

УДК 004.9:519.768

DOI:10.25729/ESI.2024.34.2.016

Оптимизационное предредактирование узкоспециальных русскоязычных текстов для их машинного перевода на английский язык

Животова Алена Анатольевна, Бердонос Виктор Дмитриевич

Комсомольский-на-Амуре государственный университет,
Россия, Комсомольск-на-Амуре, *zhivotova.aa@gmail.com*

Аннотация. В статье исследуется процесс перевода текстов, а именно метод оптимизационного предредактирования, как способ повышения качества машинного перевода на английский язык при работе с русскоязычными узкоспециальными текстами. Авторы рассматривают математическую модель процесса перевода и постановку задачи машинного перевода, предлагают новую теорию вероятностной оценки сложности задачи перевода, приводят постановку и решение задачи оптимизационного предредактирования, описывают методику подготовки данных для обучения модели автоматического оптимизационного предредактирования. В результате исследования реализован программный комплекс оптимизационного предредактирования русскоязычных текстов. При разработке программного комплекса использованы ресурсы Центра коллективного пользования научным оборудованием «Центр обработки и хранения научных данных ДВО РАН». Данные для обучения и валидации моделей предоставлены ООО «Агентство переводов “ФИАС-Амур”». Тестирование программного комплекса показало эффективность предложенных методик для повышения качества машинного перевода узкоспециальных русскоязычных текстов на английский язык.

Ключевые слова: машинный перевод, оптимизационное предредактирование, сложность задачи перевода, качество перевода

Цитирование: Животова А.А. Оптимизационное предредактирование узкоспециальных русскоязычных текстов для их машинного перевода на английский язык / А.А. Животова, В.Д. Бердонос // Информационные и математические технологии в науке и управлении. – 2024. – № 2(34). – С. 169-182. – DOI:10.25729/ESI.2024.34.2.016.

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. Несмотря на выдающиеся прорывы нейросетевых, гибридных и больших языковых моделей машинного перевода (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 [1, 2]. Эту проблему активно освещают зарубежные исследователи A. Lear, C. Quinci, C. Canfora, A. Ottman, D. Kenny, P. Sanchez-Gijon. Предоставляя пользователю средства обработки текста на языке, носителем которого он является, на любом из этапов перевода, можно повысить его качество. Именно на этот принцип опирается концепция интерактивного перевода, широко описанная в литературе [3]. Одним из направлений интерактивного перевода является перевод с предредактированием, когда исходный текст предварительно редактируется с целью его адаптации для более легкого «понимания» системой МП.

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 [4]; 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 пред-

редактированию [5, 6, 7]; L. O'Brien, D. Folaron, W. Aziz, M. Toledo – в области контролируемых и упрощенных языков [7]. Среди российских авторов и для русского языка данная тема освещена незначительно, однако известны работы И. В. Оборневой [8], А.Д. Дмитриевой, А. Н. Лапошиной, М. Ю. Лебедевой [9] и др. в области оценки восприятия текста и упрощения русскоязычных текстов в соответствии с квалификацией реципиента.

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

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

В основе проведенного исследования лежит применение апробированных математических методов, включая теорию множеств, численные методы оптимизации, такие, как метод наименьших квадратов и метод градиентного спуска, статистические методы, в том числе метод максимального правдоподобия.

В работе 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; 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; предложена новая архитектура и реализован программный комплекс для повышения качества МП 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 МП текстов с русского языка на английский язык.

1. Математическая модель процесса перевода. Перевод – передача смысла текста на языке оригинала в соответствии с языковыми и культурными традициями языка перевода согласно требованиям конечного реципиента.

Процесс перевода формально может быть представлен совокупностью множеств и операций над этими множествами. В дальнейшем слово «множество» опускается для удобства повествования. Пусть $txt_{iТХТ}$ – исходный текст, характеризующийся множеством выраженных смысловых единиц $СМ_i$, $яз_{вх}$ – язык исходного текста, $яз_{вых}$ – язык перевода, $ДП_iДП$ – предметная область (домен приложения) исходного текста, $ТР|txt_{iТХТ} : \exists ОК|ТР \in R$ – требования к переводу исходного текста на язык перевода, такие, что существует ОК – аналитически вычисляемая оценка качества перевода, соответствующая критериям, зафиксированным в требованиях к переводу, перевод выполняется переводчиком $пер_{iПЕР}$, обладающим компетенцией $К_{перЛПЕР}$ и специализацией $С_{перЛПЕР, дпДП}$, тогда качественный перевод – это сгенерированный текст $txt_{jТХТ}$ на языке $яз_{вых}$, характеризующийся множеством выраженных смысловых единиц $СМ_j$: $(СМ_j \rightarrow СМ_i)$. В общем виде соотношения множеств функция перевода $F_{пер}$ может быть описана как

$$F_{пер} : (txt_{iТХТ}, К_{перЛПЕР}, С_{перЛПЕР, дпДП}) \rightarrow txt_{jТХТ} | (СМ_j \rightarrow СМ_i, ТР|txt_{iТХТ}). \quad (1)$$

Обобщенно модель процесса перевода представлена на рисунке 1. Подробное описание полной модели было опубликовано авторами ранее [10].

