

УДК 81'33, 81.512.1

**Б.Э. Хәкимов,  
М.Р. Шәехов**

## **РУСЧА-ТАТАРЧА МАШИНА ТӘРЖЕМӘЧЕСЕ ӨЧЕН ПАРАЛЛЕЛЬ ТЕСТ КОРПУСЫНДА ЖӨМЛӘЛӘР ТӘНГӘЛЛЕГЕ**

При подготовке обучающих и тестовых корпусов для систем машинного перевода осуществляется поиск готовых параллельных текстов на двух языках в различных доступных источниках. Обычно считается, что пары предложений из оригинальных текстов и их переводов, если они были опубликованы и имеют официальный статус, являются эквивалентными. Однако в результате анализа таких предложений установлено, что степень их эквивалентности в определенных случаях может не соответствовать задачам разработки систем машинного перевода. В данной статье изучаются типичные случаи, в которых проявляются различия в требованиях к эквивалентности перевода предложений в исходных текстах и в качестве отдельных единиц параллельного тестового корпуса, используемого для оценки качества результатов машинного перевода. В результате исследования установлено, что предложения, извлекаемые из параллельных текстов на двух языках, и их переводы даже при полной контекстной эквивалентности могут потребовать значительного количества лексических, грамматических и других трансформаций.

**Ключевые слова:** машинный перевод, качество перевода, параллельный корпус, эквивалентность перевода, татарский язык, русский язык.

The current practice of preparing training and test corpora for machine translation systems usually involves searching for parallel texts in real use. Pairs of sentences from the original texts and their published translations are considered a priori equivalent. However, a detailed qualitative analysis of such sentences shows that this does not always meet the requirements of machine translation systems. This paper discusses typical differences between the equivalence of texts and the equivalence of parallel corpus units. The results show that sentences from the authentic texts and their translations, despite of their complete parallelism and contextual equivalence, may require a significant amount of lexical, grammatical and other transformations.

**Keywords:** machine translation, translation quality, parallel corpus, translation equivalence, Tatar language, Russian language.

Нәйрочелтәрле машина тәржемәсе системалары өчен өйрәтүче корпуслар һәм тест корпусларын эзерләүнең төп ысулы мөмкин кадәр күбрәк чыганақлардан ике телдә эер параллель жөмлөләр туплаудан гыйбарәт. Русча-татарча машина тәржемәсе сыйфатын бәяләү өчен, махсус параллель корпус төзелде [Хәкимов, Шәехов, с. 283–292]. Гадәттә мондый төр корпусларны эзерләгәндә оригиналь текстлардан һәм рәсми, басылып чыккан тәржемәләрдән алынган жөмлө парлары тулысынча эквивалент, тәңгәл дип санала. Эмма андый жөмлөләргә анализ ясау нәтижәсендә аларның тәңгәллек дәрәжәсе машина тәржемәсе системалары таләпләренә җавап бирмәскә дә мөмкин икәнлегә ачыклана.

Эквивалентлык – тәржемә гыйлеменәң төп мәсьәләләреннән берсе. Тәржемә теориясендә ул төрле лингвистик һәм экстралингвистик факторларга бәйлә булган чагыштырма характеристика буларак карала. В.Н. Комиссаров эквивалентлыкны тәржемә белән төп нөсхәнең максималь мөмкин булган лингвистик якынлыгы дип билгели һәм төп нөсхә белән тәржемә арасында аның берничә дәрәжәсен аерып күрсәтә: коммуникация максаты дәрәжәсе, ситуацияне тасвирлау дәрәжәсе, ситуацияне тасвирлау ысулы дәрәжәсе, әйтелгәннәренәң структур төзелеше дәрәжәсе, лексик берәмлекләр семантикасы дәрәжәсе [Комиссаров, б. 167]. Л.С. Бархударов фикеренчә, эквивалентлык – ул төп нөсхә һәм тәржемә телендә текстларның мәгънәви яктан туры килүендә чагыла торган семантик категория [Бархударов, б. 150]. Шулай итеп, тәржемәнең реципиентка йогынтысы төп нөсхәгә тиндәш булсын өчен, аларның тәңгәллегә лексик, грамматик һәм синтаксик дәрәжәдә генә түгел, ә коммуникатив ситуация һәм максат дәрәжәсендә дә тәмин ителергә тиеш. Тагын да катлаулырак дәрәжәләрдәге эквивалентлыкка ирешү өчен, тел чараларының гади тәңгәллегенә игътибар итмәскә мөмкин дип санала.

