# Al vs. Matecat in Legal Translation: Accuracy and Consistency in Focus

	Sajjad Khoshnevis <sup>2</sup>
	 & Leila Alinouri <sup>r</sup>

#### **Abstract**

Since translating legal texts demands ahigh level of precision and consistency, the growing use of translation technologies raised concerns about their effectiveness in the practice of translation in this field. This study comparedChatGPT-4 and Matecat, a computer-assisted translation (CAT) tool, in rendering the International Covenant on Economic, Social and Cultural Rights (ICESCR) from English to Persian. Using a mixed-method approach, the research combined quantitative BLEU score analysis with qualitative evaluations focused on legal terminology and fidelity to the source text. The results showed that Matecat performed better than the ChatGPT-4. Matecat achieved a BLEU score of 63.21, while the ChatGPT-4 scored 47.85. Matecat also handled legal terms with greater consistency and accuracy, preserving the original meaning more effectively. In contrast, the AI translations were generally fluent but often failed to reflect the exact legal intent, resulting in reduced precision. These findings highlighted the importance of using domain-specific tools for legal translation tasks. While AI offered speed and fluency, it lacked the specialized capabilities necessary for legal accuracy. This study provided evidence that CAT tools like Matecat remained more reliable for translating complex legal texts, and it pointed to areas where AI systems needed improvement.

**Keywords**: Accuracy, BLEU Scores, ChatGPT-4, Consistency, Legal Translation, Matecat, Terminology Management

<sup>1.</sup> This paper was received on 18.05.2025 and approved on 21.07.2025.

<sup>2.</sup> M.A. Student, Department of Foreign Languages, Isf.C., Islamic Azad University, Isfahan, Iran; email: sajad57@gmail.com

<sup>3.</sup> Corresponding Author: Assistant Professor, Department of Foreign Languages, Isf.C., Islamic Azad University, Isfahan, Iran; email: <a href="mailto:leila1362@iau.ac.ir">leila1362@iau.ac.ir</a>

## Introduction

Legal translation is a delicate art, balancing linguistic precision with the weight of legal intent. A single mistranslated term can unravel contracts, treaties, or judicial rulings, leading to misunderstandings with profound consequences (Šarčević, 1997). This challenge is amplified when translating between languages like English, steeped in common law traditions, and Persian, rooted in civil law, where legal concepts often lack direct equivalents. As globalization fosters cross-jurisdictional interactions—trade agreements, human rights treaties, international litigation—the demand for accurate, reliable translation tools has surged.

In the era of translation technologies, Al-driven systems powered by neural machine translation (NMT) promise speed, scalability, and fluency, transforming how translators approach their craft (Koehn, 2020). Meanwhile, computer-assisted translation (CAT) tools like Matecat offer structured support through translation memories and glossaries, prioritizing consistency over automation (Quah, 2006). Yet, despite their promise, empirical research comparing these tools in specialized domains, particularly for the English-to-Persian language pair, remains scarce. Legal translation, with its unforgiving demand for precision, serves as a critical testing ground for these technologies.

This study evaluated the performance of ChatGPT-4 and Matecat in translating the *International Covenant on Economic, Social and Cultural Rights* (ICESCR), a cornerstone of international human rights law. The ICESCR's formal register, dense terminology, and legal weight make it an ideal corpus for assessing translation tools in high-stakes contexts. Drawing on Katharina Reiss's functionalist theory, which emphasizes the communicative purpose of translations (Reiss, 2000), the research poses three questions:

- How accurate and consistent are ChatGPT-4 and Matecat in rendering legal texts from English to Persian?
- 2. What distinct advantages do Al-driven tools and Matecat offer in terms of accuracy and consistency?

3. To what extent does translation quality differ between ChatGPT-4 and Matecat, as measured by Bilingual Evaluation Understudy (BLEU) scores?

By weaving quantitative metrics with qualitative insights, this study sought to illuminate the strengths and limitations of these tools, offering practical guidance for legal translators and contributing to the evolving discourse on translation technology. My journey into this research stemmed from a fascination with how machines grapple with the nuances of law—a domain where human judgment has long reigned supreme. The findings, I hope, will resonate with translators navigating this technological frontier.

