметных областей/тематик.

Нейронный МП «без учителя»: система строится на нейронных сетях, архитектура которых позволяет производить обучение «без учителя» на моно-язычных корпусах текстов, то есть без эталонных данных для проверки перевода.

 Γ ибридный $M\Pi$: система, в которую интегрированы разные технологии $M\Pi$ для достижения лучшего качества перевода. Наиболее часто в таких системах интегрируют нейронный и статистический $M\Pi$.

В таблице 1.2 представлен перечень выявленных в рамках анализа главных параметров систем МП.

Таблица 1.2 – Описание главных параметров систем МП

Обозначение	Описание
MP1	Качество перевода (МР1.1 – с точки зрения лексической точности; МР1.2 –
	с точки зрения грамматической точности; МР1.3 – с точки зрения стилисти-
	ческой точности; МР1.4 – с точки зрения единообразия; МР1.5 – с точки
	зрения смысловой точности; МР1.6 – гладкость перевода; МР1.7 – длинных
	предложений; МР1.8 – коротких предложений; МР1.9 – текстов узких тема-
	тик)
MP2	Количество возможных смысловых значений одной лексической единицы
MP3	Работоспособность системы
MP4	Количество языковых пар
MP5	Время на подготовку обучающих данных (МР5.1 – двуязычных словарей;
	МР5.2 – параллельных корпусов; МР5.3 – правил; МР5.4 – моно-корпусов)
MP6	Сложность структуры текста оригинала
MP7	Объем обучающих данных (МР7.1 – правил; МР7.2 – корпусов текста)
MP8	Время на разработку системы
MP9	Трудозатраты на сопровождение системы
MP10	Количество вариантов сочетаний слов
MP11	Трудозатраты на пред-, пост-редактирование
MP12	Время на обучение системы
MP13	Вероятность грубых ошибок
MP14	Объем исходного текста
MP15	Время на поиск ошибок в переведенном тексте
MP16	Количество переводимых тематик
MP17	Сложность системы

Для каждой из систем произведен подробный анализ с выявлением ключевых проблем, ограничивающих рост идеальности, на основании которых сформулированы противоречия. Спецификация выявленных проблем и противоречий в системах МП представлена в таблице 1.3.

Таблица 1.3 – Спецификация противоречий систем машинного перевода

Orwoodyna wa ofizona y	Пи отуго почило	Символьное	
Описание проблемы	Противоречие	представление	
Д	ословный МП	,	
Практически невозможно язык представить в виде лишь набора правил, следовало бы также учесть и все возможные исключения из правил (неправильные глаголы в английском, плавающие приставки в немецком, суффиксы, диалекты, сленг и т.д.	Противоречие 2.1: при повышении качества перевода недопустимо увеличивается объем обучающих данных (правил)	Con2.1: MP1↑, MP7.1↑	
Количество правил в каждом языке огромно и для качественной проработки необходимых обучающих данных необходимо огромное количество человеко-часов	Противоречие 2.2: при повышении качества перевода недопустимо увеличивается время на подготовку обучающих данных	Con2.2: MP1↑, MP5↑	
Необходимо постоянно поддерживать лингвистическую базу в актуальном состоянии, так как язык - динамическая система	Противоречие 2.3: при повышении качества перевода недопустимо увеличиваются трудозатраты на сопровождение системы	Con2.3: MP1↑, MP9↑	
Тр	ансферный МП		
С одной стороны можно задать общие правила переноса грамматической структуры, что упрощает задачу пере-	Противоречие 3.1: при сокращении времени на разработку системы недопустимо увеличивается объем обучающих данных (правил)	Con3.1: MP9↓, MP7.1↑	
вода, с другой стороны, сочетаний слов намного больше, чем самих слов, и каждый вариант почти невозможно учесть.	Противоречие 3.2: при сокращении количества правил недопустимо снижается качество перевода возможных сочетаний слов	Con3.2: MP7.1↓, MP1.1↓	
Интерлингвистический МП			
Сложность реализации и отсутствие методов и моделей поиска закономерностей и классификации атрибутов текста для создания унифицированного языка и его структуры.	Противоречие 4.1: при повышении качества перевода недопустимо увеличивается время на разработку системы	Con4.1: MP1†, MP8†	
МП на примерах			
Примеры содержат слова, словосочетания и даже предложения, но, фактически, мы находим дословно схожие части, не учитывая особенности синтаксиса, морфологии, грамматического	Противоречие 5.1: при повышении качества перевода объем обучающих данных (корпусов текста) недопустимо увеличивается	Con5.1: MP1↑, MP7.2↑	

Описание проблемы	Противоречие	Символьное
	12perimepe me	представление
строя и т.д. Чтобы учесть все возмож-		
ные варианты, необходимо больше обучающих данных.	Противоречие 5.2: при увеличении количества вариантов сочетаний слов в тексте оригинала качество перевода с точки зрения единообразия недопустимо снижается	Con5.2: MP10↑, MP1.4↓
Даже если объем обучающих данных		
достаточный, система не делит предложение на структурные части из-за чего, например, служебные части речи, влияющие на контекст, могут отразится на качестве перевода.	Противоречие 5.3: при повышении качества с точки зрения единообразия недопустимо снижается качество с точки зрения передачи смысла	Con5.3: MP1.4↑, MP1.5↓
Из-за разницы структуры языков перевода опущенные (нулевые) части предложений не учитываются в переводе, либо переводятся части предложения, которые ввиду грамматических правил языка перевода должны быть опущены.	Противоречие 5.4: при повышении качества с точки зрения единообразия недопустимо снижается качество перевода с точки зрения грамматической точности	Con5.4: MP1.5↑, MP1.2↓
Инт	герактивный МП	
Система не может работать автоматически без участия человека	Противоречие 6.1: при повышении качества перевода недопустимо увеличиваются трудозатраты на пред-, пост-редактирование	Con6.1: MP1↑, MP11↑
	ческий МП по словам	T
При подготовке обучающих данных необходимо максимально точное соответствие оригинала и перевода. Однако, не всегда перевод может быть строго формализован, есть еще литературные или вольные переводы, которые также необходимо учитывать.	Противоречие 7.1: при повышении качества перевода недопустимо увеличивается время на подготовку обучающих данных (корпусов)	Con7.1: MP1↑, MP5.2↑
Из-за отсутствия двуязычных словарей между некоторыми языками, система переводит текст сначала на английский, а затем на язык перевода из-за чего возникают «двойные потери» качества.	Противоречие 7.2: при повышении количества языковых пар недопустимо снижается качества перевода	Con7.2: MP4↑, MP1↓
	ческий МП по фразам	
Статистические аномалии	Противоречие 8.1: при увеличении объема обучающих данных недопустимо снижается качество перевода с точки зрения смысловой точности	Con8.1: MP7↑, MP1.5↓
Отдельные фразы плохо согласуются между собой, в итоге переведенное	Противоречие 8.2: при повышении качества с точки зрения лексической точности недопустимо снижается гладкость перевода	Con8.2: MP1.1↑, MP1.6↓
предложение – набор фраз, иногда не связных по смыслу	Противоречие 8.3: при повышении качества с точки зрения лексической точности недопустимо снижается смысловая точность	Con8.3: MP1.1↑, MP1.5↓

Описание проблемы	Противоречие	Символьное
Описание прослемы	Прогиворечие	представление
Статистический МП на основе синтаксиса		
Даже для корпуса с простейшими 2-3 уровневыми деревьями, время обучения слишком велико, а значит на практике система не может быть применима.	Противоречие 9.1: при повышении качества перевода с точки зрения грамматической точности время на обучение системы недопустимо увеличивается	Con9.1: MP1.2↑, MP12↑
Не для всех языков разработаны методы синтаксического разбора. Не для всех языков, для которых разработаны методы синтаксического анализа, они работают достаточно качественно.	Противоречие 9.2: при повышении количества языковых пар недопустимо снижается качество перевода	Con9.2: MP4↑, MP1↓
	Іейронный МП	
В целом «гладкий» перевод может содержать грубые лексические ошибки.	Противоречие 10.1: при повышении гладкости перевода недопустимо повышается вероятность наличия в переводе грубых ошибок	Con10.1: MP1.6↑, MP13↑
Почти любой перевод требует понимания контекста нескольких предложений, иногда это имеет решающее значение для точного перевода с точки зрения используемой лексики. Нейронная система не может анализировать и хранить информацию о контексте текста большого объема и эффективно ее запоминать.	Противоречие 10.2: при увеличении объема исходного текста недопустимо снижается качество перевода с точки зрения смысловой точности	Con10.2: MP14↑, MP1.5↓
Зависимость от состава обучающих данных: если тренировать нейронный перевод только на длинных парах предложений, система будет неспособна перевести корректно короткое предложение или даже отдельное слово.	Противоречие 10.3: при повышении качества перевода длинных предложений недопустимо снижается качество перевода коротких предложений	Con10.3: MP1.7↑, MP1.8↓
Аномалии в переводе: пропущенные отрицания, отдельные слова или целые фразы. Аномалии непредсказуемы и непоследовательны, что затрудняет их автоматическое выявление и исправление.	Противоречие 10.4: при повышении гладкости перевода недопустимо увеличивается время на поиск возможных ошибок	Con10.4: MP1.6↑, MP15↑
Низкое качество перевода исходных текстов, которые сильно отличаются от	Противоречие 10.5: при увеличении количества тематик исходного текста недопустимо увеличивается время на подготовку обучающих данных	Con10.5: MP16↑, MP5↑
данных, использованных для машинного обучения.	Противоречие 10.6: при увеличении количества тематик исходного текста недопустимо снижается качество перевода	Con10.6: MP16↑, MP1.9↓

Описание проблемы	Противоречие	Символьное представление
Адаптирог	ванный нейронный МП	•
По узким тематикам чрезвычайно	Противоречие 11.1: при повышении	
сложно производить сбор и обработку	качества перевода текстов узкой те-	Con11.1:
релевантных двуязычных корпусов для	матики и/или редкой языковой пары	MP1.9↑,
обучения нейронной сети в достаточ-	недопустимо увеличивается время	MP5.2↑
ном объеме.	на подготовку обучающих данных	
Нейронн	ный МП «без учителя»	
Необходимость подготовки моно-корпусов большого объема по редким языкам и/или тематикам в векторном представлении.	Противоречие 12.1: при повышении качества перевода редких языковых пар или тематик объем необходимых обучающих данных недопустимо увеличивается	Con12.1: MP1↑, MP7↑
Гибридный МП		
При комбинировании разных систем	Противоречие 13.1: при повышении качества перевода недопустимо увеличивается сложность конечной системы	Con13.1: MP1↑, MP17↑
достигается более высокое качество перевода, но это означает усложнение системы, которая наследует не только преимущества, но и недостатки систем,	Противоречие 13.2: при увеличении количества тематик недопустимо увеличивается время на разработку системы	Con13.2: MP16↑, MP8↑
входящих в ее состав.	Противоречие 13.3: при увеличении количества языковых пар недопустимо увеличивается время на разработку системы	Con13.3: MP4↑, MP8↑

В таблице 1.4 произведен анализ выявленных противоречий и инструментов ТРИЗ, при помощи которых были разрешены противоречия в описанных системах, описаны итерации ТРИЗ-эволюции систем МП.

Таблица 1.4 – Спецификация итераций ТРИЗ-эволюции систем МП

№ п/п	Переход	Итерация	Описание решения
		Перевод «слово-в-слов	B0»
1	Перевод «слово-в- слово» → Дословный МП	Con1.2: MP1 \uparrow , MP3 \downarrow $\xrightarrow{IP:10}$ Sol1:MP3 \uparrow	Предварительно сформировать набор правил, в соответствии с которыми будет обрабатываться текст.
2	Перевод «слово-в- слово» → Трансферный МП	Con1.4: MP6 \uparrow , MP1 \downarrow $\xrightarrow{IP:1}$ Sol2:MP6 \downarrow	Разбить сложное предложение на набор синтаксических структур.
3	Перевод «слово-в- слово» → Интерлингвистический МП	Con1.3: MP4 \uparrow , MP5 \uparrow Sol3:MP5 \downarrow	Создать универсальный язык, переводить сначала на универсальный язык, а затем на целевой язык перевода. Для каждого языка прописать правила конвертации только

№ п/п	Переход	Итерация	Описание решения
			в одной языковой паре с промежуточным представлением.
4	Перевод «слово-в- слово» → Интерактивный МП	$\frac{\text{Con1.1: MP2}\uparrow, \text{MP1.1}\downarrow}{\overset{\text{IP:15,16,23}}{\longrightarrow}} \text{Sol4:MP1}\uparrow$	Переводить с погрешностью, предлагать пользователю в процессе перевода выбрать необходимый перевод слова, ориентируясь на контекст. Перестроить оставшуюся часть перевода в соответствии с внесенными пользователем изменениями.
		$\begin{array}{c} \text{Con1.2: MP1}\uparrow, \text{MP3}\downarrow^{\stackrel{\text{IP:16}}{\longrightarrow}} \\ \text{Sol5:MP1}\uparrow \end{array}$	Переводить как есть с погрешностями, связанными с грамматикой и стилистикой, готовый перевод пост-редактировать.
		$ \begin{array}{c} \text{Con1.4: MP6}\uparrow, \\ \text{MP1}\downarrow^{\text{IP:10, 24}} \text{Sol6:MP6}\downarrow \end{array} $	Перед переводом либо в процессе перевода редактировать исходный текст таким образом, чтобы он стал более «понятным» для машины.
5	Перевод «слово-в- слово» → МП на примерах	$ \begin{array}{c} \text{Con1.1: MP2}\uparrow, \text{MP1.1}\downarrow \\ \xrightarrow{\text{IP:1,26,33}} \text{Sol7:MP1.1}\uparrow \end{array} $	Разработать систему, которая будет делить текст на смысловые части, находить в базе и копировать перевод ранее переведенных частей, генерируя перевод только для нового фрагмента.
6	Перевод «слово-в- слово» → СМП по словам	Con1.3: MP4 \uparrow , MP5 $\uparrow \longrightarrow$ Sol8: MP5.1 \downarrow	Вместо словарей в формате словослово обучать систему на параллельных корпусах в формате текстэталонный перевод
7	Перевод «слово-в- слово» → СМП на основе син- таксиса	$Con1.2: MP1\uparrow,$ $MP3\downarrow \xrightarrow{IP:10,17} Sol9:MP3\uparrow$	Перед переводом производить полный синтаксический разбор предложения, структурируя части текста при помощи деревьев, затем переводить по словам.
8	Перевод «слово-в- слово» → СМП по фразам	Con1.1: MP2 \uparrow , MP1.1 \downarrow & Con1.4: MP6 \uparrow , MP1 \downarrow $\xrightarrow{IP:1, 16, 17}$ Sol10: MP1 \uparrow	Разбить предложение на фразы во всех возможных сочетаниях по пслов в каждом и проанализировать перевод для каждого из них, выбрав в последствии только статистически наиболее верные сочетания, т.е. учитывая контекст.
		Дословный МП	
9	Дословный МП → Интерактивный МП	Con2.1: MP1↑, MP7.1↑ ——————————————Sol11:MP7.1↓	Переводить с погрешностью, предлагать пользователю скорректировать перевод с целью улучшения качества.
10	Дословный МП → МП на примерах	Con2.1: MP1 \uparrow , MP7.1 \uparrow $\xrightarrow{\text{IP:}2,26}$ Sol12: {MP1.4 \uparrow , MP1.3 \uparrow , MP7.1 \downarrow }	Не описывать правила вообще, загрузить в систему большое количество примеров, из которых система

№ п/п	Переход	Итерация	Описание решения	
			будет брать готовые части перевода.	
11	Дословный МП → Трансферный МП	Con2.2: MP1 \uparrow , MP5 \uparrow $\xrightarrow{\text{IP}:5,16} \text{Sol13: MP5}\downarrow$	Выделить ключевые синтаксические конструкции. Заложить в систему правила перевода каждого слова и подстановки в соответствии с синтаксическими конструкциями языка перевода.	
12	Дословный МП → СМП по словам	Con2.2: MP1 \uparrow , MP5 \uparrow $\xrightarrow{\text{IP:2,24}} \text{Sol14: MP5}\downarrow$	Создать промежуточный алгоритм обучения системы: на основе статистических методов система будет выбирать наиболее статистически вероятный перевода, при этом нет необходимости описывать правила.	
		Con2.3: MP1 \uparrow , MP9 \uparrow $\xrightarrow{\text{IP}:19}$ Sol15: MP9 \downarrow	Периодически подгружаем в систему новые параллельные корпуса текста, на которых система доучивается.	
		Трансферный МП		
13	Трансферный МП → МП на примерах	Con3.1: MP9 \downarrow , MP7.1 $\uparrow \xrightarrow{\text{IP:3,26}}$ Sol16:MP7.1 \downarrow	Не описывать правила, загрузить в систему большое количество примеров, из которых система будет брать готовые части перевода.	
14	Трансферный МП → СМП по словам	Con3.1: MP9 \downarrow , MP7.1 \uparrow $\xrightarrow{\text{IP:3,24}}$ Sol17: MP7.1 \downarrow	Создаем промежуточный алгоритм обучения системы: на основе статистических методов система будет выбирать наиболее статистически вероятный вариант перевода, и не нужно описывать правила.	
15	Трансферный МП → СМП по фразам	Con3.2: MP7.1 \downarrow , MP1.1 \downarrow $\xrightarrow{\text{IP:22}}$ Sol18: MP1.1 \uparrow	Найти все возможные сочетания слов в тексте, при помощи статистических методов найти наиболее вероятный перевод.	
16	Трансферный МП → Интерактивный МП	Con3.2: MP7.1 \downarrow , MP1.1 \downarrow $\xrightarrow{\text{IP}:10,16}$ Sol19: MP1.1 \uparrow	Предварительно обрабатывать оригинал для того, чтобы он был более простым и «понятным» машине, переводить с погрешностью, привлекая пост-редактора для повышения качества перевода.	
Интерлингвистический МП				
17	Интерлингвистический $M\Pi \rightarrow CM\Pi$ по словам	Con4.1: MP1 \uparrow , MP8 \uparrow $\xrightarrow{\text{IP}:24,27}$ Sol20: MP8 \downarrow	Для редких языков и/или тематик из-за сложности создания и разметки универсальной модели и сбора обучающих данных использовать английский язык в качестве промежуточного.	

№ п/п	Переход	Итерация	Описание решения				
18	Интерлингвистический $M\Pi \rightarrow$	Con4.1: MP1 \uparrow , MP8 \uparrow $\xrightarrow{\text{IP:17}}$ Sol21:	Для кодировки текста в универ- сальное представление использо-				
	Нейронный МП	{MP1↑,MP8↓}	вать многослойную нейронную сеть.				
	МП на примерах						
		Con5.1: MP1↑, MP7.2↑	Перед копированием перевода вы-				
		$\xrightarrow{\text{IP:10}} \text{Sol22:MP1.1} \uparrow$	брать из множества вариантов наиболее вероятный на основе статистических методов.				
		Con5.2: MP10↑, MP1.4↓	Записывать порядок слов в предло-				
19	МП на примерах → СМП по словам	$\xrightarrow{\text{IP:}10} \text{Sol23:MP1.4} \uparrow$	жении и относительно друг друга, при переводе выбирать наиболее вероятный вариант на основе ста-				
			тистических методов				
		Con5.4: MP1.5↑, MP1.2↓	Заранее промаркировать места, где				
		$\stackrel{\text{IP:11}}{\longrightarrow} \text{Sol24:MP1.2}\uparrow$	необходимо подставить служебное				
			слово, при переводе заполнить маркеры.				
		Con5.2: MP10↑, MP1.4↓	Найти все возможные сочетания				
		$\stackrel{\text{IP:22}}{\longrightarrow} \text{Sol25:} \{\text{MP1.4}\uparrow, \\ \text{MP1.1}\uparrow\}$	слов в тексте, при помощи стати-				
		MP1.1↑}	стических методов найти наиболее				
	$M\Pi$ на примерах \rightarrow	.,	вероятный перевод.				
20	СМП по фразам	Con5.3: MP1.4↑, MP1.5↓	При помощи статистических мето-				
	11	$\xrightarrow{\text{IP:5,10}} \text{Sol26:MP1} \uparrow$	дов находить среди множества ва-				
		·	риантов обучающего корпуса				
			наиболее вероятный перевод це-				
			лых фраз, а не отдельных слов.				
		{Con5.2: MP10↑,	Попольного п				
	МП на примерах →	MP1.4 \downarrow , Con5.4:	Перед переводом производить				
21	СМП на основе син-	IP:10.17	полный синтаксический разбор				
	таксиса	$MP1.5\uparrow, MP1.2\downarrow\} \xrightarrow{a.i.s,i.}$	предложения, структурируя части текста при помощи деревьев.				
		Sol27:{MP1.4↑,	текста при помощи дереввев.				
		MP1.2↑}	T				
22	Интерактивный МП 22 Интерактивный МП → Con6.1: MP1↑, Использовать в качестве посред-						
	Гибридный МП	$MP11\uparrow \xrightarrow{IP:0, 3, 5}$	ника между машинным переводом				
	т поридный тип		и человеком среду автоматизации				
		Sol28:{MP1 \uparrow ,MP11 \downarrow }	переводов, объединив все цифро-				
			вые инструменты, облегчающие				
			перевод в одной среде. Использо-				
			вать МП для тех фрагментов, для				
			которых не нашлось совпадений по				
			памяти перевода, при этом выби-				
			рать модель перевода исходя из				
		C	особенности задачи.				
22		Статистический МП по с					
23	СМП по словам →	Con7.1: MP1↑, MP5.2↑	Объединить отдельные слова в				
	СМП по фразам	$\xrightarrow{\text{IP:5,16,17}} \text{Sol29:MP5.2} \downarrow$	предложении в фразы во всех воз-				
			можных сочетаниях по п-слов в				

№ п/п	Переход	Итерация	Описание решения
			каждом и проанализировать перевод для каждого из них, выбрав в последствии только статистически наиболее вероятные сочетания, т.е. учитывать контекст.
24	СМП по словам → Нейронный МП	Con7.2: MP4 \uparrow , MP1 \downarrow $\xrightarrow{\text{IP:6,17}}$ Sol30: $\{\text{MP1}\uparrow,\text{MP4}\uparrow\}$	В качестве алгоритма вывода использовать нейронные сети, состоящие из энкодера и декодера. Энкодер сначала кодирует исходный текст в универсальное представление. Декодер распознает универсальное представление и переводит его на язык перевода.
		Статистический МП по ф	рразам
25	СМП по фразам → Гибридный МП	$ {Con8.2: MP1.1↑, MP1.6↓, Con8.3: } MP1.1↑, MP1.5↓} \xrightarrow{IP:3,5} Sol31:MP1↑ $	Объединить в одной системе несколько типов систем, воспользовавшись преимуществами каждой из них, разделив задачи. Например: перевод отдельных слов по словарю и использование правил, а «выравнивание» текста доверить СМП или НМП, НМП переводить длинные предложения, СМП пере-
26	СМП по фразам →	{Con8.2: MP1.1↑,	водить короткие фразы. Учитывать контекст всего предло-
20	Нейронный МП	$\begin{array}{c} \text{MP1.6}\downarrow, \text{Con8.3}:\\ \text{MP1.1}\uparrow, \text{MP1.5}\downarrow\rbrace \xrightarrow{\text{IP:17,20}}\\ \text{Sol32: MP1}\uparrow \end{array}$	жения в переводе, а не только отдельных фраз при помощи многослойных нейронных сетей с долгосрочной памятью.
	СМП по фразам →	Con8.1: MP7↑, MP1.5↓	Решение отсутствует
	**	стический МП на основе	· · ·
27	Syntax -based SMT → Нейронный МП	{Con9.1: MP1.2 \uparrow , MP12 \uparrow , Con9.2: MP4 \uparrow , MP1 \downarrow } $\xrightarrow{\text{IP:2, 17}}$ Sol33:MP1 \uparrow	Использовать многослойную нейронную сеть, не использовать методы синтаксического разбора и анализа.
		Нейронный МП	
28	Нейронный МП → Гибридный МП	Con10.3: MP1.7 \uparrow , MP1.8 \downarrow $\xrightarrow{\text{IP:3,5}}$ Sol34: {MP1.7 \uparrow , MP1.8 \uparrow }	Объединить в одной системе два типа систем, выбрав наиболее подходящую технологию отдельно для длинных предложений, отдельно для коротких фраз.
29	Нейронный МП → НМП «без учителя»	Con10.5: MP16 \uparrow , MP5 \uparrow $\stackrel{\text{IP:2}}{\longrightarrow} \text{Sol35: MP5.2}\downarrow$	Обучать сеть алгоритмом без учителя, исключив тем самым необходимость подготовки эталонной выборки данных.
30	Нейронный МП → Адаптированный НМП	Con10.6: MP16 \uparrow , MP1.9 $\downarrow \xrightarrow{\text{IP:5}}$ Sol36:MP1.9 \uparrow	Для каждой из тематик готовить корпуса обучающих данных отдельно.