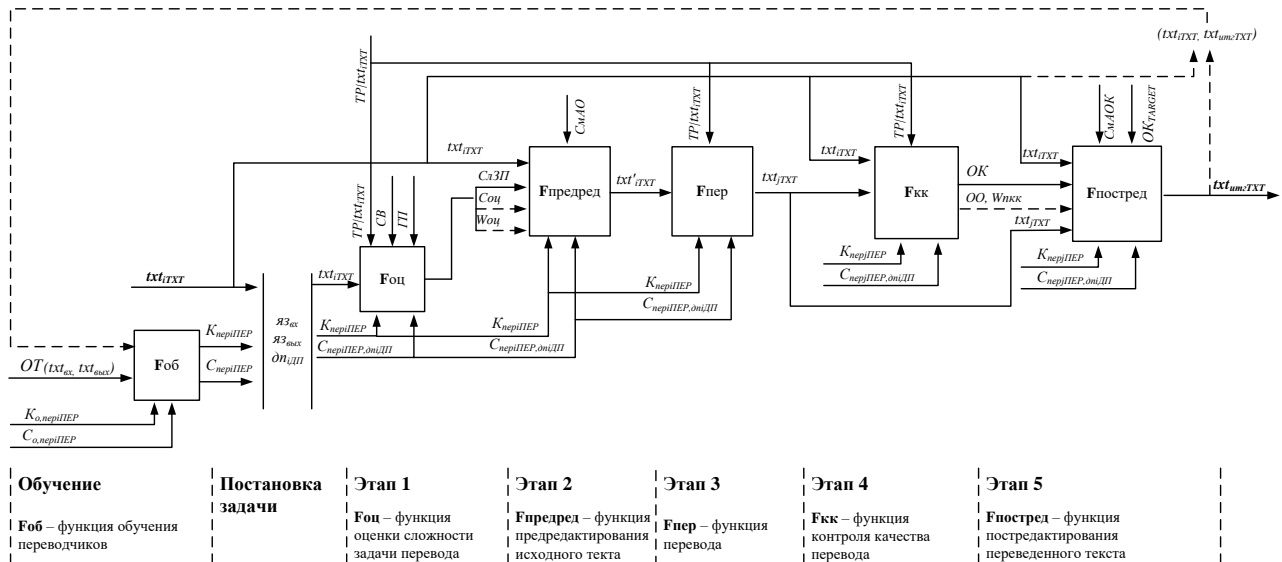


Рис. 1. Модель процесса перевода

Математически, для систем МП задачу перевода можно формализовать через метод максимизации функции правдоподобия [11] следующим образом. Опишем условия задачи, пусть:

$txt_{iТХТ}$ – исходный текст на языке $ЯЗ_{ВХ}$.

$txt_{jТХТ}$ – переведенный текст на языке $ЯЗ_{ВЫХ}$.

TXT_{itrg} – множество всех возможных вариантов перевода текста $txt_{iТХТ}$ на язык $ЯЗ_{ВЫХ}$:

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

где j – общее число вариантов перевода, $txt_{jТХТ} \in TXT_{itrg}$.

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

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

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

Каждому варианту перевода соответствует одна оценка качества перевода, то есть множества TXT_{itrg} и OK_{itrg} биективны: $TXT_{itrg} \leftrightarrow OK_{itrg}$.

$[minOK; maxOK]$ – диапазон значений оценок качества перевода OK_{itrg} .

$OK_{доп}$ – минимально допустимое значение критерия «Высокая оценка качества перевода» при допущении, что чем выше значение OK_j , тем лучше.

$KП_i$ – нечёткое подмножество множества OK_{itrg} , определяющее принадлежность элементов множества OK_{itrg} классу «Высокая оценка качества перевода текста $txt_{iТХТ}$ »:

$$KП_i = \{(OK, \mu_{KП_i}(OK)) | OK \in OK_{itrg}\}.$$

$\mu_{KП_i}(OK)$ – функция принадлежности, указывающая, в какой степени текст txt с оценкой OK принадлежит нечеткому множеству $KП_i$.

$\mu_{KП_i}(OK) \in [0; 1]$ и имеет вид логистической кривой:

$$\mu_{KП_i}(OK) = \frac{1}{1 + e^{-\left(\frac{OK - OK_{доп}}{maxOK - OK_{доп}}\right)2\pi}}. \tag{2}$$

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

$$F_{МП}(\theta, \mu_{KП_i}(OK)) = \ln P_{\theta}(\max \mu_{KП_i}(OK)) \rightarrow \max_{\theta}, \tag{3}$$

где θ – параметры системы МП из множества исполнителей перевода, или переводчиков: $пер_{iПЕР} \in ПЕР$, максимизирующие вероятность P получить максимальное значение функции принадлежности $\mu_{KП_i}(OK)$.

Решение поставленной задачи лежит в области оптимизации и совершенствования алгоритмов генерации переведенного текста.

В результате теоретического моделирования определено, что 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 оптимизации предредактирования исходного текста с целью повышения вероятности получения качественного перевода.

2. Сложность задачи перевода. Для разработки стратегии и методики оптимизационного предредактирования требуется определить критерий оптимизации исходного текста. В качестве такого критерия была выбрана оценка сложности задачи перевода. При оценке сложности задачи перевода переводчик обращает внимание на неизвестные ему слова и сочетания слов на языке $яз_{вх}$, для которых он не может идентифицировать значение смысловой единицы, либо смысловые единицы, для которых он не может найти аналог на языке перевода $яз_{вых}$ среди известных ему слов и сочетаний слов. Множества свойств и параметров исходного текста $СВ|txt_{ITXT}$ и $ГП|txt_{ITXT}$, и то, обладает ли переводчик достаточной компетенцией $К_{перЛЕР}$ относительно языков $яз_{вх}$ и $яз_{вых}$ и специализацией $С_{перЛЕР, дпДП}$, т.е. навыками описания семантических единиц на языке перевода в рамках заданной предметной области исходного текста, определяет вероятность создания переводчиком переведенного текста на таком уровне качества, который определяется требованиями $ТР|txt_{ITXT}$.

Оценки сложности задачи перевода включает следующие шаги:

Шаг 1. Исходя из домена приложения текста $дп_{дП}$, формируется множество оценок текста $ОЦ = СВ \cup ГП$. Свойства текста СВ условно можно разделить на группы признаков: общие (количество символов/слов/строк и т.д., стиль, язык, домен приложения и пр.); лексические (процент покрытия текста лексическими минимумами, частотными списками и др.); морфологические (количество различных частей речи и грамматических форм); синтаксические (глубина глагольных и именных групп, связи между глаголами в предложениях); признаки, основанные на базовых подсчетах (средняя длина слов и предложений и пр.). Совокупность свойств и признаков определяет главные параметры текста ГП, к которым относится целостность, связность, удобочитаемость, сложность и другие.

Шаг 2. Для каждого значения $св_{iСВ}, гп_{iГП} \in ОЦ$, на основе требований к переводу $ТР|txt_{ITXT}$, компетенций переводчика относительно языковой пары $К_{перЛЕР}$ и специализации переводчика относительно домена приложения текста $С_{перЛЕР, дпДП}$ формируется значение значимости $w_{оцk}$, множество нормированных значений $w_{оцk}$ значимости формирует матрицу значимости оценок сложности $W_{оц}$ 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 $i = \overline{1, N}$ формируется матрица оценок фрагмента исходного текста $С_{оци}$ размерностью $1 \times k$, где k – общее число оценок.