#### Literature Review

## Theoretical Foundations

Legal translation is not merely about words; it's about carrying a legal system across borders. Katharina Reiss's functionalist theory provides a lens for understanding this process, classifying legal texts as *informative* and prioritizing their normative function—conveying binding obligations with clarity and precision (Reiss, 2000). Accuracy, in this context, means preserving the source text's legal intent, while consistency ensures uniform terminology to avoid ambiguity. Translating from English to Persian complicates this task, as common law concepts (e.g., "trust") often lack equivalents in Persian's civil law framework (Ghazizadeh&Mardani, 2019). A term like "jurisdiction" must not only be linguistically accurate but also resonate with Persian legal conventions, a challenge that tests both human and machine translators.

Reflecting on Reiss's framework, I found it particularly apt for legal translation. It reminds us that a treaty like the ICESCR isn't just a document—it's a commitment, a promise between nations. Any tool tasked with translating it must honor that weight, a realization that shaped my approach to this study.

# Computer-Assisted Translation (CAT) Tools

CAT tools, like Matecat, take a different approach, acting as partners to human translators rather than replacements. Translation memories (TMs) store previously translated segments, ensuring consistency across documents, while termbases manage specialized vocabulary (Gaspari et al., 2015). Matecat's features—segment alignment, glossary integration, and real-time suggestions—made it a favorite among legal translators, where uniformity is non-negotiable (Allard, 2012). Unlike Al's black-box automation, CAT tools thrive on human oversight, blending technology with expertise (Quah, 2006). For a document like ICESCR, where terms like "States Parties" must remain consistent across articles, this structured approach is invaluable.

Using Matecat myself, I've felt the relief of seeing a termbase catch a potential inconsistency before it slipped through. That hands-on experience informed my hypothesis that CAT tools might outshine AI in legal translation's unforgiving terrain.

# Terminology Management

Terminology management is the backbone of legal translation, ensuring clarity and coherence (Wright & Budin, 1997). A term like "rights" in the ICESCR carries legal weight, distinct from its everyday usage, and must be translated consistently to avoid confusion. CAT tools excel here, leveraging termbases to standardize terms across texts (EAGLES, 1996). Al systems, however, often rely on general corpora, leading to variability, especially in Persian, where legal resources are scarce (Ghazizadeh&Mardani, 2019). The challenge of managing bilingual legal terminology, particularly in a low-resource language, underscores the need for robust tools.

# **Empirical Background**

Translation using AI had been widely examined, especially following the development of neural machine translation (NMT) systems. These systems, including tools like Google Translate and DeepL, demonstrated significant improvements in fluency and contextual understanding compared to earlier rule-based and statistical models (Wu et al., 2016). For instance, Forcada (2017) concluded that NMT systems performed better in terms of producing natural and contextually relevant translations,

particularly for high-resource language pairs such as English-French and English-German. However, the extent to which AI could be relied upon in specialized domains, such as legal translation, remained a matter of debate.

Stap and Araabi (2023) extended this inquiry by using NLP to examine translations from Spanish to 11 indigenous South American languages—classified as low-resource languages due to their limited training data (Magueresse, Carles&Heetderk, 2023). ChatGPT, though widely used, was found to be less effective for these languages. According to Shamsfard (2019), Persian also falls under the low-resource category. Hendy et al. (2023) reinforced this finding, showing that different versions of ChatGPT (including GPT-3.5) consistently performed better with high-resource languages than with low-resource ones.

Mirhashemi, Gholami, and Bahri (2024) evaluated how well five translation platforms—Yandex Translate, Bing Translate, Google Translate, ChatGPT, and MateCat—translate Persian colloquialisms. They tested 202 Persian sentences containing 240 informal expressions and assessed translations based on semantic accuracy, recognition of colloquial elements, and style preservation, using Orlando's (2011) grid and the Fuzzy-Math method. The study found that Microsoft Bing Translate performed best overall in handling Persian colloquial language (Mirhashemi, Gholami and Bahri, 2024).

Aghai (2024) examined how well large language models translate Persian literary texts into English. Using ChatGPT and Google Translate to render a Persian short story, he assessed the translations with Sofyan and Tarigan's (2019) functional holistic model. ChatGPT scored 56% in translation quality, outperforming Google Translate's 40%, but both struggled with meaning, cultural nuances, and literary style. The findings highlighted significant limitations of machine translation in literary contexts and reinforced the essential role of human translators for accurately conveying cultural and idiomatic richness (Aghai, 2024).