Өлеге мәсьәлә машина тәржемәсе максатларында параллель корпуслар әзерләү процессында үзенчәлекле рәвештә чагылыш таба. Машина тәржемәсенәң теләсә кайсы конкрет нәтижәсенә карата төрле дәрәжәләрдәге тәңгәллек күзлегеннән барлык таләпләр кулланырлык булса да, мондый прагматик юнәлешле параллель корпусларда урын алган жөмлә парлары аерым үзенчәлекләргә ия.

Төп үзенчәлек контекст күләме белән бәйлә. Тәржемә ителә торган эчтәлекнең контекстуаль мәгънәсә бөтен текст яки аның бер өзегә дәрәжәсендә аерылып тора. Кагыйдә буларак, машина тәржемәсе максатларына хезмәт итә торган параллель корпусларга аерым жөмләләр контекст бәйләнешләрен сакламыйча гына урнаштырыла, чөнки текстлардан барлык жөмләләр түгел, ә билгеле бер критерийларга туры килгәннәре генә сайлап алына. Димәк, чыганак текстлардагы параллель жөмләләр ике телдә дә коммуникатив ситуация һәм контекст дәрәжәсендә эквивалент булырга мөмкин (мәсәлән, әдәби эсәрләрдә төп эчтәлек – образлар аша, ә матбугаттагы яңалык текстларында фактлар аша бирелә). Параллель корпуска кертелгәч исә, контекстан максималь бәйсезлек шартларында шул ук жөмләләр структур төзелеш яисә лексик берәмлекләр дәрәжәләрендә тәңгәллек таләпләренә инде туры килмәскә дә мөмкин. Башкача әйткәндә, тәржемә ителгән жөмлә эквивалентлыгының кайбер факторлары өлеге жөмлә конкрет контекстан аерылгач юкка чыга. Шуңа бәйлә рәвештә мондый корпусларда параллель жөмләләренә бәяләү һәм аларга төзәтмәләр кертү өчен, билгеле бер кагыйдәләр билгеләү зарурлыгы туа. Мәсәлән, жөмләләр тәңгәллегенә ирешү тагын да төгәлрәк дәрәжәдә мөмкин булган очракта, тиешле төзәтмәләр кертеп, мөмкин кадәр күбрәк контекстларга туры килерлек вариант барлыкка китерү максатка ярашлы була.

Жөмлә парларының контексттан бәйсез эквивалентлыгы аларның төп нөсхәдәнме, тәрҗемәдәнме булуына караганда мөһимрәк, шуна күрә дә төзәтмәләр параллель жөмләләрнең икесенә дә кертелә ала. Без аерым очрақларда хәтта әдәби әсәрләрдән алынган оригиналь жөмләләрне дә үзгәртүгә мөмкин дип саныбыз, чөнки матур әдәбият тәрҗемәсендә гадәттә югарырак дәрәҗәдәге эквивалентлык турында сүз бара. Бер яктан, корпуста реаль тексттан булган оригиналь жөмләне саклап калу өчен, тәрҗемәгә үзгәрешләр кертү кулайрак. Икенче яктан, әгәр корпус лингвистик күренешләрнең төрлелеге (шул исәптән грамматик формалар һәм лексик берәмлекләр) ягыннан репрезентативлырак булса, машина тәрҗемәсе сыйфатын бәяләүнең нәтиҗәлелеге дә арта. Шулай итеп, төзәтү варианты һәр очракта аерым сайлана. Моннан тыш, һәр вариантны редакцияләүгә сарыф ителә торган вакыт һәм хезмәт күләме дә игътибарга алынырга мөмкин.

Гомумән алганда, машина тәрҗемәсе нәтиҗәләрен бәяләү өчен параллель жөмләләр корпусын эзерләгәндә, экспертка параллель жөмлә парларының максималь эквивалентлыгы һәм тәрҗемәнең табигыйлеге арасында баланс саклау бурычы йөкләнә.

Русча-татарча машина тәрҗемәсе сыйфатын анализлау өчен хезмәт итә торган параллель жөмләләр корпусын төзү барышында, текст составындагы һәм корпусның аерым берәмлекләре буларак жөмләләр эквивалентлыгына карата булган таләпләрдәге аермаларны истә тотып, түбәндәге төр төзәтмәләр кертелде. Әлеге мәкаләдә параллель русча-татарча тест корпусынан [Хусаинов, Хәкимов, Шәехов, 2022] татарча жөмләләрне төзәтү мисаллары китерелә: башта чыганак текстлардан үзгәрешсез алынган жөмләләр, жәя эчендә татарча жөмләнең төзәтелгән варианты бирелә.