Шаг 4. На основании оценок $С_{оци}$ и значимости $W_{оц}$ формируется уравнение поиска теоретического результирующего фактора, т.е. качества перевода $\widehat{КП}$:

$$\widehat{КП}_i = w_0 + w_{оц1}C_{оци1} + w_{оц2}C_{оци2} + \dots + w_{оцk}C_{оциk}. \quad (4)$$

Для системы МП веса значимости оценок рассчитываются на основании тренировочных данных с использованием численного метода наименьших квадратов [12], при котором минимизируется сумма квадратов отклонений эмпирических значений результирующего признака от теоретических, полученных по уравнению (4):

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

$$S(w) = \sum_{i=1}^R (OK_i - w_0 + w_{оц1}C_{оци1} + w_{оц2}C_{оци2} + \dots + w_{оцk}C_{оциk})^2 \rightarrow \min, \quad (5)$$

где R – объем тренировочной выборки.

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

$$\sum_{i=1}^R (OK_i - \widehat{OK}_i(C_{оци}, w)) \frac{\partial \widehat{OK}_i(C_{оци}, w)}{\partial w} = 0. \quad (6)$$

Получаем систему нормальных уравнений с k неизвестными:

$$\begin{cases} \sum OK = R w_0 + w_{оц1} \sum C_{оци1} + w_{оц2} \sum C_{оци2} + \dots + w_{оцk} \sum C_{оциk} \\ \sum OK \cdot C_{оци1} = w_0 \sum C_{оци1} + w_{оц1} \sum C_{оци1}^2 + w_{оц2} \sum C_{оци2} C_{оци1} + \dots + w_{оцk} \sum C_{оциk} C_{оци1} \\ \dots \\ \sum OK \cdot C_{оциk} = w_0 \sum C_{оциk} + w_{оц1} \sum C_{оци1} C_{оциk} + w_{оц2} \sum C_{оци2} C_{оциk} + \dots + w_{оцk} \sum C_{оциk}^2 \end{cases}$$

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

$$W_{оц} = (C_{оц}^T \cdot C_{оц})^{-1} \cdot C_{оц}^T \cdot OK = \left(\frac{1}{R} C_{оц}^T \cdot C_{оц}\right)^{-1} \frac{1}{R} C_{оц}^T \cdot OK. \quad (7)$$

Шаг 5. Для каждого i -го фрагмента текста рассчитывается вероятность получения переведенного текста на таком уровне качества, который определяется требованиями $TP|txt_{iТХТ}$, применив к уравнению (4) логит-преобразование:

$$p_i = \frac{1}{1 + e^{-k \Pi_i}}. \quad (8)$$

Шаг 6. Сложность задачи перевода i -го фрагмента текста оценивается по формуле:

$$СлЗП_i = \frac{1}{p_i}. \quad (9)$$

Шаг 7. Результирующая сложность задачи перевода текста – это наибольшее значение сложности задачи перевода $СлЗП_i$ среди N фрагментов исходного текста, то есть

$$СлЗП_{txt_{iТХТ}} = \max СлЗП_i \quad (10)$$

В зависимости от значения $СлЗП_{txt_{iТХТ}}$ определяется стратегия дальнейшей обработки текста, в том числе необходимость применять оптимизационное предредактирование.

3. Задача оптимизационного предредактирования. Задача оптимизационного предредактирования состоит в том, чтобы максимизировать правдоподобие, то есть вероятность того, что при параметрах Ψ предредактора, текст $txt'_{iТХТ}$ на языке $яз_{вх}$ будет эквивалентен $txt_{iТХТ}$ по смыслу, понятен системе МП $пер_{iПЕР}$ и оценка качества OK_j перевода $txt_{jТХТ}$ относительно $txt_{iТХТ}$ при генерации перевода из $txt'_{iТХТ}$ будет максимальной. Далее опишем задачу более подробно.

Опишем условия задачи, пусть:

- 1) $txt'_{iТХТ}$ – текст на языке $яз_{вх}$, созданный системой автоматического оптимизационного предредактирования, такой, при котором $СМ'_i \rightarrow СМ_i$ и $txt'_{iТХТ} \neq txt_{iТХТ}$, где $СМ$ – смысл или упорядоченный набор семантических единиц, описываемый текстом: $СМ = \{(ce_0, ce_1, \dots, ce_{ncm}) : ce_{icm} \in CE\}$, $СМ'_i$ и $СМ_i$ – смыслы $txt'_{iТХТ}$ и $txt_{iТХТ}$, соответственно;
- 2) $ТХТ_{isrc}$ – множество всех возможных вариантов предредактированного текста, т.е. выражения смысла $СМ_i$ текста $txt_{iТХТ}$ на языке $яз_{вх}$;
- 3) $ТХТ_{isrc} = \{txt_0, txt_1, \dots, txt_k\}$,
- 4) где k – общее число вариантов предредактированного текста, причем $txt_{iТХТ}, txt'_{iТХТ} \in ТХТ_{isrc}$;
- 5) $мСлЗП_{isrc}$ – множество оценок сложности задачи перевода вариантов предредактирования текста $txt_{iТХТ}$ в соответствии с компетентностью $К_{перПЕР}$ и специализацией $С_{перПЕР, дпдП}$ системы МП $пер_{iПЕР}$ для всех возможных вариантов предредактированного текста $ТХТ_{isrc}$;
- 6) $мСлЗП_{isrc} = \{СлЗП_0, СлЗП_1, \dots, СлЗП_k\}$;

- 7) Каждому варианту предредактированного текста соответствует одна оценка сложности задачи перевода для системы МП $\text{пер}_{i\text{ПЕР}}$, то есть множества $TXT_{i\text{SRC}}$ и $\text{мСлЗП}_{i\text{SRC}}$ биективны: $TXT_{i\text{SRC}} \leftrightarrow \text{мСлЗП}_{i\text{SRC}}$;
- 8) $\text{СлЗП}_k \in \text{мСлЗП}_{i\text{SRC}}$ – оценка сложности задачи перевода варианта предредактированного текста $\text{txt}'_{i\text{TXT}}$ для системы МП $\text{пер}_{i\text{ПЕР}}$;
- 9) $[\text{minСлЗП}; \text{maxСлЗП}]$ – диапазон значений оценок сложности задачи перевода $\text{мСлЗП}_{i\text{SRC}}$;
- 10) $\text{СлЗП}_{\text{доп}}$ – максимально допустимое значение критерия «Низкая сложность задачи перевода» при допущении, что чем ниже значение СлЗП_k , тем лучше;
- 11) нСлЗП_i – нечёткое подмножество множества $\text{мСлЗП}_{i\text{SRC}}$, определяющее принадлежность элементов множества $\text{мСлЗП}_{i\text{SRC}}$ и соответствующих элементов множества $TXT_{i\text{SRC}}$ классу «Низкая сложность задачи перевода»;
- 12) $\text{нСлЗП}_i = \{(\text{СлЗП}, \mu_{\text{нСлЗП}_i}(\text{СлЗП})) | \text{СлЗП} \in \text{мСлЗП}_{i\text{SRC}}\}$;
- 13) $\mu_{\text{нСлЗП}_i}(\text{СлЗП})$ – функция принадлежности, указывающая в какой степени текст txt с оценкой СлЗП принадлежит нечеткому множеству нСлЗП_i ;
- 14) $\mu_{\text{нСлЗП}_i}(\text{СлЗП}) \in [0; 1]$ и имеет вид логистической кривой:

$$\mu_{\text{нСлЗП}_i}(\text{СлЗП}) = \frac{1}{1 + e^{\frac{\text{СлЗП} - \text{СлЗП}_{\text{доп}}}{(\text{minСлЗП} - \text{СлЗП}_{\text{доп}})^{2\pi}}}}. \quad (11)$$

Требуется максимизировать правдоподобие сгенерированного системой оптимизационного предредактора текста $\text{txt}'_{i\text{TXT}}$, то есть вероятность того, что $\text{txt}'_{i\text{TXT}}$ примет такое значение, при котором $\mu_{\text{нСлЗП}_i}(\text{СлЗП})$ будет максимальна.

В дискретном случае функция правдоподобия $F_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ – вероятность выборки $\mu_{\text{нСлЗП}_i}(\text{СлЗП}) = \{\mu_0, \mu_1, \dots, \mu_l\}$ в рассматриваемой серии экспериментов будет равняться $\{\max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_0, \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_1, \dots, \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_l\}$. Эта вероятность меняется в зависимости от Ψ :

$$F_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП})) = \prod_{l=1}^L F_{\text{АОПР}}(\mu_{\text{нСлЗП}_i}(\text{СлЗП})_l) = P_{\Psi}(\mu_0 = \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_0) \cdot \dots \cdot P_{\Psi}(\mu_l = \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_l) = P_{\Psi}(\mu_0 = \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_0, \dots, \mu_l = \max \mu_{\text{нСлЗП}_i}(\text{СлЗП})_l), \quad (12)$$

где l – номер экземпляра в обучающей выборке объемом L .

Тогда логарифмическая функция правдоподобия автоматического оптимизационного предредактирования $F_{\text{АОПР}}$ имеет вид:

$$L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП})) = \ln P_{\Psi}(\max \mu_{\text{нСлЗП}_i}(\text{СлЗП})), \quad (13)$$

где Ψ – параметры системы автоматического оптимизационного предредактирования, максимизирующие вероятность P получить максимальное значение функции принадлежности $\mu_{\text{нСлЗП}_i}(\text{СлЗП})$.

Поскольку $\ln(y)$ монотонна, то точки максимума $F_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ и $L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ совпадают, и оценкой максимального правдоподобия можно назвать точку максимума функции $L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ по Ψ . Задача оптимизации, таким образом, заключается в поиске оценки максимального правдоподобия $\hat{\Psi}$ вектора параметров Ψ , или:

$$\hat{\Psi} = \arg \max_{\Psi} L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП})) \quad (14)$$

Решение поставленной задачи оптимизации выполняется методом градиентного спуска (подъема) [13]. Для этого необходимо найти градиент логарифмической функции правдоподобия $L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ – вектор, который показывает направление возрастания функции.

Учитывая, что Ψ – вектор параметров системы автоматического оптимизационного предредактирования и $\Psi = \{\psi_1, \psi_2, \dots, \psi_m\}$, где m – количество параметров модели, градиент функции $L_{\text{АОПР}}(\Psi, \mu_{\text{нСлЗП}_i}(\text{СлЗП}))$ может быть найден по формуле:

$$\nabla L_{\text{АОПР}}(\Psi) = (\partial L_{\text{АОПР}} / \partial \psi_1, \partial L_{\text{АОПР}} / \partial \psi_2, \dots, \partial L_{\text{АОПР}} / \partial \psi_m), \quad (15)$$

где $\partial L_{\text{АОПР}} / \partial \psi_m$ – частная производная функции правдоподобия по m -ному параметру.

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

$$\Psi^{[s+1]} = \Psi^{[s]} + \alpha \cdot \nabla L_{\text{АОПР}}(\Psi^{[s]}), \quad (16)$$

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 $\nabla L_{\text{АОПР}}(\Psi)$:

$$\|\nabla L_{\text{АОПР}}(\Psi)\| = \sqrt{\left(\frac{\partial L_{\text{АОПР}}}{\partial \psi_1}\right)^2 + \left(\frac{\partial L_{\text{АОПР}}}{\partial \psi_2}\right)^2 + \dots + \left(\frac{\partial L_{\text{АОПР}}}{\partial \psi_m}\right)^2}. \quad (17)$$

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:

$$\|\nabla L_{\text{АОПР}}(\Psi^{[s]})\| \leq \varepsilon. \quad (18)$$

Используя описанную математическую модель, создадим модель автоматического оптимизационного предредактирования.

4. Обучающие данные для модели оценки сложности задачи перевода. Для обучения модели оценки сложности задачи перевода заданного текста $txt_{\text{ТХТ}}$ с языка $\text{яз}_{\text{ВХ}}$ на язык $\text{яз}_{\text{ВЫХ}}$ переводчиком $\text{пер}_{\text{ИПЕР}}$ в соответствии с формализованными требованиями к переводу $\text{TR}|txt_{\text{ТХТ}}$ требуются исходные данные в виде корпуса параллельных текстов следующей структуры:

TranslatorExpCor: [*src*; *trg*; *ref*], где *src* – это оригинал, т.е. текст на языке $\text{яз}_{\text{ВХ}}$, *trg* – это перевод (текст на языке $\text{яз}_{\text{ВЫХ}}$), выполненный переводчиком или системой МП $\text{пер}_{\text{ИПЕР}}$; *ref* – это контрольный перевод (текст на языке $\text{яз}_{\text{ВЫХ}}$), т.е. проверенный эталон.