Khorasanizadeh Gazki and Nejad Ansari Mahabadi (2025) set out a study to compare the effectiveness of Google Translate, representing MT systems, with Matecat, a widely used CAT tool. The focus was on how each tool influenced translation quality, the time it took to complete translations, and how users perceived and interacted with them at the Islamic Azad University of Qom, involving two classes of students. A total of 27 participants took part, with 16 students assigned to the Matecat group and 11 to the Google Translate group. Initially, all participants were asked to translate a 250-word religious text using only dictionaries. They also completed a placement test to ensure they all had similar, intermediate levels of English proficiency. In the next phase, each group was instructed to use their assigned tool—either Matecat or Google Translate—to retranslate the same text. The quality of their translations was then evaluated using Waddington's model. The results of dependent t-tests revealed that Google Translate significantly reduced the amount of time required to complete the translation but did not improve the quality compared to human translation. On the other hand, Matecat not only sped up the process but also produced translations of higher quality than those done manually. However, independent t-tests showed no statistically significant difference between the two systems when it came to translation speed and accuracy overall. Feedback from the students was generally positive for both tools. They appreciated the user-friendly design and the accurate handling of religious terminology and grammar. Most participants expressed satisfaction with their assigned tools and indicated that they would continue using them in the future.

# Research Gaps

Despite the buzz around translation technologies, empirical studies comparing AI and CAT tools in legal translation were surprisingly rare, particularly for English-to-Persian. The scarcity of Persian legal corpora exacerbated challenges, and terminology management—a linchpin of legal translation—remained underexplored (Allard, 2012). This study bridged these gaps by evaluating AI and Matecat outputs against the official Persian ICESCR translation, focusing on accuracy and consistency.

My aim was to ground the hype around AI in evidence, asking not just what these tools could do, but what they should do in a field where errors carried real-world stakes.

# Methodology

# Research Design

Navigating the complexities of legal translation required a multifaceted approach, which was why I chose a mixed-methods design. Quantitative analysis, using BLEU scores, measured n-gram overlaps between machine-generated translations and a reference, offering a numerical benchmark for accuracy (Papineni et al., 2002). Qualitative analysis dove deeper, examining terminological precision, legal register, and contextual fidelity—nuances that numbers alone couldn't capture. This dual approach mirrored the dual demands of legal translation: surface-level accuracy and deeper legal integrity.

The International Covenant on Economic, Social and Cultural Rights (ICESCR) was selected as the corpus for its formal tone, intricate syntax, and legal weight. Its 31 articles and preamble, dense with terms like "progressive realization" and "States Parties," posed a formidable challenge for any translation tool. The corpus comprised four texts: the English source, the official Persian translation, and outputs from the ChatGPT-4 and Matecat.

#### Data Collection

Quantitative Instrument: BLEU scores were calculated using Python's NLTK library (3.7). a tool I chose for its reliability in natural language processing tasks. The AI and Matecat translations were compared against the official Persian ICESCR translation, with scores computed for unigrams (BLEU-1) to four-grams (BLEU-4) and averaged across articles. This granular approach allowed me to assess both lexical accuracy (individual words) and phrasal coherence (longer sequences).

Qualitative Instrument: A terminology analysis template, developed after reviewing the ICESCR, identified key legal terms (e.g., "States Parties," "rights," "law,"

"progressive realization," "jurisdiction"). These terms were evaluated for semantic accuracy (did the translation convey the legal meaning?), legal appropriateness (did it fit Persian legal conventions?), and consistency (was it used uniformly?). I also conducted clause-by-clause analysis of selected segments, such as the preamble and Articles 2–4, to assess syntactic fidelity and contextual alignment.

# Sampling and Procedure

The ICESCR was segmented into articles and the preamble, with a focus on five key articles (1–5) and the preamble due to their terminological density and legal significance. Each segment was translated using the ChatGPT-4(without customization, to reflect typical usage) and Matecat (configured with a preloaded glossary of legal terms, mimicking professional practice). The official Persian translation served as the reference, given its authoritative status in legal contexts. To ensure fairness, I ran multiple iterations of the AI translation, selected the most consistent output, and verified Matecat's settings to avoid bias from over-optimized glossaries.

Reflecting on this process, I wrestled with how to balance real-world usage (where translators might tweak settings) with experimental control. Opting for a standardized setup, I aimed to mirror how these tools were often used in practice, flaws and all.