1. Контексттан бәйсез вариантка үзгәртү: *Бушueva сидела в кабине второй машины. – Бушueva икенче машинаның кабинасында **бара иде** (Бушueva икенче машинаның кабинасында **утыра иде**).*

Әлеге мисалда транспорт чарасы пассажирына карата *сидеть* – *барырга* фигыльләренең контекст эквивалентлыгы күзәтелә. Татарча жөмләгә төзәтмә кертү барышында фигыль туры тәңгәлләккә алмаштырыла, югарырак дәрәҗәдәге эквивалентлыкка да зыян килми, жөмләнең табигыйлеге югалмый. Димәк, бу очракта төзәтмә кертү үзен аклый. Тагын бер мисалда тәрҗемә төгәлрәк вариантка үзгәртелә: *Все знают – и никто не понимает. – Барысы да белә – һәм беркем **аңлата алмый** (Барысы да белә – һәм беркем дә **аңламый**).*

2. Параллель жөмләдә туры тәңгәллеге булмаган сүзләрне бетерү: *Вспомнила Суок из сказки «Три толстяка». – «Три толстяка» әкиятендәге Суокны **да** искә төшерде. («Три толстяка» әкиятендәге Суокны **искә төшерде**).*

Мисаллардан күренгәнчә, параллель жөмләләр еш кына анык-лаучы яки көчәйтүче функция башкара торган тел чаралары йөрткән мәгънә төсмерләре белән аерылып торырга мөмкин: *По словам Тяминава Артура, картофель хорошо хранится. – Тяминов Артур*

әйтүенчә, бәрәңге **бик** әйбәт саклана (Тяминов Артур әйтүенчә, бәрәңге әйбәт саклана).

3. Грамматик форманы үзгәртү: *Минздрав ответил на предложение вернуть курилки в аэропорты.* – Сәламәтлек саклау министрлыгы тәмәке тарту урыннарын **аэропортка** кайтару тәкъдименә җавап бирде. (Сәламәтлек саклау министрлыгы тәмәке тарту урыннарын **аэропортларга** кайтару тәкъдименә җавап бирде).

Китерелгән мисалда, жөмлөләр эквивалентлыгын арттыру өчен, төп нөсхәдә булган күплек сан формасы тәржемәдә дә торгызыла.

4. Стилистик хаталарны төзәтү һәм тәржемәне стилистик яктан үзгәртү: *Убежище в январе в ЕС получили только 13 граждан Грузии, а 1288 гражданам дали отказ.* – Гыйнварда ЕСтә **сыенуны** Грузиянең 13 гражданы гына алган, ә 1288 **гражданга кире кагылган**. (Гыйнварда ЕСтә **сыену мөмкинлеген** Грузиянең 13 гражданы гына алган, ә 1288 **гражданның үтенече кире кагылган**).

Бу төрдәге мисаллар күрсәткәнчә, чыганак буларак рәсми нәшер ителгән текстлар һәм аларның тәржемәләре кулланылуга карамастан, жөмлөләрдә хаталар да булырга мөмкин (ешрак стилистик характерда). Тәңгәллек дәрәжәләре күзлегеннән караганда, мондый хаталарны төзәтү югарыдагы мисаллар белән чагыштырганда капма-каршы нәтижәгә китерә, ягъни төзәтмәләр керткәч, бигрәк тә билгеле бер контекст белән турыдан-туры бәйлә булмаган очракта, тәңгәллекнең түбәнрәк дәрәжәсенә генә түгел, ә югарырагына да ирешергә мөмкин.

Түбәндәге мисалда исә төзәтүнең төрле юнәлештәгә ике төрөн күрергә була: бер яктан, *бумага* сүзенең тәржемәсе контекстка бәйлә булмаган вариантка алмаштырыла һәм, шулай итеп, тәңгәллек түбәнрәк дәрәжәгә күчә. Икенче яктан, сүзгә-сүз тәржемәне үзгәртү исәбенә (*со словами* – *мондый сүзләр белән*), сүзтезмәнең гомуми эквивалентлык дәрәжәсе күтәрелә: *Они должны за подписью руководителя присылать в Роскомнадзор бумагу со словами:...* – *Алар җитәкче имзасы белән «Роскомнадзор»га мондый сүзләр белән хәбәр җибәрергә тиешләр:...* (Алар җитәкче имзасы белән Роскомнадзорга **мондый эчтәлекле кагазь** җибәрергә тиешләр:...)