Путем обработки корпуса *TranslatorExpCor* модулями оценки качества МП и препроцессинга текстовых данных для взвешенной оценки параметров русскоязычного текста формируется база данных структурного анализа предложений, содержащая расчет вещественных параметров по морфологическим, синтаксическим, лексическим и прочим признакам [14].

5. Обучающие данные для модели автоматического оптимизационного предредактирования. Для обучения модели, которая будет преобразовывать текст на языке $\text{яз}_{\text{ВХ}}$ в текст требуемой структуры для повышения качества перевода, необходимо создать корпус тренировочных текстов в паре $\text{яз}_{\text{ВХ}} - \text{яз}_{\text{ВХ}}$. Для оптимизации временных затрат на подготовку исходных данных для тренировки модели предредактирования текста предлагается методика с использованием обратного перевода для генерирования эталонного предредактированного текста.

Структура параллельного корпуса исходных данных: *RefCor*: [*src_ref*; *tgt_ref*], где *src_ref* – это оригинал, т.е. текст на языке $\text{яз}_{\text{ВХ}}$, *tgt_ref* – это перевод (текст на языке $\text{яз}_{\text{ВЫХ}}$).

Методика сбора корпуса обучающих текстов для модели оптимизационного предредактирования включает следующие шаги:

1. Настраиваем системы МП $\text{MT}: \text{tgt} - \text{src}$, $\text{MT}: \text{src} - \text{tgt}$.

При помощи системы $\text{MT}: \text{tgt} - \text{src}$ переводим текст *tgt_ref* на язык $\text{яз}_{\text{ВХ}}$, получим массив текстовых данных *pre_src* (массив условно предредактированных текстов).

При помощи системы $\text{MT}: \text{src} - \text{tgt}$ переводим текст *src_ref* на язык $\text{яз}_{\text{ВЫХ}}$, получим массив текстовых данных *tgt1* (*src_ref* \rightarrow *tgt1*).

При помощи системы $\text{MT}: \text{src} - \text{tgt}$ переводим текст *pre_src* на язык $\text{яз}_{\text{ВЫХ}}$, получим массив текстовых данных *tgt2* (*pre_src* \rightarrow *tgt2*).

Оцениваем качество выполненного перевода на язык $\text{яз}_{\text{вых}}$ tgt1 и tgt2 относительно эталона tgt_ref , получаем массивы оценок $\text{QC_score}(\text{tgt1})$ и $\text{QC_score}(\text{tgt2})$.

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 $\Delta\text{QC_score}_i$ является условно значимой d_{max} для выбранного типа оценки:

$$\text{TrainCor} = \{(\text{src_ref}_i; \text{pre_src}_i): \exists (\text{tgt1}_i, \text{tgt2}_i) | \text{QC_score}(\text{tgt2}_i) > \text{QC_score}(\text{tgt1}_i) \ \& \ \Delta\text{QC_score}_i \geq d_{\text{max}}\} \quad (19)$$

Полученный корпус обучающих текстов TrainCor объемом выборки L будем использовать для обучения языковой модели LM: $\text{src} - \text{pre_src}$ для решения задачи автоматического оптимизационного редактирования текстов на языке $\text{яз}_{\text{вх}}$.

6. Оценка качества выполненного перевода. Для реализации программного комплекса определим критерий качества МП QC_score . Традиционно, для оценки качества перевода используются алгоритмы, которые сравнивают выполненный перевод с одним или несколькими эталонными переводами при помощи числовой метрики. В рамках исследования была выбрана метрика $h\text{LEPOR}$ [15], которая имеет наивысший балл корреляции Пирсона с человеческими суждениями по языковой паре английский-русский.

$h\text{LEPOR}$ – гармоническое среднее между штрафом за длину перевода в сравнении с эталоном, штрафом за различие в позициях и степенью перекрытия n -грамм; рассчитывается в диапазоне от 0 до 1, где 0 – полное несовпадение; 1 – полное совпадение:

$$h\text{LEPOR} = \frac{(\text{LP} \cdot \text{PP} \cdot \text{OR})^{\frac{1}{3}}}{\left(\frac{1}{\text{LP}} + \frac{1}{\text{PP}} + \frac{1}{\text{OR}}\right)^{\frac{1}{3}}} \quad (20)$$

где LP — штраф за длину, который учитывает различие в длине между переводом и эталонным переводом; PP — штраф за различие в позициях, который учитывает различие в расположении слов в переводе и эталоне; OR — степень перекрытия, которая учитывает совпадение n -грамм в переводе и эталоне.

Для расчета каждого из компонентов используются следующие формулы:

$$\text{LP} = f(x) = \begin{cases} e^{1 - \frac{\text{ref}}{\text{cand}}}, & \text{cand} > \text{ref} \\ 1, & \text{cand} \leq \text{ref} \end{cases} \quad (21)$$

где cand — длина переведенного текста, а ref — длина эталонного перевода;

$$\text{PP} = \frac{1}{1 + p_{\text{diff}}}, \quad (22)$$

где p_{diff} — среднее арифметическое различий позиций одинаковых n -грамм в переводе и эталоне;

$$\text{OR} = \text{BP} \cdot \text{ng}_{\text{prec}} \cdot (1 - \text{pen}), \quad (23)$$

где BP — штраф за слишком короткий перевод, pen — дополнительный штраф за ошибки, ng_{prec} — точность n -грамм, которая рассчитывается как отношение количества совпадающих n -грамм в переводе и эталоне к общему количеству n -грамм в переводе.

6. Программный комплекс оптимизационного предредактирования русскоязычных текстов. Опираясь на описанные математические модели и методики, была разработана архитектура программного комплекса, который состоит из трех основных подсистем: подсистемы тренировки языковой модели, подсистемы оценки сложности задачи перевода, подсистемы оптимизационного предредактирования русскоязычного текста и генерации МП на английский язык. При реализации приняты 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 CAT, процесс накопления тренировочных корпусов, включающих исходный текст, перевод, выполненный системой МП и проверенный перевод, утвержденный редактором, происходит автоматически в режиме реального времени.

Схематично архитектура программного комплекса представлена на рисунке 2.