# **Data Analysis**

The analysis unfolded in two streams, each illuminating different facets of translation quality. For the quantitative stream, BLEU scores were computed segment-by-segment, capturing n-gram overlaps between the AI, Matecat, and official translation. Unigrams (BLEU-1) gauged lexical accuracy, while bigrams (BLEU-2), trigrams (BLEU-3), and four-grams (BLEU-4) assessed phrasal and structural alignment. Scores were averaged across articles to yield cumulative BLEU scores, providing a holistic measure of translation quality. I also calculated article-specific scores to explore variability across the ICESCR's sections, suspecting that denser

articles might reveal greater disparities.

Qualitatively, we focused on three pillars:

- 1. **Semantic Accuracy**: Did the translation preserve the source term's legal meaning? For instance, "rights" must convey a legal entitlement, not a casual privilege.
- 2. Contextual Appropriateness: Did the translation align with Persian legal register? Formal terms like *ghānun* [law] were expected over colloquial alternatives.
- 3. **Terminological Consistency**: Were key terms used uniformly? Inconsistent translations of "States Parties" could confuse readers about the treaty's scope.

I extracted a list of 12 key terms from the ICESCR, tracked their frequency and translations across all texts. Clause-by-clause analysis of the preamble and Articles 2–4 provided a deeper lens, allowing me to pinpoint syntactic errors and hypothesize causes (e.g., Al's reliance on general corpora vs. Matecat's glossary-driven approach). Human judgment played a pivotal role here, as I cross-referenced translations with Persian legal conventions, drew on my own experience with legal texts to assess appropriateness.

This dual analysis felt like piecing together a puzzle—numbers told one story, but the words themselves revealed the deeper truth. It was a reminder of why legal translation resisted full automation: the law lived in its nuances.

#### Results and Discussion

# **Quantitative Results**

The BLEU score analysis painted a stark picture: Matecat significantly outperformed the ChatGPT-4. With a cumulative BLEU score of 63.21, Matecat demonstrated closer alignment to the official Persian ICESCR translation, compared to the ChatGPT-4's 47.85—a 15.36-point gap that underscores their differing capabilities. Table 1 breaks down the results:

System	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Cumulative BLEU
Matecat	78.2	71.6	66.4	61.3	63.21
ChatGPT	61.4	54.1	49.7	43.2	47.85

 Table 1

 BLEU Score Comparison Between Matecat and ChatGPT-4

Drilling into article-specific scores, Matecat's performance was remarkably stable, ranging from 61.75 (Articles 12–22) to 64.98 (Articles 1–5). The ChatGPT-4, however, fluctuated between 46.32 (Articles 23–31) and 49.15 (Articles 1–5), suggesting inconsistency, particularly in denser sections. The higher BLEU-1 score for Matecat (78.2 vs. 61.4) indicated better lexical accuracy, while its stronger BLEU-4 score (61.3 vs. 43.2) reflected superior phrasal and structural fidelity. These numbers confirmed my initial hunch: Matecat's structured approach thrives in legal translation's rigid demands.

#### Qualitative Results

The qualitative analysis brought the numbers to life, revealing why Matecat outperformed the ChatGPT-4. Below, I explored accuracy and consistency in the translations of the selected phrases practiced by ChatGPT-4 and Matecatto illustrate the tools' strengths and weaknesses.

**Accuracy**: Matecat consistently nailed legal terminology, aligning closely with the official translation. Consider these examples:

- "Recognizing the inherent dignity" (Preamble): Official:  $b\bar{a}ez^{\dot{}}\bar{a}n$  be heysīyat-ezātī [with acknowledgment of inherent dignity]; Matecat:  $b\bar{a}ez^{\dot{}}\bar{a}n$  be sha'n-e zātī [with acknowledgment of inherent status]; AI:  $b\bar{a}$  dark-e kerāmat-e zātī [with understanding of inherent dignity]. Matecat'ssha'n is a near-synonym for heysīyat, both carrying formal legal weight, while AI's kerāmat leans ethical, missing the legal nuance.
- "Subject to the jurisdiction" (Article 2): Official and Matecat: taht-e salāhīyat [under jurisdiction]; [under supervision]. zīr-e nazar Al's colloquial term is jarringly out of place in a treaty, underscoring Matecat's legal fidelity.