№ п/п	Переход	Итерация	Описание решения		
	Нейронный МП →	Con10.1: MP1.6↑, MP13↑	Решение отсутствует		
	Нейронный МП →	Con10.2: MP14↑, MP1.5↓	Решение отсутствует		
	Нейронный МП →	Con10.4: MP1.6↑, MP15↑	Решение отсутствует		
Адаптированный НМП					
31	Адаптированный НМП → НМП «без учителя»	Con11.1: MP1.9 \uparrow , MP5.2 $\uparrow \xrightarrow{\text{IP:2, 17}}$ Sol37: MP5.2 \downarrow	Обучать нейронную сеть алгоритмом «без учителя» на моно-корпусах текста вместо параллельных корпусов.		
НМП «без учителя»					
	$\mathrm{HM}\Pi$ «без учителя» $ ightarrow$	Con12.1: MP1↑, MP7↑	Решение отсутствует		
Гибридный МП					
	Гибридный МП $ ightarrow$	Con13.1: MP1↑, MP17↑	Решение отсутствует		
	Гибридный МП $ ightarrow$	Con13.2: MP16↑, MP8↑	Решение отсутствует		
	Гибридный МП $ ightarrow$	Con13.3: MP4↑, MP8↑	Решение отсутствует		

Полная ТРИЗ-эволюционная карта представлена на рисунке 1.3.

Анализ ТРИЗ-эволюции систем МП [44, 45] показывает в первую очередь качественные скачки развития, которые позволили сформировать новые парадигмы: переход к трансферному МП привел к созданию парадигмы «перевод, основанный на правилах»; переход к СМП по словам – к созданию парадигмы «статистический МП»; переход к НМП – «к созданию парадигмы нейронный МП». По улучшаемым параметрам в ходе развития систем МП мы видим, что оно происходило по следующим ключевым направлениям: повышение качества перевода, причем с развитием МП и сам параметр качества трансформировался и детализировался в подпараметры; сокращение времени на сбор и подготовку обучающих данных; развитие способов обработки оригинала с целью упрощения его структуры для более точного понимания системой семантики текста; совершенствование технической реализации МП. Часть противоречий существующих систем не разрешены (8.1, 10.1, 10.2, 10.4, 12.1, 13.1, 13.2, 13.3), и, следовательно, задачи по разрешению этих противоречий являются перспективными направлениями исследования в области МП.

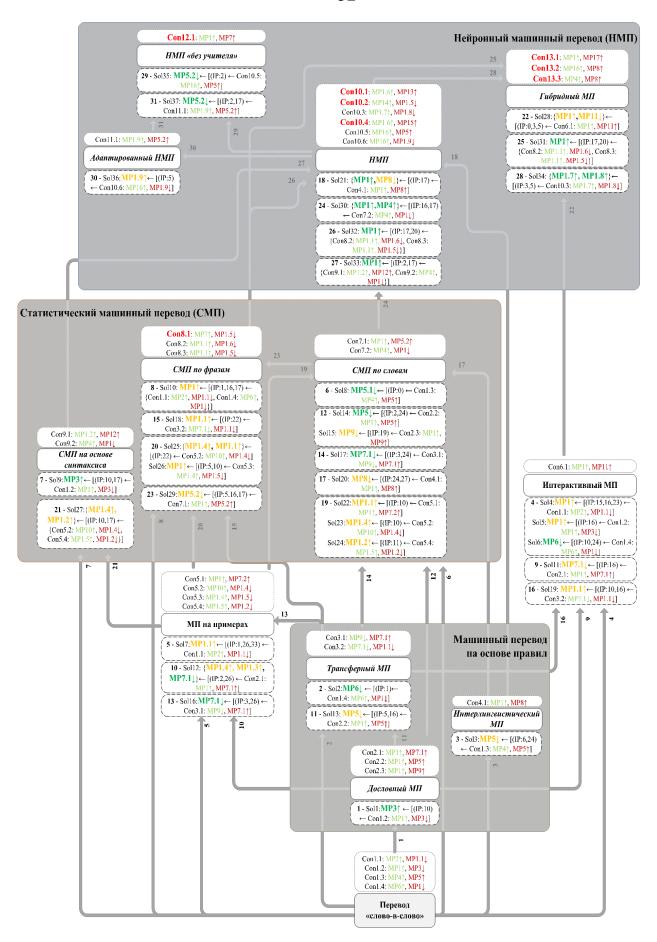


Рисунок 1.3 – ТРИЗ-эволюционная карта систем машинного перевода

По ТРИЗ-эволюционной карте и неразрешенным противоречиям мы видим, что одним из направлений совершенствования систем МП является развитие методов интерактивного перевода и нативного взаимодействия с пользователями МП. В системах МП не развиты методы обратной связи с пользователем, соответственно пользователь не может оценить риски некачественного перевода текста. из-за того, что система не может достаточно точно определять контекст текста целиком, а не отдельного предложения, оттенки смысла и т.д.

1.2 Интерактивный перевод и предредактирование

Несмотря на выдающийся прорыв МП в вопросах качества перевода, его применимость в прикладных задачах перевода все еще вызывает вопросы у пользователей и исследователей [46, 47], особенно при переводе узкоспециальных текстов [48]. В частности, С. Canfora и А. Оttmann в своем исследовании [49] поднимают вопрос рисков, связанных с высокими ожиданиями пользователей и индустрии переводов от внедрения технологий МП и ответственности за некачественный перевод в узкоспециальных предметных областях.

Очевидно, что для корректного использования систем МП необходимо прямое участие человека на всех этапах работы системы — от загрузки оригинала, сбора релевантных обучающих данных до оценки качества выполненного перевода [50]. Справедливо правило, что для достижения оптимального результата работы, пользователь системы МП должен знать хотя бы один язык из языковой пары перевода, и чем лучше знание языка перевода пользователя, тем точнее будет оценка качества МП и его пост-редактура. Другими словами, конечное качество перевода улучшается, если пользователь системы МП является носителем языка перевода.

Принимая во внимание этот принцип и рассматривая систему МП как средство автоматизированной поддержки работы переводчика, которое предоставляет переводчику модели и алгоритмы поддержки принятия решения относительно перевода, можно заключить, что, предоставляя пользователю средства обработки текста на языке, носителем которого он является, на любом из этапов перевода, можно

повысить качество перевода. Ключевая идея данного подхода [51] в том, что на разных стадиях перевода привлекается участие человека. Такое участие может быть выражено в разных формах:

- с постредактированием: человек редактирует результат МП;
- *с предредактированием*: человек редактирует входной текст, приспосабливая его для более легкого понимания машиной;
- частично автоматизированный перевод: человек и машина взаимодействуют в процессе перевода. Например, использование человеком электронных словарей при переводе; участие человека в процессе машинного перевода для разрешения трудностей.

Работы по частично-автоматизированному переводу [52, 53, 54] обычно рассматривают варианты совершенствования автоматизированного инструментария переводчиков с целью повышения скорости и качества их работы.

Множество работ и исследований посвящено постредактированию машинного перевода. В литературе описываются стратегии постредактирования, вопросы подготовки постредакторов, типичные ошибки машинного перевода, автоматизации постредактирования т.д. [55, 56, 57, 58]. Предредактированию посвящено меньше работ, однако существующие практические статьи указывают на эффективность такого подхода к повышению качества машинного перевода [59, 60].

Seretan, Bouillon и Gerlach в своем исследовании [61] показали, что использование даже простых полуавтоматических правил предредактирования текста повышает качество МП. Гипотеза была проверена на сообщениях пользователей технического форума и медицинских текстах и показала одинаковую эффективность. Авторы указывают, что реализация методов автоматического предредактирования возможна только для языков, для которых уже созданы поверхностные синтаксические анализаторы.

Ряд исследований сравнивают подходы постредактирования и предредактирования. Так, исследование С. Shei [62] в языковой паре китайский-английский показало преимущества и позитивное влияние предредактирования в сравнении с

постредактированием на качество МП для пользователей с низким уровнем владения языком перевода.

Y. Liang и W. Нап в задаче перевода медицинских текстов [63] сравнили оба подхода по критериям снижения количества ошибок НМП и сохранения экономической эффективности. Результаты показали, что предредактирование повышает качество перевода и позволяет повысить экономическую эффективность по сравнению с пост-редактированием.

Gerlach, Porro, Bouillon и Lehmann [64] показали, что в прикладных задачах оба подхода могут использоваться совместно, при этом предредактирование позволяет сократить время пост-редактуры МП и в 65% случаев повышает качество результата МП.

Исследование, проведенное коллективом российских ученых в задаче перевода новостей с английского языка на русский [65], также подтвердило, что предредактирование исходного текста сокращает время пост-редактуры МП на 30-40%.

Несмотря на подтвержденную исследования эффективность предредактирования как способа повышения качества перевода, следует учитывать и тот факт, что при автоматизации эффект от предварительного редактирования может быть непредсказуемым, включая отсутствие положительного эффекта и даже преобладание отрицательного [66].

А. Агепаз в широкомасштабном анализе обоих подходов [67] показала, что способ предредактирования и методология его применения должны адаптироваться под особенности языковой пары перевода, специфику предметной области, а также используемую технологию МП. Отмечается, что необходим анализ исходного текста с целью определить «уязвимости» с точки зрения применяемой технологии МП, на основе которого становится возможной выбор стратегии предварительного редактирования.

1.3 Методы предредактирования для повышения качества машинного перевода

1.3.1 Правила для авторов по написанию исходных текстов

Один из способов повышения качества МП заключается в том, чтобы изначально написать текст «правильно», минимизировав вероятность полисемии и типичных ошибок [68, 69]. Каждый язык имеет грамматические правила. Не существует таких правил, которые давали бы подходящие для всех языков результаты. Однако существуют правила, которые снижают уровень двусмысленности в большинстве текстов на многих языках, например, «писать короткими и грамматически простыми предложениями», «использовать существительные вместо местоимений», «использовать определяющие слова», «использовать активный залог вместо пассивного».

Это позволяет создавать тексты, которые легче читать, понятнее и легче запоминать, а также с лучшим словарным запасом и стилем. В связи с этим многие компании, использующие МП, а так же разработчики систем МП разрабатывают и выпускают различные рекомендации для пользователей МП по написанию исходных текстов [70, 71]. Существенным недостатком данного подхода является то, что автор исходного текста и пользователь системы МП – не всегда один и тот же человек, поэтому зачастую авторы игнорируют эти рекомендации, так как не заинтересованы в получении перевода написанного ими текста [72].

1.3.2 Концепция контролируемого языка

Контролируемый язык — это подвид естественного языка, полученный ограничением в использовании грамматики, терминологии и речевых оборотов посредством регламентирующих правил с тем, чтобы снизить или искоренить его многозначность и сложность [73].

Традиционно контролируемые языки подразделяются на две группы: в одной все усилия направлены на повышение удобочитаемости для человека (например, для тех, кому язык текста не родной) [74]; в другой эти меры направлены на создание языка, надежного в плане автоматического семантического анализа, в частности, для повышения эффективности систем МП [75].

Контролируемый язык особенно полезен в задачах машинного перевода, где требуется высокая точность и аккуратность, например, перевода юридических документов, медицинских текстов или технических спецификаций [76], так как позволяет получить более единообразный и стандартизированный исходный текст, что обеспечивает как более высокую частоту совпадений в памяти переводов, так и более высокое качество генерируемого МП. Однако, несмотря на преимущества концепции контролируемого языка, она сложна и трудоемка в реализации и внедрении.

1.3.3 Предредактирование, основанное на правилах

В отличии от правил для авторов и концепции контролируемого языка, предредактирование предполагает, что редактор подготавливает и оптимизирует исходный текст под задачу МП. Выделяют следующие принципы предредактирования:

- Удаление неоднозначных слов и выражений, которые могут иметь несколько значений на языке перевода.
- Трансформация сложных грамматических конструкций в более простые, к которым проще найти соответствие в памяти переводов и которые проще перевести.
- Замена терминологии на стандартизированную, расшифровка аббревиатур, что помогает обеспечить согласованность перевода и избежать ошибок.

Так, Нігаока и Yamada [77] в своей работе предприняли попытку сформулировать основные правила предредактирования текста при переводе с японского языка на английский. При этом, они разделяют предредактирование на два вида: одноязычное и двуязычное. При одноязычном предредактировании редактор знает только язык оригинала и, следуя рекомендация или «правилам» вносит правки с исходный текст. При двуязычном предредактировании редактор оценивает оригинальный текст, выполненный машинный перевод и корректирует оригинальный текст с целью добиться желаемого перевода. Их стратегия показала статистически значимые результаты, в том числе и для китайского и корейского языков. Подобные исследования проводились для перевода с японского на английский и восточные языки [78], с индонезийского на английский [79], с английского на испанский [80]. Все указанные выше исследования показали, что предредактирование и переписывание исходного текста, опираясь на правила, повышает качество МП. В качестве

недостатков такого подхода можно выделить трудоемкость, зависимость результата от профессионализма редактора.

1.3.4 Решение задачи перефразирования как способ предредактирования

Разработки в области обработки естественного языка и возросшая доступность текстовых корпусов сделали возможным развитие методов и алгоритмов автоматической адаптации (упрощения) текста, которые могут быть использованы и в задачах МП [81]. Одним из перспективных направлений в этой области является автоматическое перефразирование текста [82] при котором создаются новые версии текста с сохранением семантики исходного текста. Путем изменения фраз, предложений или структуры текста текст адаптируется текст под требования реципиента [83].

Существует несколько подходов к решению задачи автоматического перефразирования текстов. Один из них основывается на использовании правил и шаблонов для замены слов или фраз на эквивалентные им варианты. Несмотря на активное изучение проблемы, оптимальные решения пока найдены далеко не для всех задач в этой области, а существующие методы, применимы к сравнительно небольшому числу языков [84, 85, 86, 87]. Для русского языка разработан ряд правил и алгоритмов по упрощению синтаксической сложности [88] и многозначности отдельных частей речи [89].

Другой подход — это использование методов глубокого обучения, таких как рекуррентные нейронные сети, которые позволяют модели генерировать новые версии текста на основе обучающих данных [90]. Это сложная задача, которая требует понимания контекста и семантики исходного текста, а также умения создавать новые, грамматически правильные и естественно звучащие предложения [91]. Анализу русскоязычных текстов посвящены работы и автоматизации оценки восприятия текста посвящены ряд работ [92, 93], однако они не предлагают алгоритмов генерации текста на основе выявленных семантик. Кроме того, для обучения моделей перефразирования требуются корпуса параллельных текстов большого объема, которые сложно собрать, особенно в узкоспециальных и высоко конкурентных предметных областях.

Выводы по первой главе

Анализ существующих исследований по теме автоматизированной предобработки исходных текстов для МП, что является перспективным направлением исследований в соответствии с определенными ключевыми направлениями развития систем МП, описанными выше, показал, что:

- не проводилось исследований по анализу ошибок нейронного МП с русского языка, связанных со спецификой той или иной предметной области или с особенностями языка;
- несмотря на то, что для русского языка разработаны синтаксические анализаторы и возможна разработка правил автоматического предредактирования для повышения качества МП, существующие исследования носят точечный несистемный характер;
- в литературе незначительно описаны правила и алгоритмы автоматического предредактирования текстов в контексте задачи повышения качества МП, существующие прикладные исследования ориентированы в большей степени на адаптацию текстов для реципиентов различной квалификации;
- существующие подходы к автоматизации предредактирования имеют ограничения, связанные со сбором и подготовкой обучающих данных.

Обобщая вышесказанное можно заключить, что являются актуальными задачи определения зависимости качества перевода от параметров исходного текста и стратегий его обработки с целью максимизации качества машинного перевода, а также разработки методов и алгоритмов оптимизационного предредактирования текстов на русском языке согласно особенностям языков оригинала и перевода, специфичных для определенной предметной области, и технологии МП.

ГЛАВА 2 МАТЕМАТИЧЕСКАЯ МОДЕЛЬ ОПТИМИЗАЦИОННОГО ПРЕДРЕДАКТИРОВАНИЯ

2.1 Обобщенная математическая модель процесса перевода

Решение задачи разработки моделей, методов и комплекса программ оптимизационного предредактирования текстов для дальнейшего перевода, лежит в области прикладной компьютерной лингвистики.

Задача компьютерной лингвистики (КЛ) – разработка методов и средств построения лингвистических процессоров для прикладных задач по автоматической обработке текстов на естественном языке (ЕЯ) [94]. Разработка лингвистического процессора для некоторой прикладной задачи предполагает формальное описание лингвистических свойств обрабатываемого текста, которое может рассматриваться как модель текста (или модель языка).

Машинный перевод – самое раннее приложение КЛ, вместе с которым возникла и развивалась сама эта область. Первые программы перевода были построены на простейшей стратегии пословного перевода, однако довольно быстро стало ясно, что машинный перевод требует гораздо более полной лингвистической модели.

Компьютерная лингвистика тесно связана с областью искусственного интеллекта (ИИ) [95], в рамках которой разрабатываются программные модели отдельных интеллектуальных функций. Несмотря на очевидное пересечение исследований в области компьютерной лингвистики и ИИ (поскольку владение языком относится к интеллектуальным функциям), ИИ не включает в себя всю КЛ, поскольку она имеет свой теоретический базис и методологию. Общим для указанных наук является компьютерное моделирование как основной способ и итоговая цель исследований, эвристический характер многих применяемых методов.

Машинное обучение в КЛ применяется для обработки коллекций текстовых документов, с использованием признаковой модели текста, при которой признаки определены для каждого документа по отдельности. Признаками могут выступать различные информационные характеристики текста: как лингвистические, так статистические и структурные: например, частота определенных слов (или их

категорий) в документе, частота использования спецзнаков, соотношение частей речи слов, наличие определенных синтаксических конструкций или разделов текста, дата создания и прочие [94].

Для постановки и формализации задачи настоящего исследования, проанализируем процесса перевода текста с одного естественного языка на другой с декомпозицией на этапы и построим обобщенную математическую модель [96]. В качестве основной методологии описания математической модели используем теорию множеств.

2.1.1 Термины и определения

Определения элементов модели сформулированы на основании анализа предметной области и фундаментальных исследований в области лингвистики [97].

Язык – сложная знаковая система, естественно или искусственно созданная и соотносящая понятийное содержание и типовое звучание (написание).

Множество существующих языков можно представить следующим образом:

$$g_3 \in \mathcal{H}_3: \mathcal{H}_3 = \{g_3, g_3, \dots g_{ng_3}\}.$$
 (3)

Понятийное содержание определяется семантическими единицами.

Семантическая единица — это абстрактная сущность, не привязанная к конкретному языку или иному средству выражения (например, языку математической логики), которая обеспечивает идентификацию определенных объектов, явлений, свойств домена приложения.

Пример: Семантическая единица – дождь. Способы выражения: дождь; атмосферные осадки в виде водяных капель; природное явление, когда с неба падают капли воды; rain (англ.); regen (нем.); и т.д.

Множество семантических единиц можно представить следующим образом:

$$ce \in CE: CE = \{ce_0, ce_1, \dots ce_{nCE}\}.$$
 (4)

Пусть множество U_{ce} — универсальное множество всех возможных описаний на всех возможных языках всех возможных семантических единиц. Ice — множество индексов-количество семантических единиц:

$$ice \in Ice: Ice = \{0, 1, ..., nCE\}.$$
 (5)

При этом множества CE и Ice являются биективными, т.е. $Ice \leftrightarrow CE$.

Каждому элементу $ice \in Ice$ однозначно сопоставлено подмножество описаний семантической единицы $O_{ice} \subseteq U_{ce}$. Тогда $CoCE = (O_{ice})_{ice \in Ice}$ — счетное семейство множеств описаний семантических единиц.

$$CoCE = (O_{ice})_{ice \in Ice} = \bigcup_{ice \in Ice} O_{ice} = \{x: \exists ice \ x \in O_{ice}\}$$
 (6)

Домен приложения или **предметная область** – абстрактное понятие, которое определяет множество объектов в пределах общего контекста. Множество доменов (предметных областей) приложения:

$$\partial n \in \Pi : \Pi = \{\partial n_0, \partial n_1, \dots \partial n_{n\Pi\Pi}\}. \tag{7}$$

Один из способов выражения семантических единиц и коммуникации в рамках некоторого домена приложения – текст.

Текст – это письменное сообщение, объективированное в виде письменного документа, состоящее из ряда высказываний, объединённых разными типами лексической, грамматической и логической связи, имеющее определённый моральный характер, прагматическую установку и соответственно литературно обработанное. Текст любого размера – это относительно автономное (законченное) высказывание. Множество текстов:

$$txt \in TXT: TXT = \{txt_0, txt_1, \dots txt_{nTXT}\}. \tag{8}$$

Уровни фрагментации текста: Символ \to Морфема (корень, суффикс и т.д.) \to Лемма \to Словоформа \to Лексема \to Словосочетание \to N-грамма \to Предложение \to Абзац \to Раздел \to Документ (книга) \to Корпус.

Текст как объект анализа и исследования обладает некоторым набором свойств или признаков. Множество признаков или свойств текста:

$$ce \in CB: CB = \{ce_0, ce_1, \dots ce_{nCB}\}. \tag{9}$$

В задачах обработки естественного языка свойства текста условно можно разделить на группы признаков: *общие* (количество символов/слов/строк и т.д., стиль, язык, домен приложения и пр.), $O\Pi \subset CB$; лексические (процент покрытия текста лексическими минимумами, частотными списками и др.), $J\Pi \subset CB$; морфологические (количество различных частей речи и грамматических форм), $M\Pi \subset CB$; синтаксические (глубина глагольных и именных групп, связи между глаголами в

предложениях), $C\Pi \subset CB$; признаки, основанные на базовых подсчетах (средняя длина слов и предложений и пр.), $B\Pi \subset CB$.

$$CB = O\Pi \cup J\Pi \cup M\Pi \cup C\Pi \cup B\Pi \tag{10}$$

Совокупность свойств и признаков определяет главные параметры текста, к которым относится целостность, связность, трудность, удобочитаемость, сложность и другие. Множество главных параметров текста:

$$\operatorname{en} \in \Gamma\Pi: \Gamma\Pi = \{\operatorname{en}_0, \operatorname{en}_1, \dots \operatorname{en}_{n\Gamma\Pi}\}. \tag{11}$$

Основные параметры текста – удобочитаемость и сложность.

Удобочитаемость позволяет оценить, соответствует ли текст читательской способности реципиента, а на уровне отдельных его единиц (например, предложений) выявить элементы, требующие упрощения (специальные термины, аббревиатуры, сокращения, отсутствие контекста (таблицы, счета, картинки и т. д.), слишком короткие или слишком длинные предложения, ошибки и пр.). Для оценки удобочитаемости разработаны различные метрики: индекс Флеша, школьный тест Флеша-Кинкейда, FOG, SMOG, Индекс Коулман-Лиау, Автоматический индекс удобочитаемости [98].

Сложность текста определяется морфологическими, лексическими, синтаксическими признаками текста (объективные оценки), а также его информационной (когнитивной) трудностью (субъективная оценка). Когнитивная трудность текста определяется на уровне реципиента его способностью идентифицировать семантические единицы в тексте [99].

Перевод — деятельность по интерпретации понятийного содержания текста на одном языке и созданию нового текста на другом языке эквивалентного понятийного содержания.

Перевод выполняется переводчиком из множества переводчиков, компетенция каждого из которых определяется набором знаний, умений и навыков в области владения языками. Множество переводчиков:

$$nep \in \Pi EP$$
: $\Pi EP = \{nep_0, nep_1, \dots nep_{n\Pi EP}\}.$ (12)

Знания языков или языковая компетенция — это владение грамматическими и словарными аспектами языка. Знания о языках можно охарактеризовать

совокупностью множеств лексем, грамматических правил, морфологических правил, правил синтаксиса и пунктуации, словарей частотности словоупотребления.

Лексема – абстрактная двусторонняя единица словарного состава языка в совокупности всех ее конкретных грамматических форм и выражающих их флексий, а также всех возможных значений во всех его употреблениях и реализациях.

Множество лексем определяется множеством лемм:

$$\Lambda M \in \mathcal{I}M: \mathcal{I}M = \{\Lambda M_0, \Lambda M_1, \dots \Lambda M_{n, \mathcal{I}M}\}, \tag{13}$$

где $лм_{iЛM}$ начальная словарная форма слова или лемма. В русском языке для существительных и прилагательных это форма именительного падежа единственного числа, для глаголов и глагольных форм — форма инфинитива.

Пусть множество $U_{c\phi\pi}$ — универсальное множество всех возможных словоформ всех возможных лемм на всех возможных языках, а $I_{\pi}M$ — множество индексов-количество лемм:

$$i$$
лм $\in I$ лм: I лм = $\{0, 1, ..., n$ Л $M\}$, I лм \leftrightarrow Л M (14)

Каждому элементу $iлм \in Iлм$ однозначно сопоставлено подмножество лексем $\mathit{ЛКC}_{iлм} \subseteq U_{\mathrm{сфл}}$. Тогда $\mathit{Cm}\mathit{ЛКC} = (\mathit{ЛКC}_{i\mathit{лм}})_{i\mathit{лм}} \in \mathit{I}_{\mathit{лм}}$ — счетное семейство множеств лексем, т.е.

$$C$$
м J K $C = (Л K C _{i Лм}) $_{i$ Лм $\in I$ Лм = \bigcup_{i Лм $\in I$ Лм J K C $_{i$ Лм = $\{x: \exists i$ Лм $x \in J$ K C $_{i$ Лм} $\}$. (15)$

Количество словоупотреблений одной леммы во всех вариантах ее лексем в некотором тексте (или в корпусе текстов, или в речевом фрагменте) называется ее частотностью, измеряется в лексической статистике и описывается в **частотных словарях**. Пусть множество $U_{\text{част}}$ — универсальное множество всех частотностей всех возможных лемм и сочетаний слов на всех возможных языках для всех доменов приложения. Каждому элементу $iлм \in Iлм$ однозначно сопоставлено подмножество частотностей $VACT_{iлм} \subseteq U_{\text{част}}$. Тогда $CMVACT = (VACT_{inm})_{inm} \in I_{nm}$ — счетное семейство множеств частотностей употребления лемм.

$$C_{M}YACT = (YACT_{i_{\mathcal{I}M}})_{i_{\mathcal{I}M} \in I_{\mathcal{I}M}} = U_{i_{\mathcal{I}M} \in I_{\mathcal{I}M}} YACT_{i_{\mathcal{I}M}} = \{x: \exists i_{\mathcal{I}M} \ x \in YACT_{i_{\mathcal{I}M}}\}$$
(16)

Интерпретация и перевод слов на другие языки содержатся в словарях. Словарь — это собрание слов, устойчивых выражений с пояснениями, толкованиями и/или с переводом на другой язык, словарь состоит из словарных статей.