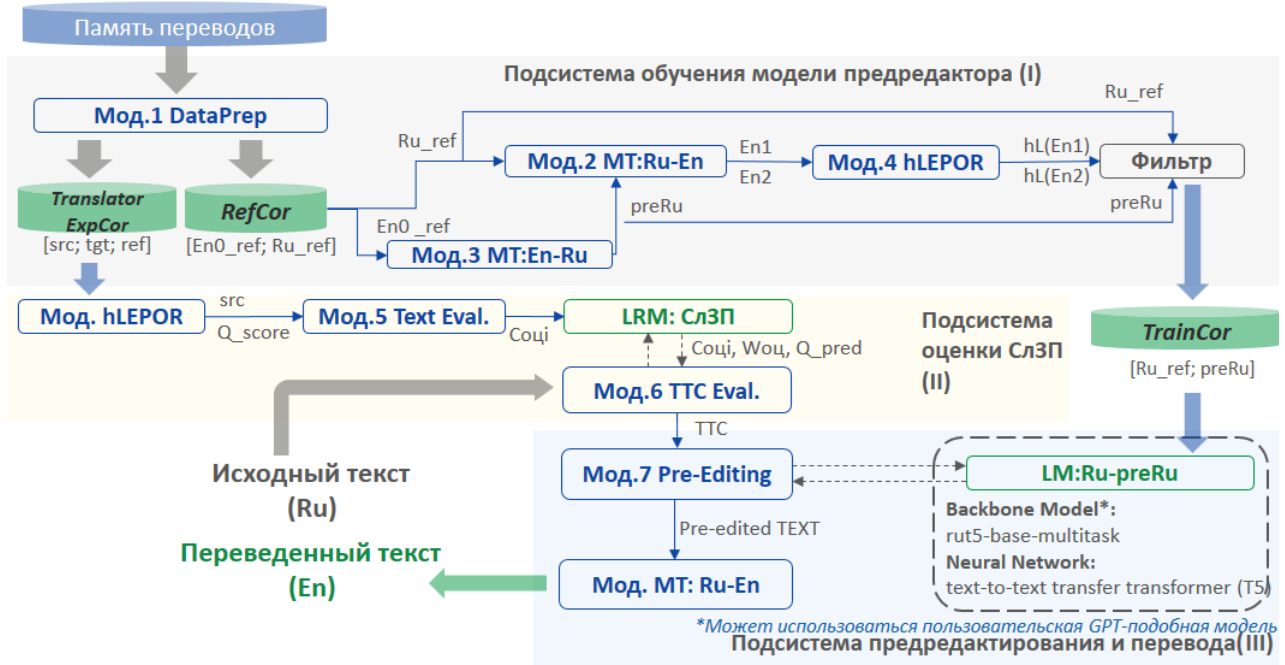


Рис. 2. Архитектура программного комплекса оптимизационного предредактирования русскоязычных текстов

Подсистема тренировки языковой модели оптимизационного предредактирования русскоязычных текстов (I) состоит из пяти программных компонентов: *Мод.1 DataPrep* – модуля подготовки «сырых данных», полученных из памяти 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 MT:Ru-En*, *Мод.3 MT: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*; *Мод.5 Text Eval.* – препроцессора для взвешенной оценки свойств русскоязычного текста, включая морфологические, синтаксические, лексические и другие, всего 96 параметров; *Мод.6 TTC Eval.* – модуля оценки сложности задачи перевода с применением модели логистической регрессии *LRM:СлЗП*, 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* (T5), представленной командой Google в 2020 г., дообученной на корпусе параллельных русскоязычных текстов подзадачу перефразирования [16]; модуля генерации МП с русского языка на английский язык на основе модели *Helsinki-NLP/opus-mt-ru-en*.

7. Тестирование программного комплекса. Программный комплекс [17] 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», созданного и функционирующего на базе ВЦ ДВО РАН – обособленном подразделении ХФИЦ ДВО РАН. Тип используемой ЭВМ: компьютер с архитектурой x86, x86_64.

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 ~60 000 ст. стр. текста (1 ст. стр. = 1800 знаков с пробелами). Из предоставленных данных сформированы корпуса RefCor объемом ~140 000 экземпляров и TranslatorExpCor объемом ~90 000 экземпляров. После обработки корпуса RefCor в подсистеме тренировки языковой модели оптимизационного предредактирования русскоязычных текстов в корпус TrainCor вошло ~84 000 экземпляров. Тестовая выборка для оценки работы системы TestCor составила ~17 000 экземпляров.

Тестирование программного комплекса включало следующие этапы:

1. Машинный перевод тестовой выборки и оценка его качества алгоритмом *hLEPOR* → получение оценки *hLEPOR(En1)*.
2. Оценка сложности задачи перевода тестовой выборки *СлЗП(Ru)* → применение оптимизационного предредактирования только к тем семплам, для которых $СлЗП(Ru) > СлЗП_{доп}$ → оценка *СлЗП* после оптимизационного предредактирования *СлЗП(preRu)* → машинный перевод → получение оценки *hLEPOR(TTC/PE)*.
3. Сравнение и анализ полученных результатов.

В рамках тестирования принято, что минимально допустимая сложность задачи перевода $СлЗП_{доп} = 1,43$. Так как сложность задачи перевода обратно пропорциональна вероятности получения перевода требуемого качества, при $СлЗП_{доп} = 1,43$ данная вероятность составляет 0,7.

С использованием оценки сложности задачи перевода из тестовой выборки было отобрано 5440 экземпляров для оптимизационного предредактирования (32,58% тестовой выборки).

Примеры оптимизационного предредактирования и его влияния на сложность задачи перевода для системы МП представлены в таблице 1. Результаты перевода текстов после оптимизационного предредактирования на английский язык представлены в таблице 2.

Таблица 1. Оптимизационное предредактирование русскоязычных текстов

Экз.	Исходный текст (Ru_ref)	Текст после предредактирования (preRU)	$\Delta СлЗП$
1	Высоковольтные испытания проводятся по отдельно разрабатываемой и утверждаемой «Программе проведения высоковольтных испытаний кабеля 110 кВ».	Испытания на высоковольтные кабели проводятся в соответствии с отдельно разработанной и утвержденной программой испытаний высоковольтных кабелей 110 кВ.	-1,356
2	Оборудование должно быть рассчитано на двойные фидеры, а если такое оборудование отсутствует, в центральном шкафу предусматривают установку контроллера автоматического ввода резерва.	Оборудование должно быть способно управлять двумя фидерами, в случае отсутствия такого оборудования в центральном шкафу должен быть установлен переключатель ввода резерва.	-0,444

3	ТУ на поставку включают в себя, помимо прочего, следующее:	Спецификация покупки должна содержать и не ограничиваться:	-0,167
4	По результатам месяца подготовка отчета (10 число) по отклонениям от намеченного графика.	Отчет о ходе месяца (10 числа) об отклонениях от плана работы.	-0,271

Таблица 2. Результаты машинного перевода на английский язык

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

Результаты оценки качества после оптимизационного предредактирования представлены в таблице 3.