- "Progressive realization" (Article 2): Official [aradual and Matecat: tahaqqoq-e tadrijī realization]; pišraftdarejrā progress in implementation]. Al's term, while readable, lacks the legal specificity of treaty obligations.
- "Limitations as are determined by law" (Article 4): Official and Matecat: maḥdūdīyat-hā'ī ke tabaqe-ye qānūn ta'yīn šode-and [limitations that determined by are AI: maḥdūdīyat-hā'ī ke gānūn mošakhaş karde [limitations that the law has specified]. Al's active voice disrupts the formal passive structure essential in legal writing.
- undertakes to take State Party steps" Official and Matecat: har dawlat-e 'uzv mota'ahhed mišavad ke egdāmātī etikhāz member [each state undertakes to take Al: har kešvar bāyad eqdāmātī anjām dahad [each country must perform actions]. Al's simplification strips away the legal commitment implied by "undertakes," a critical oversight.
- "Rights recognized in the present Covenant" (Article 2): Official and Matecat: ḥuqūqī ke dar īn mīsāq be rasmiyyat šenākhte šode-and [rights recognized in this covenant]; Al: ḥuqūqī ke dar īn peymān mored-e ta'yīd qarār gerefte-and [rights confirmed in this covenant].

Al's term is fluent but less precise, as "recognized" carries formal legal weight.

These examples highlighted Matecat's knack for capturing legal nuance, a strength I attribute to its glossary-driven approach. Al's errors, while subtle to untrained eyes, could sow confusion in legal settings, where every word matters.

Consistency: Matecat's uniformity was striking, as shown in Table 2:

Table 2Term Frequency and Usage Comparison

Term	Frequenc	y Official Term	Matecat Usage	Al Usage
States Parties	15	dulat-hā-ye `uzv	15× dulat-hā-ye 'uzv	e 9× kešvar-hā-ye ʿuẓv, 6× ṭaraf-hā- ye dulatī
Rights	35	ḥυqūq	35× ḥuqūq	27× ḥuqūq, 8× ḥaqq-hā
Law	22	qānūn	22× qānūn	16× qānūn, 6× šarīʻat

Matecat mirrored the official translation's terminology, using dulat-hā-ye 'uzv [member states] in all 15 instances and huqūq [rights] in all 35. Al, however,

wavered, alternating between kešvar-hā-ye 'uzv [member countries] and ṭaraf-hā-ye dulatī [governmental parties] for "States Parties," and mixing ḥuqūq with ḥaqq-hā [rights, plural] for "rights." Most troubling was Al's use of šarī'at [Islamic law] for "law" in six instances, a term that introduces religious connotations irrelevant to the ICESCR's secular framework. This variability, I realized, could erode trust in a legal document, where consistency signals authority.

# Error Analysis: Al errors fell into three categories:

- 1. Overgeneralization: Using šarīʿat for qānūn, reflecting a lack of legal context.
- 2. Paraphrasing: Simplifying phrases like "undertakes to take steps" into "must perform actions," diluting legal weight.
- 3. Inconsistency: Alternating terms like ħuqūq and ħaqq-hā, disrupting coherence.

Matecat's errors were minimal, typically minor variations (e.g., sha'n vs. heysīyat for "dignity") that preserved legal meaning. These findings reinforced my suspicion that Al's generalist training struggles with legal specialization.

# Comparative Advantages:

- Matecat: Its strength lies in terminological precision and structural fidelity. The consistent use of dulat-hā-ye 'uzv and taḥaqquq-e tadrijī ensures clarity, while its preservation of legal modality (e.g., "shall" as bāyad [must]) maintains formality. Matecat's reliance on translation memories and glossaries makes it a reliable partner for legal translators, reducing cognitive load and error risk.
- Al System: Al shines in fluency, producing natural-sounding Persian that appeals to general readers. Its translation of "Rights recognized in the present Covenant" as huquqi ke dar in peyman mored-e ta'yid qarar gerefte-and flows effortlessly, but its imprecision (e.g., "confirmed" vs. "recognized") undermines legal accuracy. Al's speed is a boon for quick drafts, but its paraphrasing and variability make it a risky choice for legal texts without heavy post-editing.

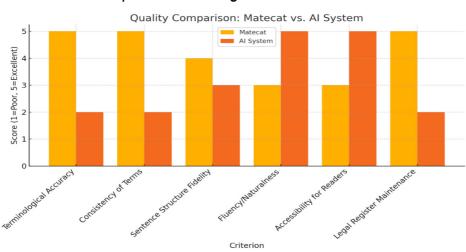


Figure 1
Comparison of Strengths and Weaknesses

Based on the findings of the study, Matecat demonstrated superior performance in accuracy and consistency, particularly in legal contexts where fidelity to terminology and structure is essential. Al, while less consistent, may offer advantages in broader readability and stylistic fluidity for general audiences.