Словарная статья — это кортеж сочетаний слов (словоформ отдельных лемм) длинной k и n (k, $n \le 4$) на разных языках, таких что описывают максимально близко друг другу одну семантическую единицу, то есть

$$CC = \{(cou. cлoв 1; cou. cлoв 2): ce_{iCE} \mid cou. cлoв 1 \rightarrow ce_{iCE} \mid cou. cлoв 2\},$$
 (17)

$$cou. cnoe 1 = (cn H31_{0c\phi}, cn H31_{ic\phi}, ..., cn H31_{mc\phi}),$$
 (18)

$$cou. cnoe 2 = (cn H320c\phi, cn H32jc\phi, ..., cn H32nc\phi),$$
 (19)

$$cлЯ31_{ic\phi}, cлЯ32_{jc\phi} \in U_{c\phiл},$$
 (20)

$$ce_{iCE}, ce_{iCE} \in CE.$$
 (21)

Пусть множество U_{cc} – универсальное множество всех возможных словарных статей. Каждому элементу множества индексов количества семантических единиц $ice \in Ice$ однозначно сопоставлено множество словарных статей $C\Pi_{ice} \subseteq U_{cc}$, то есть $CMC\Pi = (C\Pi_{ice})_{ice \in Ice}$ – счетное семейство множеств словарей.

$$C_{\mathcal{M}}C_{\mathcal{I}} = (C_{\mathcal{I}_{ice}})_{ice \in I_{ce}} = U_{ice \in I_{ce}}C_{\mathcal{I}_{ice}} = \{x: \exists ice \ x \in C_{\mathcal{I}_{ice}}\}$$
 (22)

Правила употребления слов можно охарактеризовать набором следующих множеств:

• множество грамматических правил:

$$\operatorname{en} \in \Gamma\Pi : \Gamma\Pi = \{\operatorname{en}_0, \operatorname{en}_1, \dots \operatorname{en}_{n\Gamma\Pi}\}; \tag{23}$$

• множество морфологических правил:

$$mop \in MOP: MOP = \{mop_0, mop_1, ..., mop_{nMOP}\};$$

$$(24)$$

• множество правил синтаксиса и пунктуации:

$$cuh \in CUH : cuh = \{cuh_0, cuh_1, ... cuh_{nCUH}\}.$$
 (25)

Умения переводчика — это способности переводчика, которые позволяют ему интерпретировать смысловую единицу, выраженную в тексте на языке оригинала, и создать эквивалентный по смыслу текст на языке перевода в соответствии с предъявляемыми к переводу требованиями. Множество умений:

$$y_M \in Y_M: Y_M = \{y_{M_0}, y_{M_1}, \dots y_{M_{nV_M}}\},$$
 (26)

где y_{MiyM} — функция обработки, преобразования или генерации текста таким образом, чтобы созданный текст максимально точно передавал смысл исходного текста, то есть

$$y_{M_{iYM}}(txt_{iTXT}) = txt'_{iTXT} : CM|txt'_{iTXT} \to CM|txt_{iTXT},$$
(27)

где CM – **смысл** или упорядоченный набор семантических единиц, описываемый текстом:

$$CM = \{(ce_0, ce_1, ..., ce_{ncm}): ce_{icm} \in CE \}.$$
 (28)

Множество всех смыслов:

$$CM \in \mathcal{M}CM: \mathcal{M}CM = \{CM_0, CM_1, \dots CM_{nCM}\}. \tag{29}$$

Умения связаны со способностью выполнять операции над текстом согласно требованиям качества. Множество требований к качеству перевода можно описать следующим образом:

$$mp \in TP: TP = \{mp0, mp1, \dots mpnTP\}, \tag{30}$$

где mp_{iTP} — это некоторое лингвистическое требование, которому должен соответствовать конечный результат перевода, может быть записано в виде условия либо утверждения.

Требования определяют: нормы языка перевода $H\mathcal{A} \subset TP$; культурные традиции реципиента $KTP \subset TP$; отраслевые стандарты согласно домену приложения $OC|\partial n_{iJII} \subset TP$; заказчик и/или реципиент перевода $3P \subset TP$.

Множество *TP* включает и все прочие требования, не вошедшие в другие подмножества требований. Подмножества требований могут пересекаться между собой. Требования к переводу от заказчика и/или реципиента — это множество сочетаний слов, используемых для выражения семантических единиц в текстах заказчика. При наличии требований от заказчика, частотность сочетаний слов, описанных в таких требованиях, при переводе текстов заказчика принимается как максимальная.

Навыки переводчика — это опыт использования языка в решении задач перевода в различных предметных областях, полученный в результате обучения либо практической деятельности. Навыки характеризуются накоплением опыта через следующие понятия.

Множество устойчивых **n-грамм**, т.е. наборов слов неограниченной длины, которые описывают некоторые смыслы (наборы семантических единиц). Взаимосвязи внутри наборов семантических единиц и их составы, а также n-граммы описывающие составленные наборы, которые обобщенно можно назвать текстами $txt \in TXT$, для различных доменов приложения определяются / дополняются / корректируются переводчиком в процессе обучения и/или практической деятельности, то есть

$$ycou \in YCOU: YCOU = \{ycou_0, ycou_1, \dots ycou_{nYCOU}\},$$
(31)

где

$$ycou_{nVCOY} = \{(txt, CM): txt \in TXT, CM \in MCM\}.$$
(32)

Семейство множеств **частотности употребления устойчивых n-грамм** для выражения смыслов вообще и/или в рамках определенного домена. Пусть *Iусоч* — множество индексов-количество устойчивых n-грамм:

$$iycou \in Iycou: Iycou = \{0, 1, ..., nVCOU\}, Iycou \leftrightarrow VCOU\}.$$
 (33)

Каждому элементу $iycou \in Iycou$ однозначно сопоставлено подмножество частотностей $VCYACT_{iycou} \subset U_{uacm}$. Тогда $C_{m}VCYACT = (VCYACT_{iycou})_{iycou \in Iycou}$ — счетное семейство множеств частотностей употребления устойчивых n-грамм.

$$C_{M}YCYACT = (YCYACT_{iycov})_{iycov} =$$

$$= \bigcup_{iycov} YACT_{iycov} = \{x: \exists iycov \ x \in YCYACT_{iycov}\}$$

$$(34)$$

Специализация переводчика $nep_{i\Pi EP}$ — это его накопленный опыт выражения семантических единиц в рамках некоторой предметной области $\partial n_{iД\Pi}$, включая обучение. Математически определим специализацию переводчика в рамках решения прикладной задачи как вектор, координатами которого являются условные подмножества:

$$\overrightarrow{Cnep_{i\Pi EP}, \partial n_{i\Pi\Pi}} = \{(OT|\partial n_{i\Pi\Pi})|nep_{i\Pi EP}, (YCOY|\partial n_{i\Pi\Pi})|nep_{i\Pi EP},$$

$$(CMYCYACT|\partial n_{i\Pi\Pi})|nep_{i\Pi EP}\}$$
(35)

Компетенция некоторого переводчика $nep_{i\Pi EP}$ оценивается всегда в рамках прикладной задачи, исходными данными которой являются язык оригинала $я3_{ex} \in Я3$ и язык перевода $я3_{выx} \in Я3$. Математически определим компетенцию переводчика в рамках решения прикладной задачи как вектор, координатами которого являются условные подмножества:

$$\overline{Knep_{i\Pi EP}} = \{ (CoCE \mid яз_{ex} \cup CoCE \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(CMЛКС \mid яз_{ex} \cup CMЛКС \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(CMЧАСТ \mid яз_{ex} \cup CMЧАСТ \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(CMСЛ \mid яз_{ex} \cup CMСЛ \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(\Gamma\Pi \mid яз_{ex} \cup \Gamma\Pi \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(MOP \mid яз_{ex} \cup MOP \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(CИН \mid яз_{ex} \cup CИН \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(VM \mid яз_{ex} \cup VM \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(VCOY \mid яз_{ex} \cup VCOY \mid яз_{ebix}) \mid nep_{i\Pi EP},$$

$$(CMVCYACT \mid яз_{ex} \cup CMVCYACT \mid яз_{ebix}) \mid nep_{i\Pi EP} \}$$

Рассмотрим применение описанных элементов модели в процессе перевода.

2.1.2 Оценка задачи перевода

Переводчик $nep_{i\Pi EP}$ получает текст txt_{iTXT} на языке $яз_{6x}$ для перевода на язык $яз_{6bix}$ и требования к переводу от заказчика $3P|txt_{iTXT}$ (если есть), при этом $3P|txt_{iTXT}$ $\subset TP|txt_{iTXT}$.

Перед выполнением перевода переводчик оценивает сложность задачи перевода. При оценке сложности задачи перевода переводчик обращает внимание на неизвестные ему слова и сочетания слов на языке g_{ax} , для которых он не может идентифицировать значение смысловой единицы, либо смысловые единицы, для которых он не может найти аналог на языке перевода g_{abx} среди известных ему слов и сочетаний слов. Множества свойств и параметров исходного текста G_{abx} и G_{abx} и то, обладает ли переводчик достаточной компетенцией G_{abx} относительно языков g_{abx} и специализацией G_{abx} и специализацией G_{abx} и специализацией G_{abx} и специализацией G_{abx} и переводчиком переведенного текста, определяет вероятность создания переводчиком переведенного текста на таком уровне качества, который определяется требованиями G_{abx} и G_{abx} и определяется требованиями G_{abx} и определяется требования G_{abx} и определяется G_{abx} и определяется G_{abx} и определяется G_{abx} и определяется G

Алгоритм оценки сложности задачи перевода следующий:

Шаг 1. Исходя из домена приложения текста $\partial n_{iД\Pi}$, формируется множество оценок текста $OU = CB \ U \Gamma \Pi$.

Шаг 2. Для каждого значения c_{BiCB} , $c_{niFII} \in OU$, на основе требований к переводу $TP|txt_{iTXT}$, компетенций переводчика относительно языковой пары $\overline{Knep_{i\Pi EP}}$ и специализации переводчика относительно домена приложения текста $\overline{Cnep_{i\Pi EP}, \partial n_{iZII}}$ формируется значение значимости w_{ouk} , множество нормированных значений w_{ouk} значимости формируют матрицу значимости оценок сложности W_{ouk} размерностью $1 \times k$, где k — общее число оценок, которые выступают коэффициентами уравнения поиска теоретического значения качества перевода.

Шаг 3. Для каждого *i*-го фрагмента текста при $i=\overline{I,N}$ формируется матрица оценок фрагмента исходного текста C_{oui} размерностью $1\times k$, где k – общее число оценок.

Шаг 4. На основании оценок C_{oui} значимости W_{ou} формируется уравнение поиска теоретического результирующего фактора, т.е. качества перевода \widehat{OK} :

$$\widehat{OK}_{i}(C_{oui}, W_{ou}) = w_{0} + w_{1ou}C_{oui} + w_{2ou}C_{oui} + \dots + w_{kou}C_{ouik}$$
(37)

Шаг 5. Для каждого *i-го* фрагмента текста рассчитаем вероятность получения переведенного текста на таком уровне качества, который определяется требованиями $TP|txt_{iTXT}$, применив к уравнению (37) логит-преобразование [100]:

$$p_i = \frac{1}{1 + e^{-\widehat{OK}_i}} \tag{38}$$

Шаг 6. Сложность задачи перевода *i-го* фрагмента текста:

$$C_{\pi}3\Pi_{i} = \frac{1}{p_{i}},\tag{39}$$

где p_i – это вероятность создания переводчиком перевода требуемого качества, рассчитанная по формуле (38).

Шаг 7. Результирующая сложность задачи перевода текста — это наибольшее значение сложности задачи перевода $Cn3\Pi_i$ среди N фрагментов исходного текста:

$$C_{\pi}3\Pi_{txtiTXT} = max C_{\pi}3\Pi_{i}$$
 (40)

В зависимости от значения $Cn3\Pi_{txt_{iTXT}}$ переводчик выбирает стратегию дальнейшей обработки текста. Возможные варианты для ручного и машинного перевода приведены в таблице 2.1.

Сложность	Действие	
задачи перевода	Ручной перевод	Машинный перевод
Низкая	Перевод	Перевод
Средняя	Предварительная редактура	
	текста / согласование требо-	Автоматическая предредак-
	ваний к переводу	тирование текста
	с заказчиком	
Высокая	Отказ от перевода / отправка	Полуавтоматическая пред-
	на доработку текста	редактирование текста с
	заказчику	привлечением пользователя

Таблица 2.1 – Действия переводчика в зависимости от сложности задачи перевода

Диапазоны значений $Cл3\Pi$, соответствующие низкому, среднему и высокому уровням определяются на основе предварительного анализа требований к переводу и с учетом способа перевода (ручной/машинный).

Формально функцию оценивания сложности переводческой задачи F_{ou} можно представить следующим образом:

$$F_{ou}: txt_{iTXT} \to C_{\pi}3\Pi|(TP|txt_{iTXT}, CB, \Gamma\Pi, \overrightarrow{Knep_{i\Pi EP}}, \overrightarrow{Cnep_{i\Pi EP}}, \overrightarrow{on_{i\Pi EP}})$$
 (41)

Схематично, функция оценки сложности переводческой задачи представлена на рисунке 2.1.

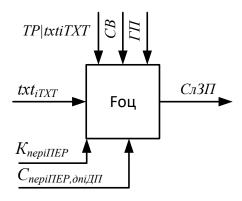


Рисунок 2.1 – Оценка сложности задачи перевода

2.1.3 Предредактирование исходного текста

Предредактирование — это процесс автоматического/полуавтоматического изменения свойств/параметров исходного текста с целью оптимизации сложности переводческой задачи с учетом влияния такого изменения на качество перевода.

Предредактирование выполняется при средней и высокой сложности переводческой задачи с использованием разнообразных алгоритмов, например, по критериям: минимизации $Cn3\Pi$, минимаксимизации отдельных свойств и параметров исходного текста и др.

Цель предредактирования — снизить сложность переводческой задачи $Cл3\Pi_{txtiTXT}$ до низкой. На начальном этапе предредактирования происходит отбор фрагментов текста, которым соответствуют высокие значения $Cn3\Pi_i$, то есть

$$txt_{i}^{*} \in txt_{iTXT} \mid Cn3\Pi_{i} > Cn3\Pi_{\partial on}$$

$$\tag{42}$$

где $Cл3\Pi_{\partial on}$ – допустимое значение сложности задачи перевода, т.е. низкое.

Для тех фрагментов текста, для которых сложность задачи перевода является средней или высокой, определяем элементы вектора C_{oui} , значения которых повышают сложность переводческой задачи данного фрагмента: $C_{oui}|C_{n}3\Pi_{i}>C_{n}3\Pi_{gon}$.

Пусть множество U_{AOT} – универсальное множество всех возможных алгоритмов обработки текста. Iou – множество индексов-количество оценок свойств и параметров текстов. При этом множества OU и Iou являются биективными, т.е. $Iou \leftrightarrow OU$. Каждому элементу $iou \in Iou$ однозначно сопоставлено подмножество алгоритмов оптимизации $AO_{iou} \subseteq U_{AOT}$. Тогда $CMAO = (AO_{iou})_{iou} \in Iou$ — счетное семейство множеств алгоритмов предредактирования с целью оптимизации оценок текста.

$$C_{MAO} = (AO_{iou})_{iou \in Iou} = \bigcup_{ice \in Ice} AO_{iou} = \{x: \exists iou \ x \in AO_{iou}\}$$
 (43)

Для каждого найденного элемента матрицы C_{oui} , соответствующего условию $C_{oui}|C_n 3\Pi_i > C_n 3\Pi_{\partial on}$, в зависимости от его значения и ограничений алгоритмов предредактирования определяется стратегия предредактирования (таблица 2.2). Выбор алгоритма зависит от критерия оптимизации выбранного значения.

Таблица 2.2 – Стратегии предредактирования

Значение параметра	Действие
В допустимых	Автоматическое предредактирование
пределах	при помощи алгоритмов
Выходит за	Полуавтоматическое предредактирование текста
допустимые пределы	с привлечением пользователя

Далее производим редактирование текста в соответствии с имеющимися методами и алгоритмами предредактирования, получаем текст txt'_{iTXT} и переходим к этапу перевода. Формально функцию предредактирования $F_{npedped}$ можно представить следующим образом:

$$F_{nedped}$$
: txt_{iTXT} , $C_{n3}\Pi$, C_{ou} , $W_{ou} \rightarrow txt'_{iTXT} \mid (C_{n3}\Pi \rightarrow min, C_{MAO}, \overline{Knep_{i\Pi EP}}, \overline{C_{nep_{i\Pi EP}}, \partial n_{i\Pi I}})$. (44)

Схематично, функция предредактирования исходного текста представлена на рисунке 2.2.

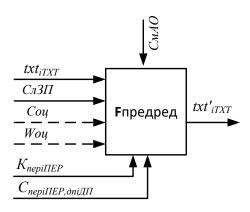


Рисунок 2.2 – Предредактирование исходного текста

2.1.4 Перевод

Этап перевода включает следующие шаги:

Далее для каждого і-го отрезка текста при $i = \overline{I,N}$:

Шаг 2. Анализ смысла исходного текста на языке g_{36x} , на основании компетенции переводчика относительно языка g_{36x} ($\overline{Knep_{i\Pi EP}} \mid g_{36x}$) и специализации переводчика в предметной области исходного текста $g_{nij\Pi}$ ($g_{nij\Pi}$) с целью определить набор семантических единиц, которые он описывает, то есть его смысл:

$$f_a: (txt'_{iTXT}, \overrightarrow{Knep_{i\Pi EP}} \mid \mathfrak{A3}_{ex}, \overrightarrow{Cnep_{i\Pi EP}, \partial n_{iJ\Pi}}) \to CM_i.$$
 (45)

Шаг 3. Генерация текста, который должен максимально воссоздавать смысл исходного текста CM_i согласно требованиям $TP|txt_{iTXT}$, но уже средствами языка $я_{36bix}$ на основании компетенции переводчика относительно языка $я_{36bix}$ ($\overrightarrow{Knep_{i\Pi EP}} \mid я_{36bix}$) и специализации переводчика в предметной области исходного текста $\partial n_{iJ\Pi}$ ($\overrightarrow{Cnep_{i\Pi EP}}, \partial n_{iJ\Pi}$):

$$f_g: (CM_i, \overrightarrow{Knep_{i\Pi EP}} \mid яз_{вых}, \overrightarrow{Cnep_{i\Pi EP}, \partial n_{i\Pi \Pi}}) \rightarrow txt'_{iTXT} \mid (CM_j \rightarrow CM_i, TP \mid txt_{iTXT}).$$
 (46)

Формально обобщенную функцию перевода F_{nep} можно представить следующим образом:

$$F_{nep}$$
: $(txt'_{iTXT}, \overline{Knep_{i\Pi EP}}, \overline{Cnep_{i\Pi EP}, \partial n_{i\Pi I}}) \rightarrow txt_{jTXT} | (CM_j \rightarrow CM_i, TP | txt_{iTXT}).$ (47)
Схематично, функция перевода представлена на рисунке 2.3.

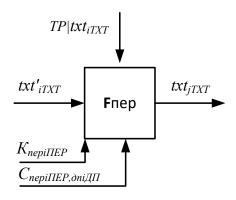


Рисунок 2.3 – Перевод текста

2.1.5 Оценка качества перевода

Проверка качества перевода — комплексная задача оценки оптимальности передачи смысла исходного текста txt_{iTXT} на языке startion startion startion startion in the startion of the star

$$\kappa\kappa \in KK: KK = \{\kappa\kappa_0, \kappa\kappa_1, ..., \kappa\kappa_{nKK}\}. \tag{48}$$

На основании требований к переводу $TP|txt_{iTXT}$ и применимости в рамках домена приложения $\partial n_{iД\Pi}$ формируется множество применимых критериев качества $\Pi KK = KK|TP|txt_{iTXT}$, $\Pi KK \subseteq KK$ и множество значений $w_{in\kappa\kappa}$ значимости для каждого из применимых критериев, которые формируют нормированную матрицу значимости критериев качества $\overline{W_{\Pi KK}}$ размерностью $1 \times M$, где M – общее число применимых критериев.

Проверка выполняется контролером качества перевода, в качестве которого выступает некоторый $nep_{j\Pi EP}$, обладающий соответствующей компетенцией $\overline{Knep_{j\Pi EP}}$ и специализацией $\overline{Cnep_{j\Pi EP}}$, $\partial n_{iД\Pi}$. Переводчик, переводивший текст может выступать контролером качества перевода, в случае если выполняется самопроверка перевода, тогда считается, что $nep_{j\Pi EP} = nep_{i\Pi EP}$, но следует учитывать, что в таком случае на оценку качества перевода будут влиять значения компетенции $\overline{K'nep_{i\Pi EP}}$ и специализации $\overline{C'nep_{i\Pi EP}}$, $\partial n_{iД\Pi}$, так как в процессе выполнения перевода происходит дообучение переводчика и значения его компетенции и специализации изменяются.

Далее для каждой i- \check{u} пары однозначно сопоставленных друг другу фрагментов текстов txt_{iTXT} и txt_{jTXT} при $i=\overline{I,N}$ производится оценка оптимальности согласно каждому k-m применимому критерию качества при $k=\overline{I,M}$ на основании CM_i и свойств и параметров текста $CB|txt_{jTXT}$ и $\Gamma\Pi|txt_{jTXT}$. Точность определения значения оптимальности зависит от компетенции и специализации контролера качества. Оценка оптимальности перевода по k-m0 критерию i-i0 фрагмента текста:

$$f_{koo}: ((CM_i)_i, (CB|txt_{jTXT}, \Gamma\Pi|txt_{jTXT})_i) \to OO_{i,k}| (\Pi KK_k, \overline{Knep_{j\Pi EP}}, \overline{Cnep_{j\Pi EP}, \partial n_{i,\Pi}})$$
 (49)

Оценка качества перевода i-го фрагмента текста — это определитель матрицы, полученной в результате перемножения матрицы оценок оптимальности фрагмента переведенного текста OO_i по критериям ΠKK и транспонированной матрицы значимости применимых критериев качества $\overline{W_{\Pi KK}^T}$:

$$OK_i = \left| OO_i \times \overline{W_{IIKK}^T} \right| \tag{50}$$

Результирующая оценка качества перевода текста txt_{iTXT} на язык $яз_{6ыx}$ — отношение суммы оценок качества перевода отдельных фрагментов текста к их количеству, то есть

$$OK = \frac{\sum_{N}^{i} OK_{i}}{N}$$
 (51)

Оценка качества перевода позволяет оценить эффективность работы переводчика и определить стратегию обработки переведенного текста, в случае если *ОК* вне допустимых значений. Возможные варианты постобработки текста приведены в таблице 2.3.

Таблица 2.3 – Действия в зависимости от оценки качества перевода

Оценка качества	Действие	
Крайне низкая	Отказ от перевода, перевод другим переводчиком	
Низкая	Постредактирование текста с привлечением эксперта	
Средняя	Самостоятельное / автоматическое постредактирование текста	
Высокая	——————————————————————————————————————	

Диапазоны значений OK, соответствующие крайне низкому, низкому, среднему и высокому уровням определяются на основе анализа применимых критериев качества перевода, их важности, а также целевых значений согласно требованиям к переводу $TP|txt_{iTXT}$.

Формально функцию оценивания качества перевода F_{kk} можно представить следующим образом:

$$F_{kk}: (txt_{iTXT}, txt_{jTXT}) \to OK|(TP|txt_{iTXT}, \overline{Knep_{j\Pi EP}}, \overline{Cnep_{j\Pi EP}, \partial n_{iД\Pi}}).$$
 (52)

Схематично, функция оценки качества перевода представлена на рисунке 2.4.

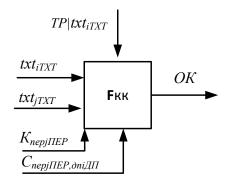


Рисунок 2.4 – Оценка качества перевода

2.1.6 Постредактирование переведенного текста

Постредактирование — это процесс автоматического/полуавтоматического изменения свойств/параметров переведенного текста с целью оптимизации оценки качества перевода. Постредактирование выполняется, если оценка качества выполненного перевода ниже допустимого либо целевого значения OK_{TARGET} .