5. Телләрнең берсендә синонимик лексик берәмлек яки грамматик форма булганда вариант өстәү: *Приглашать в школы героев войны.* – *Мәктәпкә сугыш геройларын чакырырга (1) Мәктәпкә сугыш геройларын чакырырга. 2) Мәктәпкә сугыш каһарманнарын чакырырга*).

Өлеге мисалдагы русча *герой* сүзенең татар телендә ике варианты бар. Димәк, әгәр тест корпусында синонимнарның берсе генә теркәлгән булса, автомат тикшерү вакытында тәржемәнең тулысынча эквивалент булган икенче, синонимик варианты төгәл түгел дип бәяләнәчәк. Шул сәбәпле параллель корпуска ике пар жөмлә өстәлә, аларда бер телдә бер генә вариант кулланыла, ә икенчесендә аерыла. Мондый хәл лексика дәрәжәсендә дә, грамматика дәрәжәсендә дә барлыкка килергә мөмкин. Түбәндәге мисалдан күренгәнчә, рус телендәге фигыльнең инфинитив формасының татар телендә ике эк-

виваленты бар: *В теории его **сосчитать** несложно. – Теориядә аны **санарга** авыр түгел (1) Теориядә аны **санарга** авыр түгел. 2) Теориядә аны **санану** авыр түгел).*

Телләрнең берсендә теге яки бу грамматик мәгънәнең формаль күрсәткече булмаган очракларда да параллель жөмлә вариантларын өстәргә кирәк була. Мәсәлән, рус телендә үткән заманда фигыльнең зат күрсәткечләре юк, татар телендә исә бу грамматик мәгънә ачык белдерелә: *Так и **скупила** все в этом маленьком магазине. – Шулай итеп, бу кечкенә кибертәгә бөтен айберне **сатып алып та бетердем** (1) Шулай итеп, бу кечкенә кибертәгә бөтен айберне **сатып алып та бетердем**. 2) Шулай итеп, бу кечкенә кибертәгә бөтен айберне **сатып алып та бетерде**.*

Тагын да төгәлрәк караганда, соңгы мисалда *сатып алып та бетердең* дигән тагын бер вариантны өстәргә кирәк. Ә татар телендә үткән заманның төп ике формасы булуын да исәпкә алсак, мөмкин булган вариантлар саны икеләтә арта. Бу очракта ике стратегия булырга мөмкин: мөмкин булган барлык вариантларны кертү яисә статистик яктан ешрак очрый торганны сайлау.

Тикшеренү нәтиҗәләре күрсәткәнчә, еш кына чыганак текстлар һәм аларның тәрҗемәләре арасында тулы контекстуаль эквивалентлык булганда да, аерым жөмлә парлары буларак махсус корпуска өстәгәндә, шактый күләмдә лексик, грамматик һәм башка трансформацияләр таләп ителә. Шундый лингвистик мәгълүмат җыелмаларының гади тел эквивалентлыгы күзлегеннән караганда төгәлрәк булуы машина тәрҗемәсенә сыйфатын арттырырга мөмкинлек бирә.

#### Әдәбият

*Бархударов Л.С.* Язык и перевод. М., 1975. 230 с.

*Комиссаров В.Н.* Общая теория перевода (лингвистические аспекты): Учеб. для ин-тов и фак. иностр. яз. М.: Высшая школа, 1990. 253 с.

*Федоров А.В.* Основы общей теории перевода. М.: Высшая школа, 1983. 300 с.

*Хакимов Б.Э., Шаехов М.Р.* К вопросу создания параллельного тестового корпуса для задачи машинного перевода в русско-татарской паре // Восьмая Международная конференция по компьютерной обработке тюркских языков «TurkLang-2020». (Труды конференции). Уфа: ИИЯЛ УФИЦ РАН, 2020. С. 283–292.

*Хусаинов А.Ф., Хакимов Б.Э., Шаехов М.Р.* База данных тестового параллельного корпуса на русском и татарском языках. Свидетельство о регистрации базы данных 2022620573, 17.03.2022.

*Koller W.* Equivalence in Translation Theory // Ed. A. Chesterman, Readings in Translation Theory. Helsinki: Oy Finn Lectura Ab., 1989. P. 99–104.

*Хакимов Булат Эрнст улы,*  
*филология фәннәре кандидаты, Казан федераль университетының*  
*билингваль һәм цифрлы белем бирү кафедрасы доценты,*  
*ТР ФА Гамәли семиотика институтының әйдәп баручы фәнни хезмәткәре*

*Шәехов Марат Рәшит улы,*  
*ТР ФА Гамәли семиотика институтының әйдәп баручы фәнни хезмәткәре*