Таблица 3. Результаты применения оптимизационного предредактирования

	h_LEPOR (En1)	h_LEPOR (En2)	$\Delta hLEPOR$	СлЗП (Ru)	СлЗП (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

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; max – максимальное значение в выборке.

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

Использование полученных результатов исследования и внедрение программного комплекса в работу агентства переводов «ФИАС-Амур» позволило повысить эффективность использования систем МП и производительность труда редакторов переводов, а также оптимизировать затраты на оказание услуг перевода узкоспециальной технической документации.

В ходе внедрения 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 увеличилась с 3,8 до 4,3 стандартных страницы в час. Таким образом, внедрение программного комплекса позволило увеличить производительность редакторов переводов на 13,16%, что является значимым показателем при переводе большого объема документации в условиях дефицита квалифицированных кадров.

Описанные модели и методы 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 инструмент объективной оценки исполнителей в рамках поставленной задачи на перевод.

Программный комплекс может быть внедрен в компаниях, генерирующих от 1000 страниц перевода в месяц, предоставив инструментарий повышения качества перевода, в том числе для редакторов без знания языка перевода.

Заключение. В ходе работы 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; 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.

Программный комплекс представляет собой ансамбль алгоритмов и моделей, включая модель классификации и генеративную модель русского языка, каждая из которых имеет потенциал к доработке с целью повышения точности, что позволит улучшить эффективность программного комплекса в целом.

Результаты исследования подтверждают эффективность применения оптимизационного предредактирования русскоязычных узкоспециальных текстов с целью повышения качества МП на английский язык. Разработанный авторами программный комплекс оптимизационного предредактирования имеет потенциал к доработке и повышению точности. Интеграция программного комплекса в контур автоматизации процессов перевода технической документации позволяет снизить затраты на постредактирование МП и организацию переводческих процессов.

Благодарности. Работа выполнена в рамках соглашения о стратегическом научно-технологическом сотрудничестве между ФГБОУ ВО «КНАГУ» (г. Комсомольск-на-Амуре), Центром обработки и хранения научных данных ДВО РАН на базе ВЦ ДВО РАН – обособленного подразделения ХФИЦ ДВО РАН (г. Хабаровск) и ООО «Агентство переводов «ФИАС-Амур» (г. Комсомольск-на-Амуре).

Список источников

1. Quinci C., Pontrandolfo G. Testing neural machine translation against different levels of specialization. *Transkom*, 2023, vol.1, pp.174-209.
2. Canfora C., Ottmann A. Risks in neural machine translation. *Translation spaces*, 2020, vol. 9(1), pp. 58–77.
3. Kumar V., Kulkarni A., Singh P., Ramakrishnan G. A machine assisted human translation system for technical documents. *Miami MT Summit XV*, 2015, vol.2, pp. 259-272.
4. Hiraoka Y., Yamada M. Pre-editing plus neural machine translation for subtitling: effective pre-editing rules for subtitling of TED Talks. *MT Summit XVII*, Dublin, Ireland, 2019, vol.2, pp. 64-74.
5. 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)*, User studies papers, Prague, Czech Republic, 2017.
6. Taufik A. Pre-editing of Google neural machine translation. *Journal of English language and culture*, 2020, vol. 10, no. 2, pp. 64-74.

7. O'Brien, S. Controlling controlled English: an analytical of several controlled language rule sets. Proceedings of EAMT-CLAW, Dublin, Ireland, 2003, pp. 105-114.
8. Оборнева И. В. Автоматизация оценки качества восприятия текста / И. В. Оборнева // ВЕСТНИК Московского городского педагогического университета, 2015. – №2(5). – С. 221–233.
9. Дмитриева А.Д. Квантитативное исследование стратегий упрощения на материале адаптированных текстов для изучающих РКИ / А.Д. Дмитриева, А.Н. Лапошина, М.Ю. Лебедева // Компьютерная лингвистика и интеллектуальные технологии: по материалам международной конференции «Диалог», 2021. – С. 191-204.
10. Zhihotova 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, pp. 1-15.
11. Barber D. Bayesian Reasoning and machine learning. Cambridge: Cambridge university press, 2012, DOI:10.1017/CBO9780511804779.
12. MacKay D.J.C. Information theory, inference, and learning algorithms. Cambridge: Cambridge University Press, 2003, DOI: 10.2277/0521642981
13. Bishop, C.M. Pattern Recognition and Machine Learning. Springer, Berlin, 2006.
14. Животова А.А. Регрессионный анализ корреляции качества машинного перевода и параметров исходного текста / Животова А.А., Бердонос В.Д. // Информатика и системы управления, 2023. – №2(76). – С.121-133.
15. Li-Feng Han A., Wong D.F., Chao L.S. et al. Language-independent model for machine translation evaluation with reinforced factors. Proceedings of machine translation Summit XIV: Posters, Nice, France, 2013, pp. 215-222.
16. Дале Д. Многозадачная модель T5 для русского языка. – URL: <https://habr.com/ru/articles/581932/> (дата обращения: 20.08.2023).
17. Свидетельство о государственной регистрации программы для ЭВМ № 2023682260 Российская Федерация. Программный комплекс для предредактирования и машинного перевода узкоспециальных русскоязычных текстов на английский язык: № 2023680875: заявл. 09.10.2023: опублик. 24.10.2023 / А. А. Животова, В. Д. Бердонос; заявитель Животова А.А. – 1 с.

Животова Алена Анатольевна. Аспирант, ассистент кафедры Прикладная математика факультета компьютерных технологий Комсомольского-на-Амуре государственного университета, ORCID: 0000-0003-1037-5503, SPIN: 3142-9290, РИНЦ AuthorID: 771774, zhihotova.aa@gmail.com, г. Комсомольск-на-Амуре, пр. Ленина, д. 27.

Бердонос В.Д. Виктор Дмитриевич. К.т.н., доцент кафедры Прикладная математика факультета компьютерных технологий Комсомольского-на-Амуре государственного университета, ORCID: 0000-0003-4093-779X, SPIN: 8730-2226, РИНЦ AuthorID: 644399, berd1946@gmail.com, г. Комсомольск-на-Амуре, пр. Ленина, д. 27.