Цель постредактирования — оптимизировать оценку качества перевода по критерию максимизации, то есть

$$OK \rightarrow max || OK_{TARGET}.$$
 (53)

Для тех фрагментов текста, для которых оценка качества OK_i ниже высокой либо целевой OK, определяем элементы матрицы OO_i , значения которых понижают качество перевода данного фрагмента:

$$OO_i|OK_i < max||OK_{TARGET}.$$
 (54)

Пусть множество U_{AOT} – универсальное множество всех возможных алгоритмов обработки текста. I_{KK} — множество индексов-количество критериев качества. При этом множества KK и I_{KK} являются биективными, т.е. $I_{KK} \leftrightarrow KK$. Каждому элементу $i_{KK} \in I_{KK}$ однозначно сопоставлено подмножество алгоритмов оптимизации свойств и параметров текста по соответствующему критерию KK: $AOK_{i_{KK}} \subseteq U_{AOT}$. Тогда $C_{MAOK} = (AOK_{i_{KK}})_{i_{KK}} \in I_{KK}$ — счетное семейство множеств алгоритмов постредактирования с целью оптимизации свойств и параметров текста по критериям качества.

$$C_{M}AOK = (AOK_{i\kappa\kappa})_{i\kappa\kappa} \in I_{\kappa\kappa} = \bigcup_{i\kappa\kappa} AOK_{i\kappa\kappa} = \{x: \exists i\kappa\kappa \ x \in AOK_{i\kappa\kappa}\}$$
 (55)

Для каждого найденного элемента матрицы OO_i , соответствующего условию (54), в зависимости от его значения и ограничений алгоритмов постредактирования определяется стратегия постредактирования (таблица 2.4). Выбор алгоритма зависит от критерия оптимизации выбранного значения.

Таблица 2.4 – Стратегии постредактирования

Значение параметра	Действие
В допустимых	Автоматическое постредактирование
пределах	при помощи алгоритмов
Выходит за	Полуавтоматическое постредактирование текста
допустимые пределы	с привлечением пользователя

Далее каждый найденный фрагмент редактируется при помощи алгоритмов постредактирования до тех пор, пока не будут достигнуты значения OO_i , при которых выполняется условие (54). Повторная оценка OO_i осуществляется по формуле (49). На выходе получаем текст txt_{umeTXT} . Эффективность применения алгоритмов постредактирования зависит от специализации и компетенции переводчика, осуществляющего постредактирование $\overline{Knep_{j\Pi EP}}$, $\overline{Cnep_{j\Pi EP}}$

Формально функцию постредактирования $F_{nocmped}$ можно представить следующим образом:

$$F_{nocmped}$$
: $(txt_{iTXT}, txt_{jTXT}, OK, OO, \overline{W_{IIKK}}) \rightarrow txt_{umeTXT} |$

$$(OK \rightarrow max | | OK_{TARGET}, CMAOK, \overline{Knep_{i\Pi EP}}, \overline{Cnep_{i\Pi EP}, \partial n_{iД\Pi}}). \tag{56}$$

Схематично, функция постредактирования представлена на рисунке 2.5.

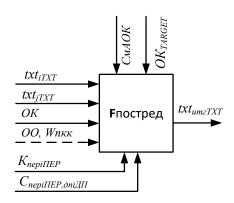


Рисунок 2.5 – Постредактирование переведенного текста

2.1.7 Обучение переводчика

Обучение переводчика – процесс формирования и/или обновления компетенций и специализаций переводчика на основе изучения обучающих текстов и/или полученного практического опыта перевода. Важно отметить, что процесс дообучения происходит непрерывно в процессе работы.

Каждый переводчик обладает множеством компетенций и специализаций в зависимости от известных переводчику языков $\mathcal{A}3|nep_{i\Pi EP}$ и доменов приложения $\mathcal{A}\Pi|nep_{i\Pi EP}$. Множество компетенций переводчика:

$$\kappa n \in K: K = \{\kappa n_0, \, \kappa n_1, \dots \, \kappa n_{nK}\}, \tag{57}$$

где $\kappa n_{iK\Pi}$ определяется по формуле (36) для некоторой языковой пары (яз_{вх}, яз_{вых}).

Множество специализаций переводчика:

$$cn \in C: C = \{cn_0, cn_1, \dots cn_{nC}\},$$
 (58)

где $cn_{iC\Pi}$ определяется по формуле (35) для некоторого домена приложения ∂n .

Обучение проводится на реальных текстах с проверкой и работой над ошибками. В зависимости от состава обучающих текстов переводчик осваивает специализацию по предметной области, например, медицина или технический перевод. Чем больше переводчик учится на специальных текстах, тем выше качество перевода текстов соответствующей предметной области, но и качество перевода других областей растет просто за счет лучшего освоения языка. В процессе обучения переводчик увеличивает объем знаний, развивает умения и научается тому, как используются отдельные слова и их сочетания в разных доменах для выражения семантических единиц.

Множество обучающих текстов:

$$om \in OT: OT = \{om_0, om_1, \dots om_{nOT}\}$$

$$(59)$$

где om_{iOT} — это пара значений txt_{ex} на языке оригинала st_{ex} и txt_{ebx} на языке перевода st_{ebx} однозначно сопоставленных друг другу по смыслу (набору выражаемых семантических единиц), принятая за эталон, то есть

$$om_{iOT} = \{ (txt_{ex}; txt_{ebix}): CM | txt_{ex} = CM | txt_{ebix} \}.$$

$$(60)$$

В общем виде функцию обучения переводчика можно записать следующим образом:

$$F_{ob}: (OT, K_{o,nepi\Pi EP}, C_{o,nepi\Pi EP}) \to K_{nepi\Pi EP}, C_{nepi\Pi EP}.$$
 (61)

Схематично, функция обучения переводчика на рисунке 2.6.

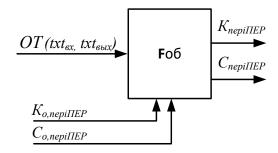


Рисунок 2.6 – Обучение переводчика

2.1.8 Обобщенная модель перевода

Обобщенно, разработанная модель процесса перевода представлена на рисунке 2.7.

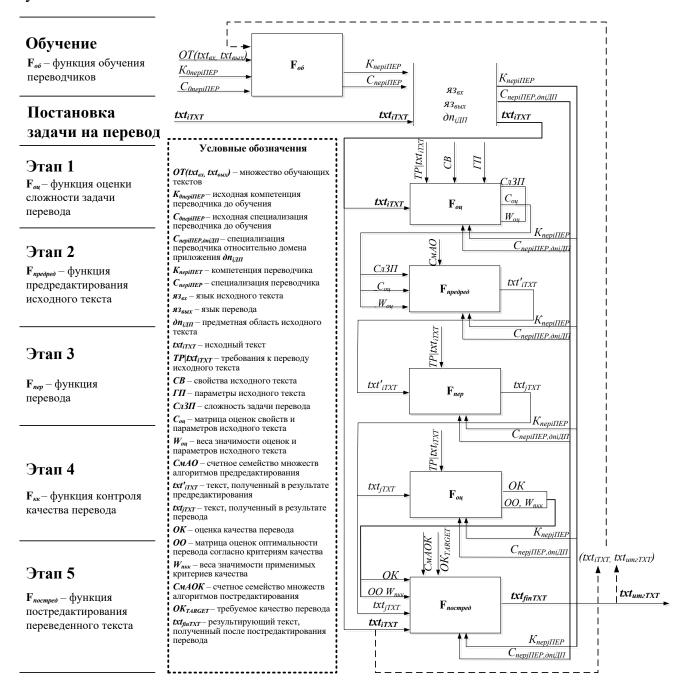


Рисунок 2.7 – Модель процесса перевода

Результаты выполненного моделирования показывают, что уже на этапе оценки исходного текста, возможно предсказать ожидаемое качество перевода на основе параметров исходного текста и компетенции и специализации переводчика [101].

Дальнейшее исследование предполагает более глубокую декомпозицию и моделирование, с четом того, что в современных системах МП эти этапы, как правило, не реализованы либо реализованы частично. Так, в системах МП не реализован этап переводческого процесса, который выполняется при «ручном переводе», а именно оценка сложности задачи перевода. На этом этапе переводчик оценивает вероятность получения качественного перевода, то есть соответствующего требованиям заказчика, и, если эта вероятность низкая, выбирает стратегию предредактирования исходного текста с целью повышения вероятности получения качественного перевода.

В соответствии с методологией определения сложности задачи перевода, описанной в п.2.1.2, для расчета сложности задачи перевода необходимо рассчитать веса значимости оценок текста $W_{oq} = \{w_{loq}, w_{2oq}, ..., w_{koq}\}$ для текстов на языке s_{loq} при переводе на язык s_{loq} . Опишем модель решения задачи поиска весов значимости признаков текста для оценки сложности задачи перевода при условии наличия некоторой экспериментальной выборки, позволяющей провести обучение модели и сравнение полученных теоретических результатов с экспериментальными.

2.2 Математическая модель поиска весов значимости признаков текста

Для поиска весов значимости признаков текста воспользуемся численным методом наименьших квадратов [102], при котором минимизируется сумма квадратов отклонений эмпирических (фактических) значений результативного признака от теоретических, полученных по уравнению (37), т.е.

$$S(w) = \sum_{i=1}^{R} (OK_i - \widehat{OK}_i (C_{oui}, w))^2,$$

$$S(w) = \sum_{i=1}^{R} (OK_i - w_0 + w_{1ou}C_{oui_1} + w_{2ou}C_{oui_2} + \dots + w_{ouk}C_{ouik})^2 \rightarrow min, \qquad (62)$$

где R — объем экспериментальной выборки.

Для решения задачи минимизации необходимо найти стационарные точки функции S(w), продифференцировав её по искомым параметрам w и приравняв производные к нулю

$$\sum_{i=1}^{R} \left(OK_i - \widehat{OK}_i(C_{oui}, w) \right) \frac{\partial \widehat{OK}_r(C_{oui}, w)}{\partial w} = 0.$$
 (63)

Получаем систему k нормальных уравнений с k неизвестными:

$$\begin{cases} \sum_{OK=Rw_{0}+w_{1ou}} \sum_{C_{oul}+w_{2ou}} \sum_{C_{ou2}+...+w_{kou}} \sum_{C_{ouk}} C_{ouk} \\ \sum_{OK:C_{oul}=w_{0}} \sum_{C_{oul}+w_{1ou}} \sum_{C_{oul}^{2}+w_{2ou}} \sum_{C_{ou2}C_{oul}+...+w_{kou}} \sum_{C_{ouk}C_{oul}} C_{ouk} C_{oul} \\ \sum_{OK:C_{ouk}=w_{0}} \sum_{C_{ouk}+w_{1ou}} \sum_{C_{oul}C_{ouk}+w_{2ou}} \sum_{C_{ou2}C_{ouk}+...+w_{kou}} \sum_{C_{ouk}^{2}} C_{ouk} C_{ouk} \end{cases}$$

Решение этой системы уравнений дает нам общую формулу поиска весов значимости W_{ou} в матричной форме:

$$W_{ou} = (C_{ou}^{T} \cdot C_{ou})^{-1} \cdot C_{ou}^{T} \cdot OK = (\frac{1}{R} C_{ou}^{T} \cdot C_{ou})^{-1} \frac{1}{R} C_{ou}^{T} \cdot OK.$$
 (64)

Предложенный алгоритм позволяет расширить область применения метода наименьших квадратов для поиска весов значимости параметров исходного текста для вероятностной оценки ожидаемого качества его перевода на целевой язык.

2.3 Математическая постановка задачи машинного перевода

Для формализации задачи машинного перевода воспользуемся методом максимального правдоподобия, который наиболее часто используется в задачах машинного обучения из-за высокой сходимости и вычислительной эффективности [103]. Метод максимального правдоподобия — статистический метод, используемый для оценки параметров распределения на основе наблюдаемых данных. Суть метода заключается в выборе таких значений параметров, которые максимизируют вероятность (или правдоподобие) получения наблюдаемых данных. Стоит отметить, что метод применим не только при традиционном выборочном подходе, но и при использовании байесовского подхода к анализу данных. Байесовские оценки обычно эквивалентны оценкам максимального правдоподобия с ограничениями [104].

Способ измерения близости к истинному параметру — вычисление квадратичной разницы между оценочными и истинными значениями параметров, где математическое ожидание вычисляется над L обучающими выборками из данных,

генерирующих распределение. Эта параметрическая среднеквадратичная ошибка уменьшается с увеличением L, и для больших L нижняя граница неравенства Крамера-Рао показывает, что среди прочих сходящихся функций у оценки максимального правдоподобия среднеквадратичная ошибка минимальна и по мере того, как число обучающих выборок приближается к бесконечности, оценка максимального правдоподобия сходится к истинному значению параметра [105].

В контексте задач МП, метод максимального правдоподобия может быть применен для оценки параметров модели языка или модели перевода. Он позволяет найти такие значения параметров, которые максимизируют вероятность генерации (или перевода) целевых предложений на основе имеющихся параллельных корпусов или наборов предложений.

Математически, для систем МП задачу перевода можно формализовать через метод максимизации функции правдоподобия следующим образом.

Опишем условия задачи, пусть:

- 1) txt_{iTXT} исходный текст на языке $яз_{ex}$;
- 2) txt_{iTXT} переведенный текст на языке $я3_{вых}$;
- 3) TXT_{itrg} множество всех возможных вариантов перевода текста txt_{iTXT} на язык $я3_{6bix}$:

$$TXT_{itrg} = \{txt_0, txt_1, ..., txt_j\},$$

где j – обще число вариантов перевода, $txt_{jTXT} \in TXT_{itrg}$;

4) OK_{itrg} – множество нормированных оценок качества перевода текста txt_{iTXT} в соответствии с требованиями к переводу $TP|txt_{iTXT}$ для всех возможных вариантов перевода TXT_{itrg} :

$$OK_{itrg} = \{OK_0, OK_1, ..., OK_j\},\$$

где OK_j — оценка качества для j-го варианта переведенного текста.

- 5) Каждому варианту перевода соответствует одна оценка качества перевода, то есть множества TXT_{itrg} и OK_{itrg} биективны: $TXT_{itrg} \leftrightarrow OK_{itrg}$;
 - 6) [minOK; maxOK] диапазон значений оценок качества перевода OK_{itrg} ;
- 7) $OK_{\partial on}$ минимально допустимое значение критерия «Высокая оценка качества перевода» при допущении, что чем выше значение OK_i , тем лучше;

9) $K\Pi_i$ — нечёткое подмножество множества OK_{itrg} , определяющее принадлежность элементов множества OK_{itrg} классу «Высокая оценка качества перевода текста txt_{iTXT} »:

$$K\Pi_i = \{(OK, \mu_{K\Pi i}(OK)) | OK \in OK_{itrg}\};$$

- $10) \mu_{K\Pi i}(OK)$ функция принадлежности, указывающая в какой степени текст txt с оценкой OK принадлежит нечеткому множеству $K\Pi_i$;
 - 11) $\mu_{K\Pi i}(OK) \in [0; 1]$ и имеет вид логистической кривой (рисунок 2.8):

$$\mu_{K\Pi i}\left(OK\right) = \frac{1}{1 + e^{-\left(\frac{OK - OK_{\partial On}}{maxOK - OK_{\partial On}}\right)2\pi}}$$
(65)

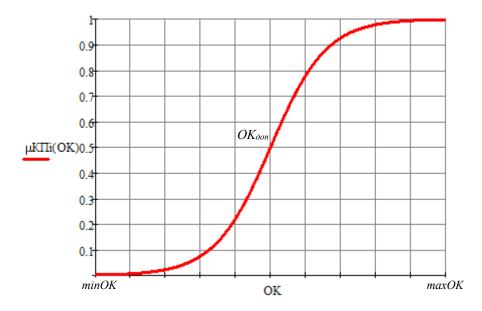


Рисунок 2.8 — График функции принадлежности $\mu_{K\Pi i}(OK)$

Требуется максимизировать правдоподобие сгенерированного системой МП текста txt_{jTXT} , то есть вероятность того, что txt_{jTXT} примет такое значение, при котором $\mu_{K\Pi i}(OK)$ будет максимальна. Тогда логарифмическая функция правдоподобия машинного перевода $F_{M\Pi}$ примет вид:

$$F_{M\Pi}(\theta, \mu_{K\Pi i}(OK)) = lnP_{\theta}(\max \mu_{K\Pi i}(OK)) \to \max_{\theta}, \tag{66}$$

где θ — параметры системы МП из множества исполнителей перевода, или переводчиков: $nep_{i\Pi EP} \in \Pi EP$, максимизирующие вероятность P получить максимальное значение функции принадлежности $\mu_{K\Pi i}(OK)$.

Решение поставленной задачи лежит в области оптимизации и совершенствования алгоритмов генерации переведенного текста. Однако, в рамках настоящего исследования стоит цель повышения качества переведенного текста на этапе подготовки к переводу до непосредственной генерации переведенного текста, поэтому далее рассмотрим задачу автоматического оптимизационного предредактирования.

2.4 Математическая постановка задачи оптимизационного предредактирования

Задача оптимизационного предредактирования состоит в том, чтобы максимизировать правдоподобие, то есть вероятность того, что при параметрах Ψ предредактора, текст txt'_{iTXT} на языке ss_{ex} будет эквивалентен txt_{iTXT} по смыслу, понятен системе МП s_{iTXT} и оценка качества s_{iTXT} относительно s_{iTXT} при генерации перевода из s_{iTXT} будет максимальной. Далее опишем задачу более подробно.

Условия задачи выглядят в следующем виде, пусть:

- 1) txt'_{iTXT} текст на языке g_{3ex} , созданный системой автоматического оптимизационного предредактирования, такой, при котором $CM'_i \rightarrow CM_i$ и $txt'_{iTXT} \neq txt_{iTXT}$, где CM смысл или упорядоченный набор семантических единиц, описываемый текстом: $CM = \{(ce_0, ce_1, ..., ce_{ncm}): ce_{icm} \in CE \}$, CM'_i и CM_i смыслы txt'_{iTXT} и txt_{iTXT} , соответственно.
- 2) TXT_{isrc} множество всех возможных вариантов предредактированного текста, т.е. выражения смысла CM_i текста txt_{iTXT} на языке startage star

$$TXT_{isrc} = \{txt_0, txt_1, ..., txt_k\},$$

где k – обще число вариантов предредактированного текста, причем txt_{iTXT} , $txt'_{iTXT} \in TXT_{isrc}$;

3) ${\it MCn3\Pi_{isrc}}$ — множество оценок сложности задачи перевода вариантов предредактирования текста ${\it txt_{iTXT}}$ в соответствии с компетентностью ${\it \overline{Knep_{i\Pi EP}}}$ и специализацией ${\it \overline{Cnep_{i\Pi EP}}}, {\it \partial n_{iZ\Pi}}$ системы МП ${\it nep_{i\Pi EP}}$ для всех возможных вариантов предредактированного текста ${\it TXT_{isrc}}$:

$$MC_{\pi}3\Pi_{isrc} = \{C_{\pi}3\Pi_{0}, C_{\pi}3\Pi_{1}, ..., C_{\pi}3\Pi_{k}\};$$

- 4) Каждому варианту предредактированного текста соответствует одна оценка сложности задачи перевода для системы МП $nep_{i\Pi EP}$, то есть множества TXT_{isrc} и $MCn3\Pi_{isrc}$ биективны: $TXT_{isrc} \leftrightarrow MCn3\Pi_{isrc}$;
- 5) $C_{n}3\Pi_{k} \in {}_{M}C_{n}3\Pi_{isrc}$ оценка сложности задачи перевода варианта предредактированного текста txt'_{iTXT} для системы МП $nep_{i\Pi EP}$;
- 6) [$minCл3\Pi$; $maxCл3\Pi$] диапазон значений оценок сложности задачи перевода $MCл3\Pi_{isrc}$;
- 7) $Cn3\Pi_{\text{доп}}$ максимально допустимое значение критерия «Низкая сложность задачи перевода» при допущении, что чем ниже значение $Cn3\Pi_k$, тем лучше;
- 8) $HCn3\Pi_i$ нечёткое подмножество множества $MCn3\Pi_{isrc}$, определяющее принадлежность элементов множества $MCn3\Pi_{isrc}$ и соответствующих элементов множества TXT_{isrc} классу «Низкая сложность задачи перевода»:

$$HC_{\pi}3\Pi_{i} = \{(C_{\pi}3\Pi, \mu_{HC_{\pi}3\Pi i}(C_{\pi}3\Pi)) | C_{\pi}3\Pi \in MC_{\pi}3\Pi_{isrc}\};$$

- 9) $\mu_{HC\pi 3\Pi i}(C\pi 3\Pi)$ функция принадлежности, указывающая в какой степени текст txt с оценкой $C\pi 3\Pi$ принадлежит нечеткому множеству $HC\pi 3\Pi_i$;
 - 10) $\mu_{HC\pi 3\Pi i}(C\pi 3\Pi) \in [0; 1]$ и имеет вид логистической кривой (рисунок 2.9):

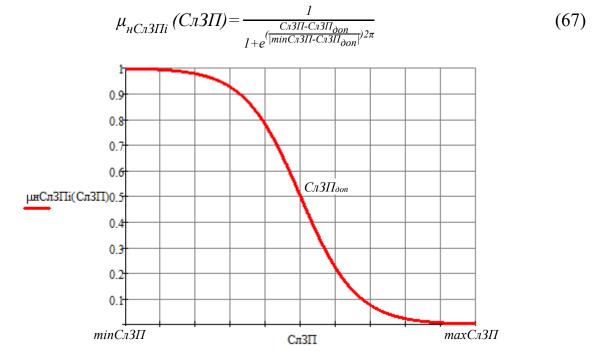


Рисунок 2.9 — График функции принадлежности $\mu_{HC\pi^3\Pi i}(C\pi^3\Pi)$

Требуется максимизировать правдоподобие сгенерированного системой оптимизационного предредактора текста txt'_{iTXT} , то есть вероятность того, что txt'_{iTXT} примет такое значение, при котором $\mu_{HC,n3\Pi i}(C,n3\Pi)$ будет максимальна.

В дискретном случае функция правдоподобия $F_{AO\Pi P}(\Psi, \mu_{HC\pi 3\Pi i}(C\pi 3\Pi))$ – вероятность выборке $\mu_{HC\pi 3\Pi i}(C\pi 3\Pi) = \{\mu_0, \mu_1, ..., \mu_l\}$ в рассматриваемой серии экспериментов равняться $\{\max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_0, \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_1, ..., \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_l\}$. Эта вероятность меняется в зависимости от Ψ :

$$F_{AO\Pi P}(\Psi, \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)) = \prod_{l=1}^{L} F_{AO\Pi P}(\mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_{l}) =$$

$$P_{\Psi}(\mu_{0} = \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_{0}) \cdot \dots \cdot P_{\Psi}(\mu_{l} = \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_{l}) =$$

$$P_{\Psi}(\mu_{0} = \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_{0}), \dots, \mu_{l} = \max \mu_{HC\pi 3\Pi i}(C\pi 3\Pi)_{l}),$$

$$(68)$$

где l – номер семпла в обучающей выборке объемом L.

Тогда логарифмическая функция правдоподобия автоматического оптимизационного предредактирования F_{AOIIP} имеет вид:

$$L_{AO\Pi P}(\Psi, \mu_{HC_{\Lambda}3\Pi i}(C_{\Lambda}3\Pi)) = lnP \ \Psi(max \ \mu_{HC_{\Lambda}3\Pi i}(C_{\Lambda}3\Pi)),$$
 (69)

где Ψ – параметры системы автоматического оптимизационного предредактирования, максимизирующие вероятность P получить максимальное значение функции принадлежности $\mu_{HCл3\Pi i}(Cл3\Pi)$.

Поскольку ln(y) монотонна, то точки максимума $F_{AO\Pi P}$ (Ψ , $\mu_{HC \Lambda 3\Pi i}(C \Lambda 3\Pi)$) и $L_{AO\Pi P}$ (Ψ , $\mu_{HC \Lambda 3\Pi i}(C \Lambda 3\Pi)$) совпадают, и оценкой максимального правдоподобия можно назвать точку максимума функции $L_{AO\Pi P}(\Psi, \mu_{HC \Lambda 3\Pi i}(C \Lambda 3\Pi))$ по Ψ . Задача оптимизации, таким образом, заключается в поиске оценки максимального правдоподобия $\widehat{\Psi}$ вектора параметров Ψ , или:

$$\widehat{\Psi} = \arg\max_{\Psi} L_{AO\PiP}(\Psi, \mu_{HC_{3}3\Pi_{i}}(C_{1}3\Pi))$$
(70)

Далее рассмотрим решение поставленной задачи оптимизации методом градиентного спуска.

2.5 Метод градиентного спуска для решения задачи автоматического оптимизационного предредактирования

Для решения задачи автоматического оптимизационного предредактирования воспользуемся методом градиентного спуска (подъема) [106]. Метод градиентного спуска и является универсальным методом решения задачи оптимизации и широко используется в задачах машинного обучения на больших объемах данных ввиду высокой скорости сходимости, гибкости в выборе шага, точности и критерия остановки, а также возможности обновлять параметры модели на отдельных подмножествах обучающей выборки данных (семпле), что делает его вычислительно эффективным.

Для решения задачи автоматического оптимизационного предредактирования необходимо найти градиент логарифмической функции правдоподобия $L_{AO\Pi P}(\Psi, \mu_{HC \Lambda 3\Pi i}(C \Lambda 3\Pi))$ — вектор, который показывает направление возрастания функции.

Учитывая, что Ψ – вектор параметров системы автоматического оптимизационного предредактирования и $\Psi = \{\psi_1, \ \psi_2, \ ..., \ \psi_m \}$, где m – количество параметров модели, градиент функции $L_{AO\Pi P}(\Psi, \mu_{HC\pi 3\Pi i}(C\pi 3\Pi))$ может быть найден по формуле:

$$\nabla L_{AO\Pi P} (\Psi) = (\partial L_{AO\Pi P} / \partial \psi_1, \partial L_{AO\Pi P} / \partial \psi_2, ..., \partial L_{AO\Pi P} / \partial \psi_m), \tag{71}$$

где ∂ $L_{AO\Pi P}$ / $\partial \psi_m$ — частная производная функции правдоподобия по m-ному параметру.

Обновление параметров Ψ происходит итеративно для каждого $\psi_m \in \Psi$:

$$\Psi^{[s+1]} = \Psi^{[s]} + \alpha \cdot \nabla L_{AODD} (\Psi^{[s]}), \tag{72}$$

где s — шаг оптимизации, $s \in [0;S]$ и S — общее число шагов оптимизации, а $\Psi^{[0]}$ — начальное приближение параметров модели; α — скорость обучения, т.е. положительное число, определяющее размер шага на каждой итерации.

Для оценки сходимости используется евклидова норма градиента функции $abla_{LAO\Pi P}(\Psi)$:

$$\|\nabla L_{AOIIP}(\Psi)\| = \sqrt{\left(\frac{\partial L_{AOIIP}}{\partial \psi_{I}}\right)^{2} + \left(\frac{\partial L_{AOIIP}}{\partial \psi_{2}}\right)^{2} + \dots + \left(\frac{\partial L_{AOIIP}}{\partial \psi_{m}}\right)^{2}}.$$
 (73)

Уменьшение нормы градиента указывает на сходимость оптимизации. Если норма градиента не снижается, это свидетельствует о медленной сходимости и необходимости изменения параметров оптимизации, например, скорости обучения α .

Оптимизация выполняется, пока норма градиента не достигла заданной точности ε , критерий остановки:

$$\|\nabla L_{AOIIP} (\Psi^{[s]})\| \le \varepsilon. \tag{74}$$

Предложенный алгоритм позволяет расширить область применения оптимизационного метода градиентного спуска путём использования элементов нечеткой логики в выражении функции правдоподобия через функцию принадлежности полученного текста низкой сложности задачи перевода для выбранной системы МП.

Выводы по второй главе

В ходе математического и процессного моделирования работы переводчика впервые обоснована целесообразность и разработана методология оценки сложности переводческой задачи. Результаты выполненного моделирования показывают, что уже на этапе оценки исходного текста, возможно предсказать ожидаемое качество перевода на основе параметров исходного текста и компетенции, и специализации переводчика.

Впервые разработана математическая модель расчета сложности задачи перевода и предложено использовать данную оценку в качестве критерия оптимизации при решении задачи предредактирования исходного текста на этапе подготовки к машинному переводу. Предложенное решение задачи оптимизационного предредактирования расширяет область применения метода градиентного спуска (подъема) за счет применения элементов нечеткой логики при формализации функции потерь.

Далее опишем методологию реализации разработанных моделей для решения задачи оптимизационного предредактирования русскоязычных текстов с целью повышения качества их перевода на английский язык.

ГЛАВА 3 МЕТОДИКА РЕАЛИЗАЦИИ ОПТИМИЗАЦИОННОГО ПРЕДРЕДАКТИРОВАНИЯ

3.1 Модель оптимизационного предредактирования

При реализации модели оптимизационного предредактирования предлагается использовать опорную языковую модель Т5. Т5 — это трансформерная архитектура моделей обработки естественного языка, которая использует подход «текст-в-текст» для решения различных задач обработки естественного языка и может быть использована для решения задач, таких как перевод, генерация текста, ответы на вопросы и т.д., путем преобразования исходного текста в целевой текст [107]. Обработка текста в модели Т5 происходит в два этапа [108]:

Этап 1. Текст токенизируется и последовательно проходит через стек энкодеров. Энкодер получает на вход список векторов (тензоров), обрабатывает его, передавая векторы в слой самовнимания для нормализации, затем передает данные в нейронную сеть с прямой связью, а далее отправляет их к следующему энкодеру. Нормализация состоит в масштабировании входных значений тензора таким образом, чтобы среднее значение было нулем, а дисперсия равнялась единице.

Этап 2. После прохождения стека энкодеров текст передается в стек декодеров для расшифровки векторных значений обратно в слова. Стек декодеров возвращает вектор чисел, который проходит через полносвязный линейный слой нейронной сети. Выходом линейного слоя являются логит-векторы, которые представляют собой распределенную вероятность выходного слова. После применения функции активации векторы заменяются конкретным словом, а матрицы преобразуются в список выходных слов.

Обучение модели Т5 проходит на корпусах параллельных текстов, состоящих из пар-примеров входных и выходных данных. Входные данные – исходный текст, а выходные данные – целевой текст.

Параметрами модели $\Psi = \{ \psi_1, \ \psi_2, \ ..., \ \psi_m \}$ являются веса нейронов сетей энкодера и декодера. Для обучения модели применяется следующий алгоритм:

1 Инициализация параметров модели случайными значениями.

Далее для каждого семпла обучающей выборки:

- 2 Кодирование входного текста в последовательность чисел.
- 3 Кодирование выходного текста в последовательность чисел.
- 4 Генерация целевого текст, используя веса модели.
- 5 Вычисление ошибки между предсказанным целевым текстом и фактическим выходным текстом.
 - 6 Обновление весов модели, используя алгоритм, описанный в п.2.5.

3.2 Обучение модели оптимизационного предредактирования

Предредактирование в переводе — это перевод с языка $яз_{ex}$ на язык $яз_{ex}$ с целью, во-первых, сделать исходный текст более понятным переводчику для адекватного подбора лексических эквивалентов, а во-вторых, минимизировать риски грамматических и стилистических ошибок. Рассматривая процесс перевода с точки зрения теории переводоведения, можно сказать, что это ряд преобразований (переводческих трансформаций), с помощью которых осуществляется перестроение от единиц оригинала к единицам перевода [109]. При генерировании переведённого текста системы МП стремятся к сохранению структуры исходного текста при условии соблюдения норм языка перевода, т.е. фактически применяют прием переводческих трансформаций — калькирование структуры предложений. Однако калькирование — одна из наиболее простых переводческих трансформаций, использование которой позволяет передать смысл текста, но значительно повышает риск стилистической и грамматической ошибки. Зная об этой особенности, возможно построить предложение на языке $яз_{ex}$ таким образом, чтобы при калькировании структуры на язык $яз_{ex}$ таким образом, чтобы при калькировании структуры на язык $яз_{ex}$ таким образом, чтобы при калькировании структуры на язык

Таким образом, создав корпус тренировочных текстов в паре $я3_{6x}$ - $я3_{6x}$, возможно настроить языковую модель, которая будет преобразовывать текст на языке $я3_{6x}$ в текст требуемой структуры для повышения качества перевода.

Для оптимизации временных затрат на подготовку исходных данных для тренировки модели предредактирования текста предлагается методика с

использованием обратного перевода для генерирования эталонного предредактированию текста. Структура параллельного корпуса исходных данных: RefCor: $[src_ref; tgt_ref]$, где src_ref — это оригинал, т.е. текст на языке $я3_{6x}$, tgt_ref — это перевод (текст на языке $я3_{6bx}$).

Методика обучения модели оптимизационного предредактирования включает следующие этапы [110]:

Этап 1. Генерация корпусов (массивов данных)

- 1) Настраиваем системы МП MT:tgt-src, MT:src-tgt.
- 2) При помощи системы MT:tgt-src переводим текст tgt_ref на язык $яз_{6x}$, получим массив текстовых данных pre_src (массив условно предредактированных текстов).
- 3) При помощи системы MT:src-tgt переводим текст src_ref на язык $яз_{вых}$, получим массив текстовых данных tgtl ($src_ref \rightarrow tgtl$).
- 4) При помощи системы MT: src-tgt переводим текст pre_src на язык $яз_{вых}$, получим массив текстовых данных tgt2 ($src_ref \rightarrow tgt1$).
- 5) Оцениваем качество выполненного перевода на язык s_{sbix} tgt1 и tgt2 относительно эталона tgt_ref , получаем массивы оценок $QC_score(tgt1)$ и $QC_score(tgt2)$.
- Этап 2. Отбор тренировочных данных и обучение оптимизационного предредактора
- 6) Для дальнейшей работы отберем тренировочный корпус TrainCor, включающий пары src_ref_i и pre_src_i , для которых наблюдается повышение оценки качества перевода на английский язык при применении предредактирования и при условии, что ΔQC_score_i является условно значимой d_{max} для выбранного типа оценки:

$$TrainCor = \{(src_ref_i; pre_src_i): \exists (tgtl_i, tgt2_i) \mid QC \ score(tgt2_i) > QC \ score(tgtl_i) & \Delta QC \ score_i \ge d_{max}\}$$

$$(75)$$

7) Настроим языковую модель $LM:src-pre_src$ для решения задачи автоматического оптимизационного редактирования текстов на языке $яз_{ex}$, в качестве тренировочного корпуса для обучения модели используем полученный корпус параллельных текстов TrainCor.

Схематично, методика обучения модели оптимизационного предредактирования представлена на рисунке 3.1.

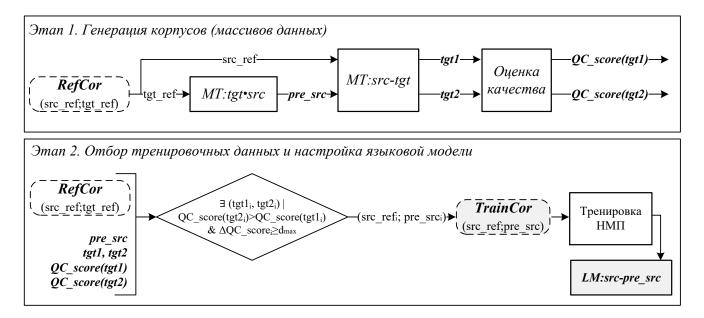


Рисунок 3.1 – Обучение модели оптимизационного предредактирования

Созданная модель может использоваться для оптимизационного предредактирования текстов как самостоятельный инструмент, так и в составе программного комплекса.

3.3 Обучение модели оценки сложности задачи перевода

Для обучения модели оценки сложности задачи перевода заданного текста txt_{iTXT} с языка st_{itXT} с языка st_{itXT} переводчиком st_{itXT} переводчиком st_{itXT} переводчиком st_{itXT} переводу st_{itXT} требуются исходные данные в виде корпуса параллельных текстов следующей структуры:

TranslatorExpCor: [src; trg; ref], где src — это оригинал, т.е. текст на языке $яз_{6x}$, trg — это перевод (текст на языке $яз_{6bx}$), выполненный переводчиком $nep_{i\Pi EP}$; ref — это контрольный перевод (текст на языке $яз_{6bx}$), т.е. проверенный эталон.

Этап 1. Оценка качества перевода. На данном этапе для каждого фрагмента текста $i=\overline{I,N}$, где N — общее число записей в корпусе TranslatorExpCor, производится оценка OK_i , при условии, что $OK_i \in R$. Полученные значения записываются

отдельным столбцом QC_score в TranslatorExpCor. Столбцы [trg; ref] удаляются как избыточные.

Этап 2. Структурный анализ предложений исходного текста. На данном этапе необходимо получить вещественные значения свойств CB текста в виде матрицы оценок C_{oui} для каждого i-го фрагмента исходного текста. В зависимости от языка исходного текста состав свойств может отличаться, однако, в целом, в задачах обработки естественного языка свойства текста условно можно разделить на группы признаков: общие (количество символов/слов/строк и т.д., стиль, язык, домен приложения и пр.), $O\Pi \subset CB$; лексические (процент покрытия текста лексическими минимумами, частотными списками, специфичность лексики и др.), $J\Pi \subset CB$; морфологические (лексические и грамматические свойства формы слова), $M\Pi \subset CB$; синтаксические (глубина глагольных и именных групп, связи между глаголами в предложениях), $C\Pi \subset CB$; признаки, основанные на базовых подсчетах (средняя длина слов и предложений и пр.), $B\Pi \subset CB$.

Для морфологического и синтаксического разбора текста используем схему Universal Dependencies [111], которая позволяет производить анализ отдельных слов в предложении и их взаимосвязей, применив преобразование, для вещественной оценки свойств всего текста.

Пусть $txt_{iTXT} = \{t_0, t_1, ..., t_m\}$, где t_m – это токен текста (слово или знак препинания), m – общее число токенов в заданном тексте; $UD = M\Pi \ UC\Pi = \{UD_0, UD_1, ..., UD_n\}$, где UD_n – это морфологическое или синтаксическое свойство токена согласно схеме Universal Dependencies, n – это общее число возможных морфологических и синтаксических свойств токена по схеме Universal Dependencies, тогда

$$\forall UD_n \exists C_{oui,k} = \frac{\sum_{l}^{m} I | f(t_l) = UD_n}{m}$$
(76)

где $t_l \in txt_{iTXT}$, $f(t_l)$ — функция морфологического/синтаксического анализа.

Общие признаки, лексические признаки и признаки, основанные на базовых подсчетах, выбираются, формализуются и рассчитываются, исходя из особенностей языка g_{sx} , и дополняют матрицу G_{oui} . Далее результаты вещественной оценки свойств текста записываются в G_{oui} Текста записываются в G_{oui} Текста удаляется как избыточный.

Этап 3. Регрессионный анализ. На данном этапе, данные TranslatorExpCor разбиваются на 2 выборки согласно принципу Парето: TranslatorExpCor_train (80%) для моделирования, TranslatorExpCor_test (20%) для валидации модели. Для выборки TranslatorExpCor_train проводится регрессионный анализ относительно целевой переменной QC score.

Коэффициенты регрессионной модели составляют матрицу весов значимости оценок свойств текста W_{ou} .

Этап 4. Оценка сложности задачи перевода. Используя данные, полученные в результате регрессионного анализа и применив формулы (37), (38), (39) рассчитывается сложность задачи перевода СлЗП для каждого *i-го* фрагмента текста.

Затем происходит отбор фрагментов текста, которым соответствуют высокие значения $Cn3\Pi_i$ по формуле (42) с учетом $Cn3\Pi_{\partial on}$ — допустимого значения сложности задачи перевода.

Оптимизационное предредактирование текста производится в соответствии с имеющимися методами и алгоритмами оптимизационного предредактирования по критерию минимизации $Cn3\Pi$.

Выводы по третьей главе

В рамках разработки методов реализации обучения моделей для оптимизационного предредактирования русскоязычных текстов впервые предложен алгоритм оценки сложности переводческой задачи для переводчика на основе его компетенции и специализации и параметров исходного текста, которая позволяет прогнозировать риски некачественного и/или несвоевременного решения задачи перевода; предложен новый алгоритм оценки русскоязычного текста по лексическим, синтаксическим и морфологическим признакам; предложена новая методика обучения модели для редактирования русскоязычных текстов, отличающаяся от существующих применением обратного перевода для сбора тренировочных данных и использованием критерия оптимизации для повышения качества машинного перевода на английский язык.

При реализации описанных методов обучения моделей следует учитывать следующие допущения:

- 1) Критерий качества перевода должен быть четко определен и формализован с возможностью получения вещественного нормированного значения. Могут применяться любые метрики оценки качества в зависимости от требований к качеству перевода.
- 2) Для тестирования МП необходим тренировочный корпус, включающий тексты на языке оригинала и перевод, принятый за эталон. В компаниях, внедривших ISO 17100 и САТ, процесс накопления тренировочных корпусов, включающих исходный текст, перевод, выполненный системой МП и проверенный перевод, утвержденный редактором, происходит автоматически в режиме реального времени.

ГЛАВА 4 ПРОГРАММНЫЙ КОМПЛЕКС ОПТИМИЗАЦИОННОГО ПРЕДРЕДАКТИРОВАНИЯ УЗКОСПЕЦИАЛЬНЫХ РУССКОЯЗЫЧНЫХ ТЕКСТОВ ДЛЯ ИХ ПЕРЕВОДА НА АНГЛИЙСКИЙ ЯЗЫК

4.1 Архитектура программного комплекса и его подсистем

Программный комплекс оптимизационного предредактирования, далее – программный комплекс, реализован на языке Python и состоит из трех основных подсистем: подсистемы тренировки языковой модели, подсистемы оценки сложности задачи перевода, подсистемы оптимизационного предредактирования русскоязычного текста и генерации машинного перевода на английский язык. Схематично архитектура программного комплекса представлена на рисунке 4.1.

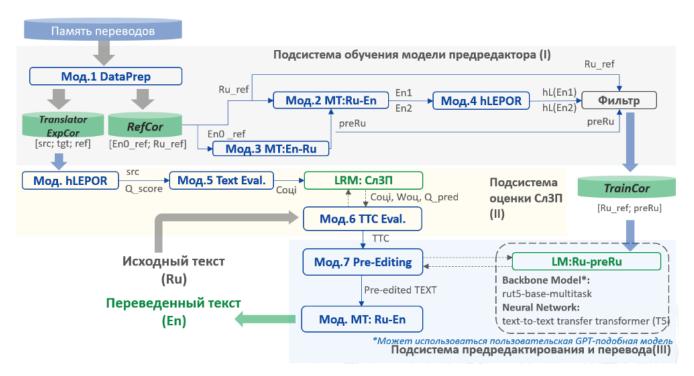


Рисунок 4.1 — Архитектура программного комплекса для повышения качества МП русскоязычных текстов на английский язык путем оптимизационного предредактирования

Подсистема тренировки языковой модели оптимизационного предредактирования русскоязычных текстов (I) состоит из пяти программных компонентов:

Мод. 1 DataPrep — модуля подготовки «сырых данных», полученных из памятей переводов «Translation Memories» поставщика лингвистических услуг, который формирует корпуса TranslatorExpCor и RefCor для тренировки моделей оценки сложности задачи перевода и оптимизационного предредактирования соответственно; Мод. 2 МТ:Ru-En, Мод. 3 МТ:En-Ru — модулей МП (генератор перевода в языковой паре русский-английский и генератор в языковой паре английский-русский); Мод. 4 hLEPOR — модуля оценки качества МП реализует алгоритм по метрике hLEPOR; модуля фильтрации данных, подходящих для тренировки модели, в котором производится отбор по условию повышения оценки качества после предредактирования. В результате обработки эталонного корпуса RefCor модулями системы 1 − 3 полученный тренировочный корпус TrainCor используется для обучения модели оптимизационного предредактирования LM:Ru-preRu.

Подсистема оценки сложности задачи перевода (II) состоит из трех модулей: модуля оценки качества перевода, выполненного системой МП, относительно эталонного по метрике hLEPOR; *Mod.5 Text Eval.* – препроцессора для взвешенной оценки свойств русскоязычного текста, включая морфологические, синтаксические, лексические и другие, всего 96 параметров; *Mod.6 TTC Eval.* – модуля оценки сложности задачи перевода с применением модели логистической регрессии *LRM:Сл3П*, которая определяет ожидаемое качество перевода полученного текста на основании взвешенной оценки его свойств и весов значимости оценок относительно системы МП.

Подсистема оптимизационного предредактирования и генерации МП (III) состоит из двух модулей: *Мод. 7 Pre-editing* — модуля автоматического предредактирования русскоязычных текстов на основе модели *LM:Ru-preRu*, которая в качестве опорной использует модель русского языка rut5-base-multitask на основе нейронной сети типа *text-to-text transfer transformer* (Т5), дообученной на корпусе параллельных русскоязычных текстов подзадачу перефразирования; модуля генерации МП с русского языка на английский язык на основе модели *Helsinki-NLP/opus-mt-ru-en*.

Программный комплекс реализован на языке Python, модули обучения моделей *LRM:СлЗП* и *LM:Ru-preRu* развернуты в Центре коллективного пользования научным оборудованием «Центр обработки и хранения научных данных ДВО РАН» на базе ВЦ ДВО РАН – обособленном подразделении ХФИЦ ДВО РАН. Тип используемой ЭВМ: ПК с архитектурой х86, х86 64.

4.1.1 Модуль автоматической очистки сырых данных из памятей переводов САТ для тренировки языковой модели

Программный модуль *Mod.1 DataPrep* реализует алгоритмы автоматической обработки параллельных корпусов текстовых данных с целью их очистки и подготовки к тренировке моделей нейронного машинного перевода. Этапы очистки данных, реализуемые в программном модуле:

- удаление дубликатов;
- удаление тегов разметки текста;
- удаление строк, где оригинал совпадает с переводом;
- удаление слишком длинных и слишком коротких строк; удаление строк, в которых менее 30% алфавитных символов.

Программный модуль может использоваться как автономное приложение в задачах прикладной компьютерной лингвистики и повышения качества перевода.

4.1.2 Модули машинного перевода

В случаях, когда сниженные требования к переводу не оправдывают затраты на привлечение высококвалифицированных исполнителей, возможно использование систем МП. Но стоит учитывать, что доступные для общего пользователя системы, такие как Google Translator, DeepL, Яндекс Переводчик и прочие, имеют ограничения по объему бесплатного использования и их политика конфиденциальности, как правило, подразумевает сбор, хранение и обработку вводимой информации при бесплатном использовании, что не всегда допустимо для пользователя. Решить данную проблему возможно путем реализации программного парсера МП на основе доступных ОрепSource моделей.

Для генерации МП с русского языка на английский язык и с английского языка на русский язык разработаны программные модули-парсеры [112] *Mod.2 MT:Ru-En*, *Mod.3 MT:En-R* с использованием модели нейронного МП на базе Transformers [113] и Ореп Source модели Helsinki-NLP от команды Hugging Face, участников

крупнейшего ежегодного соревнования разработчиком систем МП WMT [114], предварительно обученной на корпусе OPUS [115].

Парсеры получают в качестве исходных данных .xls-файл, в котором содержится построчно сегментированный текст. Используя настроенную модель МП, парсер анализирует текст и генерирует текст на требуемом языке. Результирующий массив строк, содержащий перевод, сохраняется в файл формата .xls. Разработка и инициализация парсеров МП производилась в среде Google Colab.

Разработанные парсеры могут использоваться как автономные риложения и позволяют генерировать до 10 000 страниц перевода с русского языка на английский язык в сутки на стороне сервера Google Colab. Использование OpenSource моделей и создание подобных парсеров МП на основе архитектуры Transformers может эффективно применяться с целью оптимизировать затраты и время на выполнение перевода в прикладных и исследовательских задачах различных областей знаний в условия сниженных требований к языковым и стилистическим нормам.

4.1.3 Модуль оценки качества машинного перевода

Модуль *Mod.4 hLEPOR* реализует алгоритм оценки качества МП. Для оценки МП выбрана метрика hLEPOR, которая является комбинацией существующих и доработанных факторов. Метрика показывает лучшие результаты оценки по сравнению с MPF, ROSE, METEOR, BLEU и TER и имеет наивысший балл корреляции Пирсона с человеческими суждениями по языковой паре английский-русский [116].

Оценка метрики производится путем сравнения сгенерированного МП (гипотезы) и эталонным переводом, выполненным человеком. Метрика изменяется в пределах от 0 до 1, где 0 – полное несовпадение гипотезы с эталоном, а 1 – полное совпадение.

4.1.4 Модуль препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста

В общем смысле текст — это письменное сообщение, объективированное в виде письменного документа, состоящее из ряда высказываний, объединённых разными типами лексической, грамматической и логической связи, имеющее определённый моральный характер, прагматическую установку и соответственно

литературно обработанное. Текст обладает морфологическими, лексическими, синтаксическими признаками (объективные оценки), а также признаками, определяемыми в отношение субъекта учитывая его информационную (когнитивную) трудность [89]. Когнитивная трудность текста определяется на уровне реципиента его способностью идентифицировать семантические единицы в тексте.

Объективные оценки текста могут быть рассчитаны и использованы для статистического анализа и сравнения текстов, их кластеризации и классификации, а также при создании и исследовании цифровых двойников в системах обработки естественного языка.

Модуль препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста [117] автоматически рассчитывает вещественную оценку признаков русскоязычного текста по 96 параметрам, разделенным на четыре группы:

1 Общие признаки текста, которые включают: длину предложения в символах; число токенов в предложении; долю слов в предложении; среднее количество слогов в слове; средняя длина слова в символах; долю токенов пунктуации; глубину дерева синтаксического разбора; индекс читабельности Флэша, адаптированной для русского языка И. В. Оборневой [88]. Индекс читабельности Флэша основан на гипотезе, что чем меньше слов в предложениях и чем короче эти слова, тем проще текст. Данный индекс принимает значение от 1 до 100, где тексты с индексом меньше 30 очень трудно читать, с индексом 70 и выше должны быть легко читаемыми [118]. Расчет индекса читабельности Флэша производится по формуле:

$$R_f = 206,835 - (1,3*mW) - (60,1*mS), \tag{77}$$

где mS – среднее количество слогов в словах, mW – количество слов в предложении.

2 Морфологические и синтаксические признаки текста, для оценки которых используется разбор по схеме универсальных зависимостей [104]. Морфологическая спецификация слова в схеме универсальных зависимостей состоит из трех частей: лемма слова, тег части речи и морфологические признаки, которые определяют лексические и грамматические свойства формы слова. В данном случае оценивается доля слов исследуемой части текста по каждому тегу части речи, как наиболее информативный морфологический признак. Универсальные зависимости имеют

фиксированный список тегов, включающий 17 частей речи. При синтаксическом разборе текста, каждая его единица (токен) маркируется соответствующим тегом отношения «rel=». Фиксированный список синтаксических связей включает 64 тега. Морфологический и синтаксический разбор выполняется следующим образом:

- 1) Разбиение текста на токены.
- 2) Определение для каждого токена значений свойств pos (часть речи) и rel (роль в предложении).
- 3) Подсчет количества частей речи и числа токенов по каждому тегу в каждой строке.
- 4) Расчет для всех тегов отношения количества токенов с соответствующим тегом к общему числу токенов в предложении.
- 3 Лексические признаки текста: доля именованных сущностей; доля слов, входящих в частотные списки наиболее употребимой лексики русского языка (для ТОП-1000, ТОП-3000, ТОП-5000 и ТОП-10000); доля слов с неизвестным значением частотности, то есть тех, которых нет в частотных списках; доля слов с коэффициентом вариации \mathcal{W} уйана $< D_{max}$, значение которого зависит от того, является ли слово термином (коэффициент Жуйана является лучшим из известных в настоящее время способов измерить, насколько общеупотребительным является слово, или, напротив, насколько оно специфично для отдельных предметных областей [119]). Для оценки лексических признаков используется Национальный частотный словарь русской лексики [110], который был предварительно обработан и отсортирован по значению частотности слов при помощи средств языка Python в среде программирования Google Colab. Ориентируясь на распределение вариации Жуайна для Национального частотного словаря русской лексики, принят коэффициент D_{max} =79. Оценка лексических признаков производится по токенам, включающим слова, которые предварительно были лемматизированы.