UDC 004.9:519.768

DOI:10.25729/ESI.2024.34.2.016

Optimization pre-editing of highly specialized Russian-language texts for its machine translation into English

Alena A. Zhihotova, Victor D. Berdonosov

Komsomolsk-na-Amure State University,

Russia, Komsomolsk-na-Amure, zhihotova.aa@gmail.com

Abstract. The authors study the process of text translation, particularly the method of optimizing pre-editing as a way to improve the quality of machine translation into English for Russian-language highly specialized texts. The paper considers the mathematical model of translation process and machine translation task formulation, proposes a new theory for probabilistic estimation of translation task complexity, provides the formulation and solution of optimizing pre-editing task, describes data preparation methodology for training automatic optimizing pre-editing model. As the research result the software package for optimizing pre-editing of Russian-language texts is developed. The software package has been developed using resources of the Center for Scientific Equipment Collective Use "Center for Processing and Storage of Scientific Data of the Far Eastern Branch of the Russian Academy of Sciences". Data for models training and validation are provided by Translation Agency FIAS-Amur Co., Ltd. Software package testing has proved the effectiveness of the proposed methods for improving the quality of machine translation of highly specialized Russian-language texts into English.

Keywords: machine translation, optimizational pre-editing, translation task complexity, translation quality

Acknowledgements: the research was conducted within the agreement on strategic scientific and technological cooperation between Komsomolsk-na-Amure State University (Komsomolsk-na-Amure), Center for Processing and Storage of Scientific Data of the Far Eastern Branch of the Russian Academy of Sciences (Khabarovsk), and Translation Agency FIAS-Amur Co., Ltd. (Komsomolsk-na-Amure).

References

1. Quinci C., Pontrandolfo G. Testing neural machine translation against different levels of specialization. *Transkom*, 2023, vol.1, pp.174-209.
2. Canfora C., Ottmann A. Risks in neural machine translation. *Translation spaces*, 2020, vol. 9(1), pp. 58–77.
3. Kumar V., Kulkarni A., Singh P., Ramakrishnan G. A machine assisted human translation system for technical documents. *Miami MT Summit XV*, 2015, vol.2, pp. 259-272.
4. Hiraoka Y., Yamada M. Pre-editing plus neural machine translation for subtitling: effective pre-editing rules for subtitling of TED Talks. *MT Summit XVII*, Dublin, Ireland, 2019, vol.2, pp. 64-74.
5. 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)*, User studies papers, Prague, Czech Republic, 2017.
6. Taufik A. Pre-editing of Google neural machine translation. *Journal of English language and culture*, 2020, vol. 10., no. 2, pp. 64-74.
7. O'Brien, S. Controlling controlled English: an analytical of several controlled language rule sets. *Proceedings of EAMT-CLAW*, Dublin, Ireland, 2003, pp. 105-114.
8. Osborneva, I. V. Avtomatizaciya ocenki kachestva vospriyatiya teksta [Automation of text perception quality assessment]. *VESTNIK Moskovskogo gorodskogo pedagogicheskogo universiteta [Herald of Moscow City Pedagogical University]*, 2015, no.2(5), pp. 221-233.
9. Dmitrieva A.D., Laposhina A., Lebedeva M. Kvantitativnoe issledovanie strategij uproshteniya na materiale adaptirovannykh tekstov dlya izuchayushhix RKI [A Quantitative Study of Simplification Strategies in Adapted Texts for L2 Learners of Russian]. *Komp'yuternaya lingvistika i intellektual'ny'e texnologii: po materialam mezhdunarodnoj konferencii "Dialog"* [Computer Linguistics and Intelligent Technologies: Proceedings of the International Conference "Dialog"], 2021, pp. 191-204.
10. Zhitovaya 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, pp. 1-15.
11. Barber D. *Bayesian Reasoning and machine learning*. Cambridge: Cambridge university press, 2012, DOI:10.1017/CBO9780511804779.
12. MacKay D.J.C. *Information theory, inference, and learning algorithms*. Cambridge: Cambridge University Press, 2003, DOI: 10.2277/0521642981
13. Bishop, C.M. *Pattern Recognition and Machine Learning*. Springer, Berlin, 2006.
14. Zhitovaya A.A., Berdonosov V.D. Regressionnyj analiz korrelyacii kachestva mashinnogo perevoda i parametrov iskhodnogo teksta [Regression analysis of the correlation between machine translation quality and source text features]. *Informatika i sistemy upravleniya [Information science and control systems]*, 2023, no.2(76), pp.121-133.
15. Li-Feng Han A., Wong D. F., Chao L.S. et al. Language-independent model for machine translation evaluation with reinforced factors. *Proceedings of machine translation Summit XIV: Posters*, Nice, France, 2013, pp. 215-222.
16. Dale, D. Mnogozadachnaya model` T5 dlya russkogo yazy`ka [T5 multitasking model for Russian language]. Available at: <https://habr.com/ru/articles/581932/> (accessed: 08/20/2023).
17. Zhitovaya A.A., Berdonosov V.D. Programmnyj kompleks dlya predredaktirovaniya i mashinnogo perevoda uzkospecial'nykh russkoyazy`ch-nykh tekstov na anglijskij yazy`k [Program complex for pre-editing and machine translation of highly specialized Russian-language texts into English]. *Software Registration Certificate of the Russian Federation no. 2023680875* (2023).

Alena Anatolievna Zhitovaya. Postgraduate student, assistant of Applied Mathematics department, faculty of computer technologies, Komsomolsk-na-Amure state university, ORCID: 0000-0003-1037-5503, SPIN: 3142-9290, AuthorID: 771774, zhitovaya.aa@gmail.com, Russia, Komsomolsk-na-Amure, 27 Lenin Ave.

Victor Dmitrievich Berdonosov. PhD in technical sciences, associate professor of Applied Mathematics department, faculty of computer technologies, Komsomolsk-na-Amure state university, ORCID: 0000-0003-4093-779X, SPIN: 8730-2226, AuthorID: 644399, berd1946@gmail.com, Russia, Komsomolsk-na-Amure, 27 Lenin Ave.

Статья поступила в редакцию 08.12.2023; одобрена после рецензирования 29.03.2024; принята к публикации 06.06.2024.

The article was submitted 12/08/2023; approved after reviewing 03/29/2024; accepted for publication 06/06/2024.