Для автоматической оценки русскоязычного текста по описанным группам параметров был разработан программный модуль на языке Python. Исходными данными для анализа являются файлы .xls, в которых собраны тексты, сегментированные по предложениям. Морфологический и синтаксический анализ выполняется при

помощи библиотеки natasha [120] для моделирования систем обработки естественного языка на основе глубокого обучения для русского языка.

Разработанный программный модуль *Mod.5 Text Eval.* является универсальным и может быть использован в качестве автономного приложения для анализа русскоязычных текстов любой тематики. Тестирование программного модуля показало, что он может использоваться для анализа корпусов текстов большого объема. Результаты оценки текстов могут быть использованы в задачах статистического анализа, сравнения текстов, выявления общих признаков, кластеризации и других задачах компьютерной лингвистики.

4.1.5 Модуль вероятностной оценки сложности переводческой задачи для систем машинного перевода

Модуль вероятностной оценки сложности переводческой задачи для систем машинного перевода *Мод.6 ТТС Eval*. позволяет рассчитывать оценку сложности задачи его перевода на английский язык с использованием модели МП [121].

Модуль обращается к модели логистической регрессии *LRM:СлЗП*, которая обучается на массиве данных, включающих исходный текст, эталонный перевод и перевод, выполненный системой машинного перевода. Обученная модель решает задачу классификации, рассчитывая вероятность попадания переведенного текста в класс «качественный перевод» в соответствии с требованиями пользователя.

Модель обучается на данных полученных в результате предварительного обучения на массиве исходных данных, обработанных в модуле препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста.

Перед обучением модели логистической регрессии *LRM:СлЗП* был проведен анализ полученных оценок параметров текста (факторов) на нормальность распределения показала наличие большого количества факторов с экспоненциальным распределением значений и факторов, большая часть наблюдений по которым равна 0 (пример представлен на рисунке 4.2). Будем учитывать это при оценке точности модели, так как нормализация и шкалирование данных не позволит устранить нулевые значения, а балансировка данных по нулевым значениям может привести к значительным потерям данных по другим факторам.

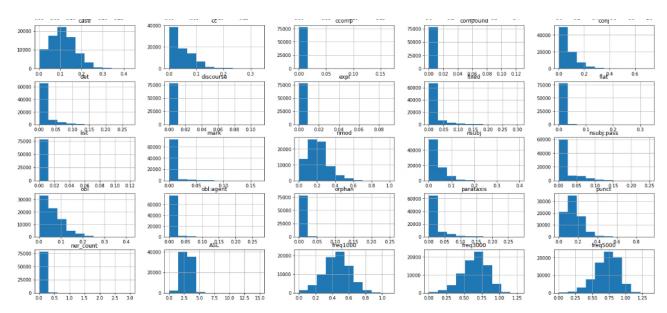


Рисунок 4.2 – Гистограммы распределения данных для части факторов

Для поиска весов значимости параметров текста и получения многофакторного уравнения поиска теоретического качества перевода использовалась модель, реализованная в библиотеках statsmodels.api и sklearn для языка Python. Результаты моделирования представлены на рисунке 4.3.

============						=======
	coef	std err	Z	P> z	[0.025	0.975]
const	-2.1237	0.119	-17.917	0.000	-2.356	-1.891
ADJ	0.8890	0.145	6.113	0.000	0.604	1.174
ADP	-1.1872	0.191	-6.222	0.000	-1.561	-0.813
ADV	-1.1501	0.323	-3.562	0.000	-1.783	-0.517
DET	1.1331	0.371	3.055	0.002	0.406	1.860
NUM	4.3565	0.227	19.228	0.000	3.912	4.801
PRON	-3.2955	0.419	-7.869	0.000	-4.116	-2.475
PROPN	-0.9230	0.136	-6.799	0.000	-1.189	-0.657
VERB	-0.9131	0.234	-3.898	0.000	-1.372	-0.454
X	2.3469	0.409	5.733	0.000	1.545	3.149
acl	-1.3931	0.356	-3.916	0.000	-2.090	-0.696
aux	5.1366	1.905	2.696	0.007	1.403	8.870
conj	2.2937	0.178	12.850	0.000	1.944	2.644
csubj	-4.6383	0.760	-6.099	0.000	-6.129	-3.148
flat	4.3092	1.322	3.261	0.001	1.719	6.899
flat:name	-4.4942	1.231	-3.650	0.000	-6.908	-2.081
list	13.0258	4.971	2.620	0.009	3.283	22.768
nmod	-1.6506	0.131	-12.597	0.000	-1.907	-1.394
nsubj:pass	-0.7066	0.310	-2.276	0.023	-1.315	-0.098
obj	-1.3491	0.261	-5.175	0.000	-1.860	-0.838
obl	-0.7080	0.224	-3.161	0.002	-1.147	-0.269
obl:agent	3.2594	0.843	3.865	0.000	1.606	4.912
parataxis	3.7789	0.360	10.506	0.000	3.074	4.484
punct	0.3628	0.168	2.166	0.030	0.034	0.691
xcomp	1.3277	0.380	3.490	0.000	0.582	2.073
freq10000	1.6783	0.064	26.415	0.000	1.554	1.803
tokens_count_log	0.3153	0.018	17.294	0.000	0.280	0.351

Рисунок 4.3 – Коэффициенты логистической регрессии

В итоговую модель вошли только те факторы (параметры исходного текста), для которых Р-значения показывают высокую значимость. Таким образом, были получены коэффициенты уравнения поиска теоретического качества перевода, выполненного тестируемым переводчиком [122]. Всего выявлено 26 параметров, из которых 12 имеют отрицательную зависимость с потенциальной оценкой качества перевода и являются мешающими (таблица 4.1).

Таблица 4.1 – Список значимых факторов модели логистической регрессии

Фактор	Коэф-т	Описание фактора (значение тега)		
		Используется для цепочек сопоставимых элементов, при-		
list	13,0258	чем в списках с более чем двумя элементами все элементы		
		списка должны модифицировать первый элемент.		
		Функциональное слово, связанное с глагольным предика-		
aux	5,1366	том и выражающее такие категории, как время, настрое-		
		ние, аспект, голос или эвиденциальность.		
csubj	-4,6383	Синтаксический субъект предложения, т.е. субъект сам яв-		
	-	ляется предложением.		
flat:name	-4,4942	Уточнение тега flat, используемое для имен.		
		Слово, функционирующее чаще всего как определитель,		
NUM	4,3565	прилагательное или местоимение, которое выражает		
1,01,1	1,5505	число и отношение к числу, такое как количество, после-		
		довательность, частота или дробь.		
		Одно из трех отношений для многословных выражений		
flat	4,3092	(MWEs) в схема универсальных зависимостей (два других		
		- фиксированные и составные), используется для (безгла-		
		вых) полуфиксированных MWE, таких как имена и даты.		
	3,7789	Связь между словом (часто главным предикатом предло-		
		жения) и другими элементами, например, сентенциальной		
parataxis		скобкой или клаузой после ":" или ";", расположенными		
		рядом без явного согласования, подчинения или аргумент-		
		ной связи с главным словом. Используется для обозначения агентов пассивных глаго-		
obl:agent	3,3594	лов.		
		Слова, заменяющие существительные или фразы суще-		
PRON	-3,2955	ствительных, значение которых можно восстановить из		
11011	-3,4933	лингвистического или экстралингвистического контекста.		
	0.01.50	Слова, которым по каким-то причинам не может быть при-		
X	2,3469	своена категория части речи.		
	2 2027	Отношение между двумя элементами, связанными коорди-		
conj	2,2937	национным союзом, таким как «и», «или» и т.д.		
frag10000	1 (792	Вхождение слова в ТОП-10000 наиболее частотных слов		
freq10000 1,6783		русского языка.		
nmad	1 6506	Используется для номинативных зависимостей от другого		
nmod	-1,6506	существительного или именной фразы.		
acl	-1,3931	Обозначает конечные и не конечные предложения, кото-		
uC1	-1,3931	рые модифицируют номинальное.		

Фактор	Коэф-т	Описание фактора (значение тега)		
obj	-1,3491	Второй по значимости аргумент глагола после субъекта.		
xcomp	1,3277	Предикативное или клаузальное дополнение без собственного субъекта.		
ADP	-1,1872	Общий термин для предлогов и постпозиций.		
ADV	-1,1501	Наречия, т.е. слова, которые обычно изменяют глаголы по таким категориям, как время, место, направление или тональность. Они также могут модифицировать прилагательные и другие наречия.		
DET	1,1331	Слова, которые изменяют существительные или фразы существительных и выражают референцию фразы существительного в контексте.		
PROPN	-0,9230	Имя собственное — это существительное (или слово номинативного содержания), которое является именем (или частью имени) конкретного лица, места или предмета.		
VERB	-0,9131	Глагол, является членом синтаксического класса слов, которые обычно обозначают события и действия, могут составлять минимальный предикат в клаузе и регулируют количество и типы других составляющих, которые могут встречаться в предложении.		
ADJ	0,8890	Прилагательные, слова, которые обычно изменяют существительные и указывают на их свойства или признаки.		
obl	-0,7080	Используется для основной части речи (существительное, местоимение, фраза существительного), функционирующего как неосновной (косвенный) аргумент или дополнение.		
nsubj:pass	-0,7066	Пассивный номинальный субъект — это фраза существительного, которая является синтаксическим субъектом пассивного предложения.		
punct	0,3628	Используется для любого знака препинания в предложении, если пунктуация сохраняется в типизированных зависимостях.		
tokens_count_log	0,3153	Натуральный логарифм количество токенов в предложении.		

Данные в таблице отсортированы по убыванию абсолютной величины коэффициента логистической регрессии. По таблице мы можем оценить какие признаки текста имеют больший вес при решении задачи бинарной классификации, а также характер зависимости (положительная/ отрицательная).

Выполняя предредактирование исходного текста с целью оптимизации данных параметров возможно повысить качество МП. Кроме того, 14 выявленных параметров имеют положительную зависимость с потенциальной оценкой качества перевода, что тоже может использоваться в оптимизационном предредактировании

для снижения влияния мешающих параметров. Проверка качества полученной модели проводилась на основе показателя ROC-AUC [123].

ROC-AUC = 0,6456 показывает, что на основании признаков исходного текста без оценки его семантики, возможно предсказывать ожидаемое качество МП. Максимальная вероятность полученных предсказаний по тестовой выборке равна 0,9765, что говорит нам о потенциале доработки данного классификатора и повышения его точности, с учетом разреженности пространства факторов и распределения их значений. Принимая во внимание допущение, что автоматическая оценка качества перевода имеет некоторую погрешность, можно говорить о состоятельности разработанной модели и предложенного алгоритма оценки сложности переводческой задачи перевода текстов нефтегазовой тематики для выбранной системы МП.

4.1.6 Предредактор русскоязычных узкоспециальных текстов для систем машинного перевода

Модуль автоматического предредактирования русскоязычных текстов Mod. 7 Pre-editing получает в качестве входных данных сегменты текста, у которых оценка сложности задачи перевода выше допустимой. Модуль работает только совместно с тренировочным модулем языковой модели для перефразирования русскоязычных технических текстов LM:Ru-preRu, обращаясь к модели оптимизационного предредактирования для генерация оптимизированных текстов.

Программный модуль *LM:Ru-preRu* реализует алгоритм инициализации и тренировки языковой модели на основе архитектуры transformers для перефразирования русскоязычных технических текстов с сохранением семантической идентичности для использования в задачах повышения качества машинного перевода.

В рамках настоящего исследования используется опорная модель расширенной архитектуры Т5-base, предобученная для решения различных задач генерации текстовых последовательностей, со следующими характеристиками: количество слоев нейросети трансформера — 16 (по 8 слоев энкодера и декодера); размер матриц весов нейронов — 1024х1024; общее число нейронов — 125 441 600; исходных объем данных предварительного обучения — более 1,56T слов.

Исходными данными для дообучения модели задаче оптимизационного предредактирования является *TrainCor* — корпус параллельных двуязычных текстов нефтегазовой тематики для тренировки языковых моделей в задачах перефразирования узкоспециальных технических русскоязычных текстов и повышения качества их перевода на английский язык, полученный в результате очистки и подготовки данных из памятей переводов средств автоматизированной поддержки переводческой деятельности.

Программный модуль *Mod.7 Pre-editing* в комплексе с *LM:Ru-preRu* может быть использован как автономное приложение в задачах оптимизационного предредактирования русскоязычных текстов либо дообучен при помощи соответствующих наборов данных для решения задач перефразирования и/или симплификации русскоязычных текстов в соответствии с заданными критериями.

4.2 Данные для обучения и тестирования программного комплекса

Исходные данные для обучения моделей и тестирования в виде русскоязычных узкоспециальных технических текстов, переведенных на английский язык, предоставлены ведущей компанией по оказанию лингвистических услуг в области технического перевода в Хабаровском крае ООО «Агентство переводов «ФИАС-Амур» в объеме \sim 60 000 ст. стр. текста (1 ст. стр. = 1800 знаков с пробелами). Из предоставленных данных сформированы корпуса RefCor объемом 139 438 семплов и TranslatorExpCor объемом 83 543 семплов. После обработки корпуса RefCor в подсистеме тренировки языковой модели оптимизационного предредактирования русскоязычных текстов в корпус TrainCor вошло 88 631 семпл. Тестовая выборка для оценки работы системы TestCor составила 16 707 семплов.

4.2.1 База данных структурного анализа предложений технических русскоязычных текстов

Путем обработки корпуса *TranslatorExpCor* модулями оценки качества машинного перевода *Мод. hLEPOR* и препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста *Мод. 6 TTC Eval.* получена база

данных структурного анализа предложений технических русско-язычных текстов, которая содержит 83543 русскоязычных предложений технического текста объемом от 50 до 300 символов, для которых выполнен расчет вещественных параметров по морфологическим, синтаксическим, лексическим и прочим признакам. Всего база данных содержит расчёты для 96 признаков. Морфологические и синтаксические признаки рассчитаны на основе универсальных зависимостей с учетом доли токенов определенного признака в предложении.

Формат: файл .xlsx. Название: DB_RUS_Tech_texts_UD. Объем: 25,4 МБ.

База данных содержит русскоязычные предложения технического текста объемом от 50 до 300 символов (Ru_test). Объем русскоязычного текста составляет ок. 5900 ст. стр. (1 ст.стр. = 1800 знаков с пробелами). Поля базы данных:

- En0_test эталонный ручной перевод, проверенный редактором;
- En1 машинный перевод;
- h_Lepor оценка машинного перевода, характеризующая близость машинного перевода эталонному, где 0 минимальная близость; 1 максимальная;
- h_Lepor_class переменная, классифицирующая предложения по признаку оценки машинного перевода на 4 класса: 1 для hLEPOR < 0,25; 2 0,25 < hLEPOR < 0,5; 3 0, 5 < hLEPOR < 0,75; 4 –hLEPOR > 0,75;
 - длина предложения в символах (len ru);
 - число токенов в предложении (tokens_count);
 - доля слов в предложении (ASW per tokens);
 - среднее количество слогов в слове (ASL);
 - средняя длина слова в символах (ASW mean len);
 - глубина дерева синтаксического разбора (unique had id count);
 - индекс читабельности Флэша (flesh readability).
- Морфологические признаки текста: ADJ, ADP, ADV, AUX, CCONJ, DET, INTJ, NOUN, NUM, PART, PRON, PROPN, PUNCT, SCONJ, SYM, VERB, X (всего 17).
- Синтаксические признаки текста: acl, acl:relcl, advcl, advmod, advmod:emph, advmod:lmod, amod, appos, aux, aux:pass, case, cc, cc:preconj, ccomp, clf, compound:lvc, compound:prt, compound:redup, compound:svc, conj,

cop, csubj; csubj:outer, csubj:pass, dep, det, det:numgov, det:nummod, det:poss, discourse, dislocated, expl, expl:impers, expl:pass, expl:pv, fixed, flat, flat:foreign, flat:name, goeswith, iobj, list, mark, nmod, nmod:poss, nmod:tmod, nsubj; nsubj:outer, nsubj:pass, nummod, nummod:gov, obj, obl, obl:agent, obl:arg, obl:lmod, obl:tmod, orphan, parataxis, punct, reparandum, vocative, xcomp (всего 64).

- доля именованных сущностей (ner count);
- доля слов, входящих в частотные списки наиболее употребимой лексики русского языка (freq1000 для ТОП-1000, freq3000 для ТОП-3000, freq5000 для ТОП-10000);
- доля слов с неизвестным значением частотности, то есть тех, которых нет в частотных списках (no_freq);
 - доля слов с коэффициентом вариации Жуйана менее 79 (D_trg).

База данных может использоваться для тренировки моделей классификации, кластеризации и регрессии в решении задачи прикладной компьютерной лингвистики.

4.2.2 Корпус параллельных двуязычных текстов нефтегазовой тематики для тренировки языковых моделей в задачах перефразирования узкоспециальных технических русскоязычных текстов и повышения качества их перевода на английский язык

Корпус *TrainCor* — это корпус параллельных двуязычных текстов нефтегазовой тематики для тренировки языковых моделей в задачах перефразирования узкоспециальных технических русскоязычных текстов и повышения качества их перевода на английский язык, который включает в себя параллельные тексты нефтегазовой тематики, сегментированные на отдельные предложения, и имеет следующую структуру: исходный текст на русском языке, отредактированный вариант исходного текста на русском языке, эталонный перевод на английский язык, проверенный редактором; машинный перевод исходного русскоязычного текста; машинный перевод отредактированного русскоязычного текста. Кроме того, корпус включает показатели качества машинного перевода, рассчитанные по методике hLEPOR и оценивающие близость машинного перевода по отношению к эталонному в диапазоне от 0 до 1. Всего корпус содержит 88631 уникальных записей.

Формат: файл .xlsx. Название: DB TrainCor RU-preRU. Объем: 26,9 МБ

Область применения корпуса – тренировка языковых моделей в задачах обработки русскоязычных текстов и повышения качества машинного перевода на английский язык.

4.3 Тестирование программного комплекса

4.3.1 Постановка задачи тестирования

Целью тестирования является проверка работоспособности программного комплекса для повышения качества машинного перевода узкоспециальных технических текстов путем автоматического оптимизационного предредактирования.

Объект тестирования: программный комплекс оптимизационного предредактирования русскоязычных текстов и их перевода на английский язык (далее – программный комплекс) [124].

Предмет тестирования: качество машинного перевода на английский язык, полученного с помощью программного комплекса.

Гипотеза тестирования: использование оценки сложности задачи перевода в качестве критерия для отбора частей текста для оптимизационного предредактирования и их последующее предредактирование позволяют повысить качество машинного перевода.

4.3.2 Начальные условия и границы проведения тестирования

Для проведения тестирования используется случайная выборка реальных данных из предварительно подготовленных и очищенных памятей переводов поставщика лингвистических услуг. *Объем тестовой выборки*: 16707 семплов. *Сложность данных*: предложения длиной от 5 до 61 токенов.

Предметная область: узкоспециальные технические русскоязычные тексты по тематикам: нефтегазодобыча и переработка, морские платформы, транспортировка нефти и газа, техническая документация проектов строительства нефтегазовых терминалов.

В рамках тестирования примем, что минимально допустимая сложность задачи перевода $Cn3\Pi_{\partial on}=1,43$. Так как, согласно уравнению (39), сложность задачи перевода обратно пропорциональна вероятности получения перевода требуемого качества, при $Cn3\Pi_{\partial on}=1,43$ данная вероятность будет составлять 0,7.

Качество перевода оценивается с использованием алгоритма hLEPOR, рассчитывающего близость сгенерированной последовательности токенов эталону в диапазоне от 0 до 1, где 0 – полное несовпадение; 1 – полное совпадение. Эталоном выступает ручной перевод, проверенный квалифицированным редактором переводов.

Характеристики вычислительной системы: процессор – Intel(R) Core(TM) i7-4770 CPU, $3.40 \mathrm{GHz}$; оперативная память – $32,0~\Gamma \mathrm{E}$; тип операционной системы – 64-разрядная OC.

4.3.3 Методология и план тестирования

Тестирование включает три этапа:

- 1 Машинный перевод тестовой выборки и оценка его качества алгоритмом $hLEPOR \rightarrow получение оценки <math>hLEPOR(En1)$.
- 2 Оптимизационное предредактирование тестовой выборки \rightarrow машинный перевод предредактированных текстов \rightarrow получение оценки *hLEPOR(En2)*.
- 3 Оценка сложности задачи перевода тестовой выборки $Cn3\Pi(Ru) \to$ применение оптимизационного предредактирования только к тем семплам, для которых $Cn3\Pi > Cn3\Pi > Cn3\Pi$ после оптимизационного предредактирования $Cn3\Pi(preRu)$ машинный перевод \to получение оценки hLEPOR (TTC/PE).
 - 4 Сравнение и анализ полученных результатов.

Схематичное представление методики тестирования языковой модели и модуля оптимизационного редактирования представлено на рисунке 4.4.

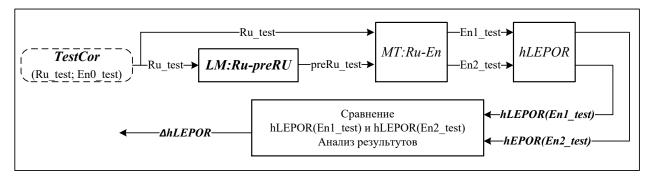


Рисунок 4.4 — Схематичное представление методики тестирования языковой модели и модуля оптимизационного редактирования

4.3.4 Результаты тестирования

Результаты, полученные в ходе тестирования, представлены на рисунках 4.5-4.7 и в таблице 4.2.

Принятые обозначения:

- mean математическое ожидание;
- std среднеквадратичное отклонение;
- min минимальное значение выборки;
- 25% значение, меньше которого 25% значений выборки;
- 50% медиана, т.е. значение, меньше и больше которого 50% значений выборки;
 - 75% значение, меньше которого 75% значений выборки;
 - тах максимальное значение в выборке.

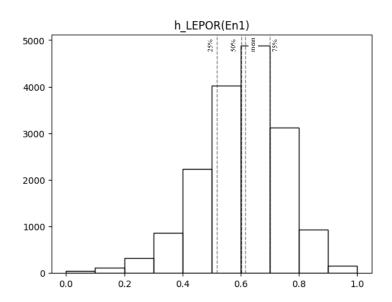


Рисунок 4.5 – Оценка качества машинного перевода до применения оптимизационного предредактирования

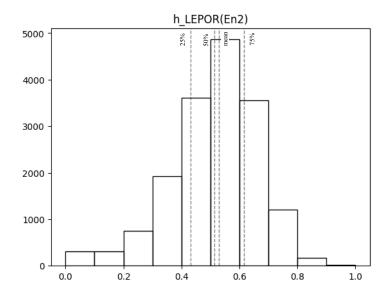


Рисунок 4.6 — Оценка качества машинного перевода после применения оптимизационного предредактирования ко всем семплам тестовой выборки

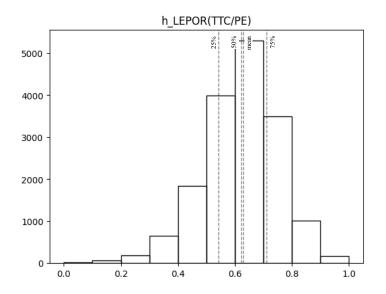


Рисунок 4.7 — Оценка качества машинного перевода после применения оптимизационного предредактирования к семплам тестовой выборки, отобранным по критерию сложности задачи перевода

	h_LEPOR (En1)	h_LEPOR (En2)	h_LEPOR (TTC/PE)	СлЗП (Ru)	СлЗП (preRu)
mean	0,603154	0,51325	0,621428	2,058798	1,082323
std	0,142844	0,151931	0,13113	0,721968	0,053624
min	0	0	0	1,009145	1
25%	0,518918	0,430738	0,541798	1,650413	1,040528
50%	0,614855	0,529115	0,629695	1,885655	1,074627
75%	0.70099	0.61654	0.710368	2 23849	1 115369

Таблица 4.2 – Результаты тестирования программного комплекса

max

Результаты тестирования показывают, что применение оптимизационного предредактирования без оценки его необходимости ведет к снижению общего качества выполненного перевода в среднем на 15%. Целесообразно использовать критерии отбора сегментов для оптимизационного предредактирования.

20,19173

1,325531

Использование оценки сложности задачи перевода в качестве критерия отбора сегментов и их последующего оптимизационного предредактирования позволяет повысить качество перевода. В рассматриваемой выборке общее среднее повышение качества перевода составило 3% при снижении сложности задачи перевода в среднем на 48%.

Рассмотрим результаты применения оптимизационного предредактирования. С использованием оценки сложности задачи перевода было отобрано 4071 семпл для оптимизационного предредактирования (24,4% тестовой выборки).

Примеры оптимизационного предредактирования и его влияния на сложность задачи перевода для системы МП представлены в таблице 4.3. Результаты перевода текстов после оптимизационного предредактирования на английский язык представлены в таблице 4.4.

Таблица 4.3 – Оптимизационное предредактирование русскоязычных текстов

Семпл	Исходный текст	Текст после предредактирования	ΔСл3Π
	(Ru_ref)	(preRU)	
1	Высоковольтные испытания про-	Испытания на высоковольтные ка-	-1,356
	водятся по отдельно разрабатывае-	бели проводятся в соответствии с от-	
	мой и утверждаемой «Программе	дельно разработанной и утвержден-	
	проведения высоковольтных ис-	ной программой испытаний высоко-	
	пытаний кабеля 110 кВ».	вольтных кабелей 110 кВ.	
2	Оборудование должно быть рас-	Оборудование должно быть спо-	-0,444
	считано на двойные фидеры, а	собно управлять двумя фидерами, в	

Семпл	Исходный текст	Текст после предредактирования	ΔC л 3Π
	(Ru_ref)	(preRU)	
	если такое оборудование отсут-	случае отсутствия такого оборудова-	
	ствует, в центральном шкафу	ния в центральном шкафу должен	
	предусматривают установку кон-	быть установлен переключатель	
	троллера автоматического ввода	ввода резерва.	
	резерва.		
3	ТУ на поставку включают в себя,	Спецификация покупки должна со-	-0,167
	помимо прочего, следующее:	держать и не ограничиваться:	
4	По результатам месяца подготовка	Отчет о ходе месяца (10 числа) о от-	-0,271
	отчета (10 число) по отклонениям	клонениях от плана работы.	
	от намеченного графика.		

Таблица 4.4 — Результаты машинного перевода на английский язык

Семпл	MП исходного текста (En1)	МП после предредактирования	$\Delta hLEPOR$
		(En2)	
1	The high voltage tests are conducted	High voltage cables shall be tested ac-	0,271
	on a separate design and approval of	cording to a separately developed and	
	the 110 kV high voltage test pro-	approved 110 kV high voltage cables	
	gramme.	programme.	
2	The equipment shall be designed for	The equipment shall be capable of con-	0,146
	double feeders, and if such equip-	trolling two feeders, in the absence of	
	ment is not available, an automatic	such equipment, a standby switch shall	
	backup controller shall be installed in	be installed in the central cabinet.	
	the central cabinet.		
3	TA for supply includes, inter alia, the	The purchase specification shall con-	0,135
	following:	tain and not be limited to:	
4	Based on the month's results, the re-	Monthly progress report (10th) on de-	0,131
	port (10 times) is based on deviations	viations from the workplan.	
	from the schedule.		

Результаты оценки качества после оптимизационного предредактирования представлены в таблице 4.5.

Таблица 4.5 — Результаты применения оптимизационного предредактирования

	h_LEPOR	h_LEPOR	AhLEPOR	Сл3П	СлЗП
	(En1)	(En2)		(Ru)	(preRu)
mean	0,508044336	0,583044424	0,075	2,111294	1,075499
std	0,140363512	0,128616684	0,0739	0,75958	0,052124
min	0	0,10348	1E-05	1,431731	1
25%	0,4253	0,506915	0,02056	1,692549	1,035022
50%	0,52027	0,59196	0,05472	1,928669	1,067745
75%	0,60549	0,67087	0,104235	2,314225	1,107176
max	0,99834	1	0,58203	20,19173	1,310526

Показано, что качество перевода отдельных сегментов, подвергшихся оптимизационному предредактированию, в среднем, возросло на 15%. Максимальное повышение качества перевода составило 30% в отдельных сегментах.

Результаты тестирования программного комплекса подтверждают работоспособность программного комплекса в части повышения качества машинного перевода узкоспециальных технических текстов путем автоматического оптимизационного предредактирования и использования оценки сложности задачи перевода в качестве в качестве критерия оптимизации. Это открывает широкие возможности для практического применения разработанной оценки при принятии решений в переводческих проектах, в том числе при использовании систем МП, а также использования оценки сложности задачи перевода в качестве критерия оптимизации. При этом программный комплекс в совокупности описанных подсистем и программных модулей имеет потенциал к доработке и повышении эффективности.

4.4 Внедрение программного комплекса в контур автоматизации процессов переводческой деятельности

Программный комплекс был внедрен в практическую деятельность ООО «Агентство переводов «ФИАС-Амур» (г. Комсомольск-на-Амуре), далее – Агентство переводов.

Целью внедрения программного комплекса в деятельность Агентства переводов является повышение качества машинного перевода, снижение затрат на его производство за счет повышения производительности редакторов переводов.

Внедрение программного комплекса проходило в несколько этапов:

- 1) Анализ бизнес-процесса осуществления перевода «как есть» с учетом средств автоматизации и поддержки процесса.
- 2) Внедрение программного комплекса в контур автоматизации бизнес-процесса перевода.
 - 3) Обучение сотрудников.
 - 4) Оценка и анализ результатов внедрения.

Схематично, процесс осуществления перевода в Агентстве переводов до внедрения программного комплекса представлен на рисунке 4.8.

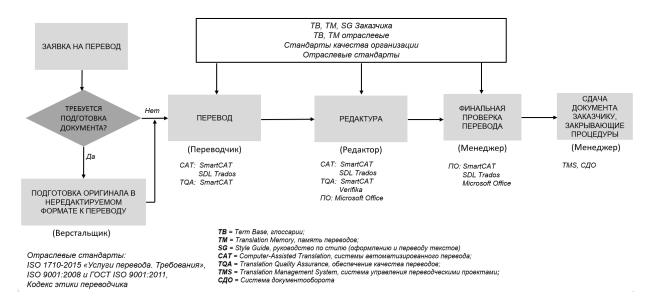


Рисунок 4.8 — Бизнес-процесс осуществления перевода до внедрения программного комплекса

Для оценки производительности переводчиков и редакторов в Агентстве переводов используется оценка количества сданных менеджеру проекта стандартных страниц текста за один час рабочего времени, 1 стандартная страница = 1800 знаков с пробелами. На момент проведения исследования средняя производительность редакторов машинного перевода составила 3,8 страниц в час.

В рамках внедрения программного комплекса в контур автоматизации бизнес-процесса перевода он был адаптирован к специфике деятельности Агентства переводов. Для этого была разработана процедура, позволяющая сотрудникам использовать функциональность программного комплекса без прямой интеграции с системами поддержки переводческой деятельности (САТ).

Процедура включает следующие шаги:

- 1 Сегментация исходного русскоязычного текста при помощи любой из используемых систем CAT.
 - 2 Экспорт сегментированного текста в файл .xlsx.

- 3 Использование программного комплекса для получения машинного перевода на английский язык.
- 4 Сохранение исходного текста и полученного машинного перевода в виде параллельного корпуса в формате таблицы файла. xslx.
 - 5 Импорт полученного корпуса в базу переводов требуемой системы САТ.
 - 6 Редактура полученного перевода в выбранной среде САТ.

В процессе внедрения программного комплекса на стороне Агентства переводов участвовали менеджер проектов и два редактора переводов. Все сотрудники прошли обучение работе с программным комплексом, для чего были разработаны учебные материалы и проведено несколько практических занятий.

Схематично, процесс осуществления перевода в Агентстве переводов до внедрения программного комплекса представлен на рисунке 4.9.

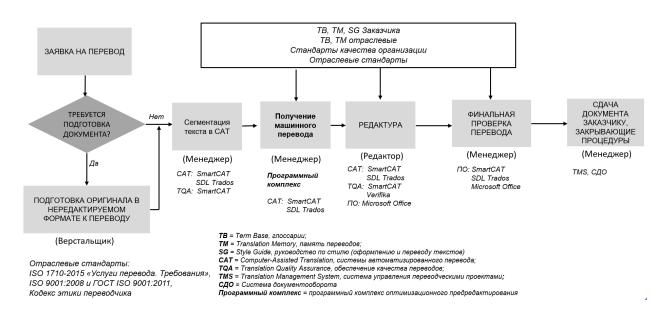


Рисунок 4.9 — Бизнес-процесс осуществления перевода после внедрения программного комплекса

Мониторинг результатов внедрения проводился в течение двух месяцев с момента интеграции программного комплекса в контур автоматизации процессов перевода и включал оценку объемов перевода, производительности редакторов переводов, опрос сотрудников Агентства переводов.

За время мониторинга при помощи программного комплекса было переведено на английский язык 677 стандартных страниц текста. При этом, средняя производительность редакторов переводов при работе с программным комплексом составила 4,3 стандартных страницы в час. Таким образом, внедрение программного комплекса позволило увеличить производительность редакторов переводов на 13,16%.

Опрос сотрудников позволил выявить две основные причины повышения производительности:

- 1 Повышение качества машинного перевода, ввиду чего сокращается количество требуемых корректировок.
- 2 Автоматизация работы с проблемными сегментами текста. Использование оценки сложности задачи перевода позволяет отфильтровать сегменты с ожидаемо низким качеством машинного перевода и применить корректирующие действия. Так, менеджер проектов в обратной связи указал, что «удобно производить оценку текста на этапе подготовки к переводу и задать клиенту уточняющие вопросы по спорным моментам, например, по расшифровке аббревиатур».

Для повышения точности оценки сложности задачи перевода и качества машинного перевода рекомендуется систематически проводить дообучение моделей, реализованных в программном комплексе, на данных памятей переводов Агентства переводов.

Выводы по четвертой главе

Результаты тестирования программного комплекса подтверждают работоспособность программного комплекса в части повышения качества машинного перевода узкоспециальных технических текстов путем автоматического оптимизационного предредактирования и использования оценки сложности задачи перевода в качестве в качестве критерия оптимизации.

Программный комплекс в совокупности описанных подсистем и программных модулей имеет потенциал к доработке и повышении точности. Интеграция созданной системы оптимизационного предредактирования в контур автоматизации процессов перевода технической документации позволит снизить затраты на постредактирование МП и организацию переводческих процессов.

Использование полученных результатов исследования и внедрение программного комплекса в работу Агентства переводов позволило повысить эффективность использования машинного перевода и производительность труда редакторов переводов, оптимизировав тем самым затраты на оказание услуг перевода узкоспециальной технической документации.

ЗАКЛЮЧЕНИЕ

В ходе проведения исследования получены следующие основные результаты:

- 1 Разработана математическая модель процесса перевода и вероятностной оценки сложности задачи перевода, отличительной особенностью которой является возможность получения результатов в аналитическом виде.
- 2 Предложен новый алгоритм оценки русскоязычного текста по лексическим, синтаксическим и морфологическим признакам, отличающийся возможностью анализа текста по 96 признакам с получением вещественных оценок по каждому из них.
- 3 Предложен алгоритм определения стратегии оптимизационного предредактирования русскоязычного текста с целью повышения качества его перевода на английский язык по критериям пользователя перевода с использованием моделей машинного перевода, отличающийся тем, что в качестве критерия оптимизации используется вероятностная оценка сложности задачи перевода.
- 4 Разработана теория вероятностной оценки сложности задачи перевода, позволяющая приближенно вычислять ожидаемое качество перевода заданного текста заданным переводчиком в соответствии с формализованными требованиями к переводу.
- 5 Предложен новый алгоритм, позволяющий расширить область применения оптимизационного метода градиентного спуска путём использования элементов нечеткой логики в выражении функции правдоподобия через функцию принадлежности полученного текста низкой сложности задачи перевода для выбранной системы МП.
- 6 Предложен новый алгоритм, позволяющий расширить область применения метода наименьших квадратов для поиска весов значимости параметров исходного текста для вероятностной оценки ожидаемого качества его перевода на целевой язык.
- 7 Реализован программный комплекс анализа русскоязычного текста и вероятностной оценки сложности задачи его перевода на английский язык с

использованием модели машинного перевода, отличающихся от существующих отсутствием необходимости оптимизации алгоритмов и моделей генерирования текста перевода.

- 8 Программно реализован алгоритм оценки русскоязычного текста по лексическим, синтаксическим и морфологическим признакам, который позволяет анализировать русскоязычные тексты и определять признаки, значимые при решении задачи перевода с учетом выбранной системы машинного перевода, отличающийся возможностью получения вещественных значений оценок текста.
- 9 Реализован программный комплекс для повышения качества машинного перевода текстов с русского языка на английский язык, отличающийся от существующих применением оптимизационного предредактирования на основе вероятностной оценки сложности задачи перевода для повышения качества машинного перевода текстов с русского языка на английский язык.

Поставленные цели и задачи исследования достигнуты, что подтверждается результатами тестирования и внедрения разработанных методов, алгоритмов и комплексов программ.

Разработанные методы, алгоритмы и комплексы программ могут быть масштабированы на различные языковые пары и способы перевода, включая ручной перевод, они намечают подходы к управлению рисками, связанными с качеством перевода в зависимости от компетенции выбранных исполнителей, и предоставит индустрии инструменты объективной оценки исполнителей в рамках поставленной задачи на перевод, автоматизированной подготовки текстов к переводу и повышения качества перевода, в том числе для редакторов без знания языка перевода.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 Search for peer-reviewed journal articles and book chapters [Электронный ресурс] // Science Direct: Platform of Peer-Reviewed Literature. 2023. Режим доступа: https://www.sciencedirect.com/
- 2 Cho K. et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation // arXiv preprint arXiv:1406.1078. 2014.
- 3 Castilho S. et al. Is neural machine translation the new state of the art? // The Prague Bulletin of Mathematical Linguistics. − 2017. − №. 108.
- 4 Wu Y. et al. Google's neural machine translation system: Bridging the gap between human and machine translation // arXiv preprint arXiv:1609.08144. 2016.
- 5 Screen B. What effect does post-editing have on the translation product from an end-user's perspective // The Journal of Specialised Translation. 2019. T. 31. C. 133-157.
- 6 Guerberof Arenas A., Moorkens J. Machine translation and post-editing training as part of a master's programme // Jostrans: The Journal of Specialised Translation. 2019. №. 31. C. 217-238.
- 7 Кислова Е. Новые смыслы профессии в условиях высокой неопределенности будущего [Электронный ресурс] / Е. Кислова, А. Мустаев, Е.Гайдерова // Пленарное заседание Translation Forum Russia 2023. 2023. Режим доступа: https://tconference.ru/provisionalprogramme/
- 8 Hutchins J. Machine translation: History and general principles // The encyclopedia of languages and linguistics. 1994. T. 5. C. 2322-2332.
- 9 Hutchins W. J. Machine translation: A brief history // Concise history of the language sciences. Pergamon, 1995. C. 431-445.
- 10 Hutchins J. Machine translation: history // Encyclopedia of languages and linguistics. 2006. C. 375-383.
- 11 Sreelekha S. et al. A survey report on evolution of machine translation // Int. J. Control Theory Appl. − 2016. − T. 9. − №. 33. − C. 233-240.

- 12 Chéragui M. A. Theoretical Overview of Machine translation // ICWIT. 2012. C. 160-169.
- 13 Митренина О. В. Назад, в 47-й: к 70-летию машинного перевода как научного направления //Вестник Новосибирского государственного университета. Серия: Лингвистика и межкультурная коммуникация. 2017. Т. 15. №. 3. С. 5-12.
- 14 Schwartz L. The history and promise of machine translation //Innovation and expansion in translation process research. 2018. C. 161.
- 15 Hutchins J. Machine translation: A concise history // Computer aided translation: Theory and practice. 2007. T. 13. №. 29-70. C. 11.
- 16 Holzinger A. From machine learning to explainable AI // 2018 World symposium on digital intelligence for systems and machines (DISA). IEEE, 2018. C. 55-66.
- 17 Al'tshuller G. S. The innovation algorithm: TRIZ, systematic innovation and technical creativity. Technical innovation center, Inc., 1999.
- 18 Gadd K. TRIZ for engineers: enabling inventive problem solving. John Wiley & Sons, 2011.
- 19 Bertoncelli T., Mayer O., Lynass M. Creativity, learning techniques and TRIZ // Procedia Cirp. 2016. T. 39. C. 191-196.
- 20 Berdonosov V. D. Fractality of knowledge and TRIZ // Procedia Engineering. 2011. T. 9. C. 659-664.
- 21 Altshuller G. 40 principles: TRIZ keys to technical innovation. Technical Innovation Center, Inc., 2002. T. 1.
- 22 Бердоносов В. Д., Животова А. А. Исследование эволюции объектно-ориентированных языков программирования // Ученые записки Комсомольского-на-Амуре государственного технического университета. 2014. Т. 1. №. 2. С. 35-43.
- 23 Berdonosov V. et al. Perspectives for development of automatized enterprise management systems // Proceedings of the TRIZfest-2015 International Conference. 2015. C. 154-163.
- 24 Berdonosov V. D., Kozlita A. N., Zhivotova A. A. TRIZ evolution of black oil coker units // Chemical engineering research and design. 2015. T. 103. C. 61-73.

- 25 Berdonosov V. D., Redkolis E. V. TRIZ evolutionary approach: Main points and implementation // Research and Practice on the Theory of Inventive Problem Solving (TRIZ) Linking Creativity, Engineering and Innovation. 2016. C. 95-111.
- 26 Hutchins J., Lovtskii E. Petr Petrovich Troyanskii (1894–1950): A forgotten pioneer of mechanical translation // Machine translation. 2000. T. 15. C. 187-221.
- 27 Троянский, П. П. Авторское свидетельство № 40995 A1 СССР, МПК В41В 13/00. Машина для подбора и печатания слов при переводе с одного языка на другой : № 134430 : заявл. 05.09.1933 : опубл. 31.01.1935.
- 28 Oswald V. A. Word-by-word translation // Proceedings of the Conference on Mechanical Translation. 1952.
- 29 Yngve V. H., Charney E. K., Klima E. S. Mechanical translation. Research Laboratory of Electronics (RLE) at the Massachusetts Institute of Technology (MIT), 1962.
- 30 Booth A. D. The practical realisation of machine translation // Methodos. 1956. T. 8. C. 23-33.
- 31 Brown A. F. R. Automatic translation of languages // Information storage and Retrieval. $-1964. -T. 2. -N_{\odot}. 1. -C. 1-28.$
 - 32 Tosh L. W. Data preparation for syntactic translation // COLING 1965. 1965.
- 33 Masterman M. Semantic message detection for machine translation, using an interlingua //Proceedings of the International Conference on Machine Translation and Applied Language Analysis. 1961. C. 438-474.
- 34 Tosh L. Stratificational grammar and interlingual mapping for automatic translation // Actes du Xe Congrès International des Linguistes. 1967. T. 4. C. 1049-1059.
- 35 Aramaki E., Kurohashi S. Example-based machine translation using structural translation examples // Proceedings of the First International Workshop on Spoken Language Translation: Evaluation Campaign. 2004. C.91-94.
- 36 Nepveu L. et al. Adaptive language and translation models for interactive machine translation // Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing. 2004. C. 190-197.
- 37 Brown P. F. et al. A statistical approach to machine translation // Computational linguistics. − 1990. − T. 16. − №. 2. − C. 79-85.

- 38 Junczys-Dowmunt M., Grundkiewicz R. Phrase-based machine translation is state-of-the-art for automatic grammatical error correction // arXiv preprint arXiv:1605.06353. 2016.
- 39 Yamada K. A syntax-based statistical translation model. University of Southern California, 2003.
- 40 Johnson M. et al. Google's multilingual neural machine translation system: Enabling zero-shot translation // Transactions of the Association for Computational Linguistics. 2017. T. 5. C. 339-351.
- 41 Lample G. et al. Unsupervised machine translation using monolingual corpora only // arXiv preprint arXiv:1711.00043. 2017.
- 42 Chu C., Wang R. A survey of domain adaptation for neural machine translation // arXiv preprint arXiv:1806.00258. 2018.
- 43 Costa-Jussa M. R., Fonollosa J. A. R. Latest trends in hybrid machine translation and its applications // Computer Speech & Language. − 2015. − T. 32. − №. 1. − C. 3-10.
- 44 Животова А. А., Бердоносов В. Д. Перспективные направления развития систем машинного перевода // Информатика и системы управления. 2022. №2(72). С.116-132.
- 45 Zhivotova A. A., Berdonosov V. D., Redkolis E. V. Machine translation systems analysis and development prospects // Proceedings of 2020 International Multi-Conference on Industrial Engineering and Modern Technologies (FarEastCon-2020). Vladivostok: Russia. 2020.
- 46 Doherty S., O'Brien S. Assessing the usability of raw machine translated output: A user-centered study using eye tracking // International Journal of Human-Computer Interaction. -2014. T. 30. N. 1. C. 40-51.
- 47 Lear A. et al. Why Can't I Just Use Google Translate? A Study on the Effectiveness of Online Translation Tools in Translation of Coas // Value in Health. − 2016. − T. 19. − №. 7. − C. A387.
- 48 Quinci C., Pontrandolfo G. Testing neural machine translation against different levels of specialization // trans-com. 2023. №16. C.174-209.

- 49 Canfora C., Ottmann A. Risks in neural machine translation // Translation Spaces. 2020. T. 9. №. 1. C. 58-77.
- 50 Животова, А. А. Автоматизация предредактирования исходного текста для повышения качества машинного перевода / А. А. Животова, В. Д. Бердоносов, И. А. Лошманова // Наука, инновации и технологии: от идей к внедрению : Материалы II Международной научно-практической конферен-ции молодых ученых, Комсомольск-на-Амуре, 14–18 ноября 2022 года / Редколлегия: А.В. Космынин (отв. ред.) [и др.]. Том 1. Комсомольск-на-Амуре: Комсомольский-на-Амуре государственный университет, 2022. С. 366-370.
- 51 Kumar V. et al. A machine assisted human translation system for technical documents // Proceedings of the 8th International Conference on Knowledge Capture. 2015. C. 1-5.
- 52 Azadi F., Khadivi S. Improved search strategy for interactive predictions in computer-assisted translation // Proceedings of Machine Translation Summit XV: Papers. 2015. T.1. C.319–332.
- 53 Federico M. et al. Machine translation enhanced computer assisted translation // Nice MT Summit XIV. 2013. T. 1. C. 425.
- 54 Alabau V. et al. User evaluation of advanced interaction features for a computer-assisted translation workbench // Machine Translation Summit XIV. 2013. T.1. C.361–368.
- 55 Yamada M. The impact of Google neural machine translation on post-editing by student translators // The Journal of Specialised Translation. $-2019. T. 31. N_{\odot}. 1.$ -C. 87-106.
- 56 Balling L. W., Carl M., O'Brian S. (ed.). Post-editing of machine translation: Processes and applications. Cambridge Scholars Publishing, 2014.
- 57 Herbig N. et al. Multi-modal approaches for post-editing machine translation // Proceedings of the 2019 CHI conference on human factors in computing systems. 2019. C. 1-11.
- 58 Toledo Báez M. C. Machine translation and post-editing: impact of training and directionality on quality and productivity // Tradumàtica. − 2018. − №. 16. − C. 24-34.

- 59 Feifei F. et al. A Study of Pre-editing Methods at the Lexical Level in the Process of Machine Translation // Arab World English Journal for Translation & Literary Studies. -2022. -T.6. $-N_{\odot}$ 2. -C.54-69.
- 60 Zheng Y., Peng C., Mu Y. Designing Controlled Chinese Rules for MT Pre-Editing of Product Description Text //International Journal of Translation, Interpretation, and Applied Linguistics (IJTIAL). -2022. - T. 4. - No. 2. - C. 1-13.
- 61 Seretan V., Bouillon P., Gerlach J. A Large-Scale Evaluation of Pre-editing Strategies for Improving User-Generated Content Translation // LREC. 2014. C. 1793-1799.
- 62 Shei C. C. Teaching MT through pre-editing: Three case studies // Proceedings of the 6th EAMT Workshop: Teaching Machine Translation. 2002.
- 63 Liang Y., Han W. Source text pre-editing versus target text post-editing in using Google Translate to provide health services to culturally and linguistically diverse clients // Science, Engineering and Health Studies. 2022. T. 16. C. 22050009-22050009.
- 64 Gerlach J. et al. Combining pre-editing and post-editing to improve SMT of user-generated content // Proceedings of the 2nd Workshop on Post-editing Technology and Practice. 2013.
- 65 Yu K. N. et al. Pre-editing English news texts for machine translation into Russian // Language Studies and Modern Humanities. 2022. T. 4. №. 1. C. 25-30.
- 66 Marzouk S., Hansen-Schirra S. Evaluation of the impact of controlled language on neural machine translation compared to other MT architectures // Machine Translation.

 2019. T. 33. №. 1-2. C. 179-203.
- 67 Arenas A. G. Pre-editing and post-editing // The Bloomsbury companion to language industry studies. -2019.-C.333-360.
- 68 Sánchez-Gijón P., Kenny D. Selecting and preparing texts for machine translation: Pre-editing and writing for a global audience // Machine translation for everyone: Empowering users in the age of artificial intelligence. 2022. T. 18. C. 81.
- 69 Zhivotova A. A., Berdonosov V. D., Redkolis E. V. Improving the Quality of Scientific Articles Machine Translation while Writing Original Text // Proceedings of 2020 International Multi-Conference on Industrial Engineering and Modern Technologies (FarEastCon-2020). Vladivostok: Russia, 2020.

- 70 Machine translation tips [Электронный ресурс] // IBM Cloud Docs. 2016. Режим доступа: https://cloud.ibm.com/docs/GlobalizationPipeline?topic=GlobalizationPipeline-globalizationpipeline tips
- 71 Writing for a global audience [Электронный ресурс] // Google developer documentation style guide. 2020. Режим доступа: https://developers.google.com/style/translation.
- 72 Miyata R. et al. Evaluating the usability of a controlled language authoring assistant // The Prague bulletin of mathematical linguistics. $-2017. T. 108. N_{\odot}. 1. C. 147.$
- 73 O'Brien S. Controlling controlled english // EAMT Workshop: Improving MT through other language technology tools: resources and tools for building MT. 2003.
- 74 Осокина С. А. Концепция «легкого языка» и перспективы ее развития в лингвистике // Филология и человек. 2022. №. 2. С. 115-133.
- 75 Muegge U. Controlled Language Optimized for Uniform Translation (CLOUT). Bepress, 2002.
- 76 Polo L. R. Controlled Language and the Implementation of Machine Translation for Technical Documentation // Translating and the Computer 27. 2005.
- 77 Hiraoka Y., Yamada M. Pre-editing plus neural machine translation for subtitling: effective pre-editing rules for subtitling of TED Talks // Proceedings of Machine Translation Summit XVII: Translator, Project and User Tracks. 2019. C. 64-72.
- 78 Miyata R., Fujita A. Dissecting human pre-editing toward better use of off-the-shelf machine translation systems // Proceedings of the 20th annual conference of the european association for machine translation (EAMT). 2017. C. 54-59.
- 79 Taufik A. Pre-editing of Google neural machine translation //Journal of English Language and Culture. 2020. T. 10. №. 2. C. 64–74.
- 80 Mercader-Alarcón J., Sánchez-Matínez F. Analysis of translation errors and evaluation of pre-editing rules for the translation of English news texts into Spanish with Lucy LT //Tradumàtica: traducció i tecnologies de la informació i la comunicació. − 2016. − №. 14. − C. 172-186.
- 81 Ниценко А. В., Шелепов В. Ю. О некоторых подходах к проблеме автоматической адаптации русскоязычных текстов //Программная инженерия: методы и

- технологии разработки информационно-вычислительных систем (ПИИВС-2020). $2020.-\mathrm{C.}\ 77\text{-}83.$
- 82 Li Z. et al. Paraphrase generation with deep reinforcement learning //arXiv preprint arXiv:1711.00279. 2017.
- 83 Zhou J., Bhat S. Paraphrase generation: A survey of the state of the art //Proceedings of the 2021 conference on empirical methods in natural language processing. 2021. C. 5075-5086.
- 84 Aluisio S. et al. Readability assessment for text simplification //Proceedings of the NAACL HLT 2010 fifth workshop on innovative use of NLP for building educational applications. 2010. C. 1-9.
- 85 Siddharthan A. A survey of research on text simplification //ITL-International Journal of Applied Linguistics. − 2014. − T. 165. − №. 2. − C. 259-298.
- 86 Drndarevic B., Saggion H. Reducing text complexity through automatic lexical simplification: an empirical study for Spanish //Procesamiento del lenguaje natural. 2012. T. 49. C. 13-20.
- 87 Сибирцева В.Г., Карпов Н.В. Автоматическая адаптация текстов для электронных учебников. Проблемы и перспективы (на примере русского языка) // Новая русистика. 2014. N 27. C.19-33.
- 88 Русское синтаксическое управление при словесных заменах. О словах с функциями наречия и существительного / А. В. Ниценко, В. Ю. Шелепов, С. А. Большакова, К. С. Ивашко // Проблемы искусственного интеллекта. 2020. № 2(17). С. 46-57.
- 89 О словесных заменах, сохраняющих смысл русского предложения / А. В. Ниценко, В. Ю. Шелепов, С. А. Большакова, К. С. Ивашко // Проблемы искусственного интеллекта. 2020. № 1(16). С. 63-74.
- 90 Fabre B. et al. Neural-Driven Search-Based Paraphrase Generation //Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. 2021. C. 2100-2111.
- 91 Sokolov A., Filimonov D. Neural machine translation for paraphrase generation //arXiv preprint arXiv:2006.14223. 2020.

- 92 Оборнева, И. В. Автоматизация оценки качества восприятия текста // ВЕСТНИК Московского городского педагогического университета, 2015. №2(5). С. 221–233.
- 93 Дмитриева А. Д., Лапошина А. Н., Лебедева М. Ю. Квантитативное исследование стратегий упрощения на материале адаптированных текстов для изучающих РКИ //Компьютерная лингвистика и интеллектуальные технологии: по материалам международной конференции «Диалог». М.: РГГУ. 2021. С. 191-204.
- 94 Church K., Liberman M. The future of computational linguistics: On beyond alchemy //Frontiers in Artificial Intelligence. 2021. T. 4. C. 625341.
- 95 Moore J. D., Wiemer-Hastings P. Discourse in computational linguistics and artificial intelligence //Handbook of discourse processes. 2003. C. 439-486.
- 96 Zhivotova A. A., Berdonosov V. D., Gordin S. A. Mathematical Modeling of the Translation Process and Its Optimization by the Criterion of Quality Maximization // Information Technologies and Intelligent Decision-Making Systems: Communications in Computer and Information Science. 2023. vol. 1821. P. 1–15.
- 97 Чернявская, В. Е. Лингвистика текста. Лингвистика дискурса / В. Е. Чернявская. М.: Общество с ограниченной ответственностью «ФЛИНТА», 2012. 208 с.
- 98 Мизернов И. Ю., Гращенко Л. А. Анализ методов оценки сложности текста // Новые информационные технологии в автоматизированных системах. 2015. N_{\odot} . 18. С. 572-581.
- 99 Солнышкина М. И., Казачкова М. Б., Харькова Е. В. Инструменты измерения сложности текстов на английском языке // Иностр. языки в школе. 2020. N₂. 3. С. 15.
- 100 Hosmer D. W., Lemeshow S. Applied Logistic Regression 2nd edn Wiley & Sons //New York. 2000.
- 101 Животова А. А., Бердоносов В. Д. Стратегия предредактирования исходного текста на основании автоматической оценки сложности задачи перевода для повышения качества машинного перевода узкоспециальных текстов на английский язык // Компьютерная лингвистика и интеллектуальные технологии: по материалам

- ежегодной международной конференции «Диалог». № 22. Доп. том. 2023. С. 1141-1149.
- 102 Каханер Д., Моулер К., Нэш С. Численные методы и программное обеспечение. М.: Мир. 1998.
- 103 MacKay D. J. C. Information theory, inference and learning algorithms. Cambridge university press, 2003.
- 104 Barber D. Bayesian reasoning and machine learning. Cambridge University Press, 2012.
- 105 Bishop C. Pattern recognition and machine learning //Springer google schola. 2006. T. 2. C. 531-537.
- 106 Bengio Y. Practical recommendations for gradient-based training of deep architectures //Neural Networks: Tricks of the Trade: Second Edition. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012. C. 437-478.
- 107 Raffel C. et al. Exploring the limits of transfer learning with a unified text-to-text transformer //The Journal of Machine Learning Research. − 2020. − T. 21. − №. 1. − C. 5485-5551.
- 108 Васильев Д. Д., Пятаева А. В. Использование языковых моделей Т5 для задачи упрощения текста //Программные продукты и системы. -2023. Т. 36. №. 2. С. 228-236.
- 109 Комиссаров В. Н., Цвиллинг М. Я. Лингвистика перевода. Изд. 5, стереотип. М.: URSS, 2020.-176c.
- 110 Животова А. А., Бердоносов В. Д. Автоматизация предредактирования русскоязычных текстов с целью повышения качества их машинного перевода на английский язык // Информационные технологии и высокопроизводительные вычисления: материалы VII Международной науч.- практ. конф., Хабаровск, 11-13 сентября 2023 г. / Редколлегия: Р.В. Намм (отв. редактор) [и др.]. ХФИЦ ДВО РАН: Хабаровск. 2023. С. 88-91.
- 111 Nivre J. et al. Universal Dependencies v2: An evergrowing multilingual tree-bank collection //arXiv preprint arXiv:2004.10643. 2020.

- 112 Животова А. А., Бердоносов В. Д. Машинный перевод корпусов текста для прикладных и исследовательских задач // Молодежь и наука: актуальные проблемы фундаментальных и прикладных исследований: Материалы VI Всероссийской национальной научной конференции молодых учёных, Комсомольск-на-Амуре, 10–14 апреля 2023 года. Том 2. Комсомольск-на-Амуре: Комсомольский-на-Амуре государственный университет, 2023. С. 467-470.
- 113 Пикалев Я. С. Обзор архитектур систем интеллектуальной обработки естественно-языковых текстов // Проблемы искусственного интеллекта. 2020. N_{\odot} 4(19). С. 45-68.
- 114 Findings of the 2019 Conference on Machine Translation (WMT19) / O. Bojar, M. R. Costajuss et al. // Proceedings of the Fourth Conference on Machine Translation. Florence, Italy: Association for Computational Linguistics. 2019. P. 1-61.
- 115 OPUS-MT Building open translation services for the World / J. Tiedemann, S. Thottingal // Proceedings of the 22nd Annual Conference of the European Association for Machine Translation. Lisboa, Portugal: European Association for Machine Translation. 2020.
- 116 Han A. L. F. et al. Language-independent model for machine translation evaluation with reinforced factors //Proceedings of Machine Translation Summit XIV: Posters. 2013. C. 215-222.
- 117 Животова А. А., Бердоносов В. Д. Автоматизированная оценка параметров русскоязычного текста // Молодежь и наука: актуальные проблемы фундаментальных и прикладных исследований: Материалы VI Всероссийской национальной научной конференции молодых учёных, Комсомольск-на-Амуре, 10–14 апреля 2023 года. Том 2. Комсомольск-на-Амуре: Комсомольский-на-Амуре государственный университет, 2023. С. 464-467.
- 118 Saggion H. Automatic text simplification: Synthesis lectures on human language technologies, vol. 10 (1) //California, Morgan & Claypool Publishers. 2017. C. 137.
- 119 Частотный словарь современного русского языка (на материалах Национального корпуса русского языка) / О. Н. Ляшевская, С. А. Шаров. М.: Азбуковник, 2009.

120 Проект Natasha — набор Руthon-библиотек для обработки текстов на естественном русском языке [Электронный ресурс] // Alexander Kukushkin Data Science Laboratory. — 2023. Режим доступа: https://natasha.github.io/.

121 Животова, А. А. Практическое применение вероятностной оценки сложности задачи перевода // Молодежь и наука: актуальные проблемы фундаментальных и прикладных исследований: Материалы VI Всероссийской национальной научной конференции молодых учёных, Комсомольск-на-Амуре, 10–14 апреля 2023 года. Том 2. – Комсомольск-на-Амуре: Комсомольский-на-Амуре государственный университет, 2023. – С. 461-464.

122 Животова А. А., Бердоносов В. Д., Регрессионный анализ корреляции качества машинного перевода и параметров исходного текста // Информатика и системы управления. – 2023. – №2(76). – С.121-133.

123 Бенгфорт Б., Билбро Р., Охеда Т. Прикладной анализ текстовых данных на Python. Машинное обучение и создание приложений обработки естественного языка. — Санкт-Петербург: Питер, 2019.

124 Животова А.А., Бердоносов В.Д. Оптимизационное предредактирование узкоспециальных русскоязычных текстов для их машинного перевода на английский язык // Информационные и математические технологии в науке и управлении. — 2024. — № 2(34). (в издании)

ПРИЛОЖЕНИЕ А

Охранные документы на результаты интеллектуальной деятельности







о государственной регистрации программы для ЭВМ

№ 2023669254

Программа для тренировки языковой модели в решении задач перефразирования русскоязычных технических текстов

Правообладатель: Жиботоба Алена Анатольебна (RU)

Автор(ы): Животова Алена Анатольевна (RU)



遊

班 班 班 班 班

班班班

遊遊

遊遊

遊遊遊

斑

遊

遊

班班班

遊遊

斑

斑斑

撥

斑

斑斑

璐

斑

攝

撥

斑斑

璐

斑斑

璐

璐

Замия № 2023667648 Дата поступления 25 августа 2023 г. Дата государственной регистрации

в Реестре програмы для ЭВМ 12 сентября 2023 г.

Руководитель Федеральной службы по интеллектуальной собственности

досман под теся упитоной подпесио свущения 4/66/кру 10/364/и/66/20/364и/ Видето Время Су и Сергения достатория Су и Сергения

Ю.С. Зубов

蓉

撡

鑗

密路路

安安安安安安安安

密

- 路路

器

密

密路路

路路路路路路路路路

嵡

RICHARILLE OF RANDENSSON



班 班 条 条 条 班

璐

猕

斑

游游游游游游游游

遊遊遊遊

撥

斑斑

班班班班班班班班班班班班

璐

斑斑斑

班班班班班班班

森

СВИДЕТЕЛЬСТВО

о государственной регистрации базы данных

№ 2023623048

Корпус параллельных двуязычных текстов нефтегазовой тематики для тренировки языковых моделей в задачах перефразирования узкоспециальных технических русскоязычных текстов и повышения качества их перевода на английский язык

Правообладатель: Животова Алена Анатольевна (RU)

Автор(м): Животова Алена Анатольевна (RU)



聚聚聚聚聚

撤

遊

推

獥

张张张

遊遊遊

遊

嶽

遊遊遊

遊

撒

遊

遊遊

巌

遊遊

斑斑

機構

極極極

斑斑斑

璨

璨

海海

森森森森

廢

森

Залика № 2023622757 Дата поступления 25 августа 2023 г. Дата государственной регистрации в Ресстре баз даниых 06 сентября 2023 г.

> Руководитель Федеральной службы по интехлектуальной собственности

досован подпис и упатичной водинских сертериях 42/90/дуту — 1945-и 66/20/104а/ Водини Бран Суна Сертевия достинения

Ю.С. Зубов

REMULACINATION RANDIMINO OCI



СВИДЕТЕЛЬСТВО

о государственной регистрации базы данных

№ 2023623022

База данных структурного анализа предложений технических русскоязычных текстов

Правообладатель: Животова Алена Анатольевна (RU)

Автор(м): Животова Алена Анатольевна (RU)



磁 磁 斑 斑 斑 斑 斑

撤

撤

班班班班班班

遊遊遊

遊遊

班班班班班

遊

遊遊遊遊

掛機機

(海海海海海海海

廢

斑

斑

操癌

璨

斑斑斑

森

Залаха № 2023622765 Дата поступления 25 августа 2023 г. Дата государственной регистрации в Ресстре баз данных 01 сентября 2023 г.

> Руководитель Федеральной службы по интеллектуальной собственности

доснице подпости учето на воднесно сертерная 42600 груп — Траний Соловой Соловой Водона повет — Соловой Соровой Соровой Добителнице — Траний Соровой Соровой

Ю.С. Зубов

斑

森





о государственной регистрации программы для ЭВМ

№ 2023668511

Предредактор русскоязычных узкоспециальных текстов для систем машинного перевода

Правообладатель: Жиботоба Алена Анатольебна (RU)

Автор(ы): Животова Алена Анатольевна (RU)



班 班 班 班 班 班

遊

遊

遊 遊 遊 遊

遊

遊 遊 遊

遊 遊

斑 遊 斑 遊 斑 斑

遊

遊 斑 斑

遊 遊 遊 斑

斑 斑 璐 斑 斑 斑

璐 斑

海森

璐 斑

斑

斑 璐 璐

璐

Замая № 2023667667 Дата поступления 25 августа 2023 г. Дата государственной регистрации в Реестре программ для ЭВМ 29 абгуста 2023 г.

> Руководитель Федеральной службы по интеллектуальной собственности

Ю.С. Зубов

蓉

RINGIAGINATION RANDININO OF



СВИДЕТЕЛЬСТВО

о государственной регистрации программы для ЭВМ

№ 2023665348

Программный модуль автоматической очистки сырых данных из памятей переводов САТ для тренировки моделей нейронного машинного перевода

Правообладатель: Жиботоба Алена Анатольебна (RU)

Автор(ы): Животова Алена Анатольевна (RU)



斑斑斑斑斑斑

班班班

班班班班班

班班

遊

遊遊

班班班班班班

遊

遊

遊

遊

遊

斑斑

斑

斑斑

班班班班

斑斑

斑

斑

璨

斑

斑

撥

撥

璐

斑

3amaa N. 2023664365

察察察察察察察察察察察<u>務務務務務務務務務務務務務務務務務</u>

Дата поступления 10 июля 2023 г. Дата государственной регистрации

в Ресстре програмы для ЭВМ 14 июля 2023 г.

Руководитель Федеральной службы по интеллектуальной собственности

досован подтаг на учетности подписано спрационат 4290 берта 2000 Сергения Васите. Вейна Сергения Дейна поста Сергения

Ю.С. Зубов

探 探 斑 斑 斑 迩

泰泰泰泰泰泰

安安安安安安安安安

路路

撡

路路路

密

密路

路路路

磜

密密密

嵡

REMULACION RANDINOSOCI



СВИДЕТЕЛЬСТВО

о государственной регистрации программы для ЭВМ

№ 2023663773

ПРОГРАММНЫЙ КОМПЛЕКС ДЛЯ АНАЛИЗА
РУССКОЯЗЫЧНОГО ТЕКСТА И ВЕРОЯТНОСТНОЙ
ОЦЕНКИ СЛОЖНОСТИ ЗАДАЧИ ЕГО ПЕРЕВОДА
НА АНГЛИЙСКИЙ ЯЗЫК С ИСПОЛЬЗОВАНИЕМ
МОДЕЛИ МАШИННОГО ПЕРЕВОДА

Правообладатель: Жиботоба Алена Анатольебна (RU)

Автор(м): Жиботоба Алена Анатольебна (RU)



班 班 班 班 班 班

班班班班班班

遊

遊遊

遊遊

斑斑

遊

班班

遊

遊遊

遊

斑

遊

遊

遊

斑斑

斑斑斑斑斑

璐

斑

斑

斑

森森

斑斑

璐

璐

璐

3amax N. 2023663302

Дата поступления 27 июня 2023 г. Дата государственной регистрации

в Реестре програмы для ЭВМ 28 июня 2023 г.

Руководитель Федеральной службы по интеллектуальной собственности

доснярн подписня упиточной подписано свучения «Анасти» (1964 годинать досня Видона Ви

Ю.С. Зубов

遊遊

密络路路

撡

经经济经济股

密

泰

撡

密路路

路路路路路路

密

瘀

路路路

森森

密

РОССИЙСКАЯ ФЕЛЕРАЦИЯ

RU

2023617748



ФЕДЕРАЛЬНАЯ СЛУЖБА по интеллектуальной собственности (12) ГОСУДАРСТВЕННАЯ РЕГИСТРАЦИЯ ПРОГРАММЫ ДЛЯ ЭВМ

Номер регистрации (свидетельства): 2023617748

Дата регистрации: 13.04.2023

Номер и дата поступления заявки: 2023613934 28.02.2023

Дата публикации: <u>13.04.2023</u>

Контактные реквизиты: patent@knastu.ru

Авторы:

Животова Алена Анатольевна (RU), Бердоносов Виктор Дмитриевич (RU)

Правообладатель:

Федеральное государственное бюджетное образовательное учреждение высшего образования «Комсомольский-на-Амуре государственный университет» (ФГБОУ ВО «КнАГУ») (RU)

Название программы для ЭВМ:

«ПРОГРАММА ДЛЯ ВЕРОЯТНОСТНОЙ ОЦЕНКИ СЛОЖНОСТИ ПЕРЕВОДЧЕСКОЙ ЗАДАЧИ ДЛЯ СИСТЕМ МАШИННОГО ПЕРЕВОДА»

Реферат:

Программа предназначена для оценки сложности переводческой задачи, обратно пропорциональной вероятности получения перевода, который соответствует требованиям качества пользователя, при его выполнении системой машинного перевода. Программа обучается на массиве данных, включающих исходный текст, эталонный перевод и перевод, выполненный системой машинного перевода. В качестве математической модели оценки вероятности используется модель логистической регрессии. Обученная модель решает задачу классификации, рассчитывая вероятность попадания переведенного текста в класс «качественный перевод» в соответствии с требованиями пользователя. Область применения программы: прикладная компьютерная лингвистика, задачи повышения качества машинного перевода. Тип ЭВМ: IBM PC - совмест. ПК: ОС: Windows, Mac OS X, Linux.

Язык программирования: Python

Объем программы для ЭВМ: 570 КБ





о государственной регистрации программы для ЭВМ

№ 2023614410

ПРЕПРОЦЕССИНГ ТЕКСТОВЫХ ДАННЫХ ДЛЯ ВЗВЕШЕННОЙ ОЦЕНКИ ПАРАМЕТРОВ РУССКОЯЗЫЧНОГО ТЕКСТА

Правообладатель: Жиботоба Алена Анатольебна (RU)

Автор(ы): Жиботоба Алена Анатольебна (RU)



密路路路路

遊

遊 遊 遊 斑 遊 遊 遊

遊 遊 遊

遊 遊

斑 斑 遊 遊 遊 斑

遊

遊

斑 斑 嶽

楽 遊 斑

斑 斑 撥 斑 斑 凝

璐 斑

斑

撥

楽 斑

凝

撥 撥

撥

斑

Замка № 2023612552

дата поступления 13 февраля 2023 г. Дата государственной регистрации

в Реестре программ для ЭВМ 01 марта 2023 г.

Руководитель Федеральной службы по интеллектуальной собственности

допумент подгателя убъртовной подгателю Engineer Bythe Court Copressor

Ю.С. Зубов

容 斑 斑 斑 斑

密





о государственной регистрации программы для ЭВМ

№ 2023613906

ПАРСЕР МАШИННОГО ПЕРЕВОДА УЗКОСПЕЦИАЛЬНЫХ ТЕХНИЧЕСКИХ ТЕКСТОВ С РУССКОГО ЯЗЫКА НА АНГЛИЙСКИЙ ЯЗЫК

Правообпадатель: Животова Алена Анатольевна (RU)

Автор(м): Животова Алена Анатольевна (RU)



斑

班班班班班

班班班

推推

班班班班班班

嶽

遊

班班班

遊遊遊

班班班班班班

斑

斑

斑

斑斑

森

斑斑斑

斑

Замаха № 2023612511

дата поступления 13 февраля 2023 г. Дата государственной регистрации

в Ресстре програмы для ЭВМ 21 фебраля 2023 г.

Руководитель Федеральной службы по интеллектуальной собственности

Appendix Traggers As Profession Incorporate
Copyrigania Silbert No. 64 (Constant Silbert Silbert Copyrigania)
Appendix Copyrigania (Copyrigania)
Applications and Copyrigania (Copyrigania)

Ю.С. Зубов

ПРИЛОЖЕНИЕ Б

Акт о внедрении (использовании) результатов кандидатской диссертационной работы



OOO sArestone repressages +0HAC-Amyps 681003, Pecces, XeSopercore spec, i Kancenore-Care-Amype, op. Revisio, 340, op. 16 Ten Khase +7 (4217) 244-344, -7 (4217) 35-37-25, o real-tendenciffes-aments, www.faccenore-tendencing Tendesicion Agency RAS-Amer Co. LTO 342, Lanin Avenue, Suite 15, Koresonolisk-on-Amer, Krabersesky Kras, Rassin, 681(13) Phone Plac. +7 (4217) 244-344, +7 (4217) 55-37-25, o-mail-tendenciffes-america, www.faccenore.



УТВЕРЖДАЮ

Генеральный директор ООО «Агентство переводов «ФИАС-Амур»

К.Н. Войтик

613/10T = 10 = WIENES 2023 T.

AKT

о внедрении (использовании) результатов

кандидатской диссертационной работы

Животовой Алены Анатольевны

Комиссия в составе: председатель Кира Николаевна Войтик, члены комиссин: Галина Юрьевна Мастевная, Сергей Петрович Светлаков, составили настоящий акт о том, что результаты диссертационной работы Животовой Алены Анатошевны «Математическая модель, алгоритмы и программный комплекс для повышения качества машинного перевода узкоспециальных технических техстов на английский язык», представленной на сонежание ученой степени кандидата технических наук, использованы в исследовательской и практической деятельности ООО «Агентство переводов «ФИАС-Амур» при разработке контура автоматизации переводческой деятельности и создания подсистемы поддержки принятия решений для менеджеров проектов и управления рисками, саятанными с компетенцией исполнителей, в том числе при использовании систем машинного перевода, а также создания баз данных для тренировки языковых моделей.

В ходе исследования Животовой Аленой Анатольевной были разработаны методика и алгоритм вероитностной оценки сложности задачи перевода, методика оптимизационного предредактирования русскоязычных текстов. Разработанные алгоритмы реализованые в программном комплексе для оценки сложности задачи перевода русскоязычных текстов на английский язык и оптимизационного предредактирования и прошли валидацию на корпусах двуязычных текстов (языковые пары русский-английский и английский-русский) объемом — 40 000 страниц, собразных на базе Translation Memory из сред автоматизации перевода, применяемых ООО «Агентство переводов «ФИАС-Амур», получены свидетельства о регистрации программ ЭВМ и баз даниых.

Завелючение: Использование полученных результатов исследования и внедрение программного комплекса в работу агентства переводов позволило повысить эффективность использования систем машининого перевода и производительность труда редакторов переводов, а также оптимизировать заграты на оказание услуг перевода узкоспециальной технической документации.

Председатель комиссии:

Генеральный директор

Члены комиссии: Гланный бухгалтер

Руховодитель службы автоматизации

Foremen

Войтик К.Н.

DOMING N.I.

Мастевная Г.Ю. Светлаков С.П.