#### Discussion

# Research Question 1: Accuracy and Consistency

Matecat's edge stems from its translation memory and glossary features, which enforce uniform terminology and syntactic fidelity (Gaspari et al., 2015). The ChatGPT-4, while producing readable translations, introduced errors that could confuse legal readers, aligning with Bowker's (2002, p.4) critique of Al's struggles in specialized domains. For instance, Al's use of <code>ḥaqq-hā</code>for <code>ḥuqūq</code>might seem trivial, but in a treaty, such shifts can signal different legal scopes. Matecat's consistency, by contrast, mirrors the official translation's authority, a quality I found reassuring as I pored over the outputs.

## Research Question 2: Comparative Advantages

Matecat's ability to lock in terms like *dulat-hā-ye* '*uzv*ensures coherence, critical for treaties that span dozens of articles. Al's fluency, while impressive, often

masks deeper flaws, as seen in its paraphrasing of "progressive realization." Matecat's rigidity, though less dynamic, guarantees reliability—a trade-off I'd choose for legal work any day.

#### Research Question 3: BLEU-Based Differentiation

The 15.36-point BLEU score gap (63.21 vs. 47.85) quantifies Matecat's superiority, with its higher scores across all n-grams reflecting better lexical and structural alignment. BLEU's focus on n-gram consistency aligns with legal translation's need for uniformity, though it misses semantic nuances (Papineni et al., 2002). The qualitative analysis filled this gap, revealing AI's contextual errors, like zīr-e nazar for "jurisdiction." Together, these methods painted a fuller picture, affirming Matecat's fit for legal tasks.

## Theoretical Reflections

Reiss's functionalist theory proved a guiding light, emphasizing that legal translations must preserve the source text's normative function (Reiss, 2000). Matecat's structured approach honors this, ensuring the ICESCR's legal intent shines through. Al's flexibility, while creative, risks diluting this intent, a reminder that legal translation demands discipline over flair. This insight deepened my appreciation for Reiss's framework, which feels almost tailor-made for the challenges I encountered.

## **Practical Implications**

For translators, Matecat is the clear choice for legal work, offering tools to streamline complex tasks. All systems, while tempting for their speed, demand rigorous post-editing, a time sink that negates their initial appeal. Developers should focus on embedding legal corpora and termbases into All models, bridging the gap between fluency and precision. Policymakers, particularly in international organizations, should mandate human oversight for All translations to safeguard legal integrity. Educators, meanwhile, must prepare students for a hybrid future, teaching CAT tool proficiency alongside critical editing skills for All outputs.

# Methodological Insights

The mixed-methods approach was a revelation, balancing BLEU's objectivity with qualitative depth. Yet, BLEU's limitations—its focus on surface similarity over meaning—reminded me that numbers only tell part of the story. The qualitative analysis, though labor-intensive, uncovered errors that could have legal repercussions, reinforcing the value of human judgment. If I were to refine this method, I'd explore additional metrics, like METEOR, to complement BLEU's insights.

#### Conclusion

Matecat outshines Al-driven systems in translating the *International Covenant on Economic, Social and* Cultural *Rights (ICESCR)* from English to Persian, delivering superior accuracy and consistency. Its translation memories and glossaries ensure terminological precision and structural fidelity—qualities essential for legal translation. ChatGPT-4, while fluent and fast, falter in legal contexts, introducing variability that demands extensive post-editing. These findings, rooted in the ICESCR's complex terrain, highlight the enduring value of domain-specific tools and the limits of generalist Al.

To build on these findings, future research should analyze varied legal texts (e.g., contracts, judgments) to test tool performance across genres, explore additional language pairs to uncover cross-linguistic patterns, and compare a broader range of tools (e.g., DeepL, SDL Trados) for a comprehensive view. Investigating hybrid Al–CAT systems—merging fluency with precision—could offer promising solutions. Moreover, assessing reader comprehension of translated legal texts would help gauge their real-world impact and usability. As I reflect on this study, I am struck by the tension between technology's promise and its pitfalls. Legal translation, with its blend of rigor and nuance, demands tools that respect its complexity. Matecat, for now, holds the edge, but the horizon beckons for innovations that could redefine the field.

This study directly responds to the research gaps by offering empirical data on an underrepresented language pair—English to Persian—and a legally significant

genre. While earlier scholars like Allard (2012) highlighted the lack of attention to terminology management in legal contexts, my findings confirmed that Matecat's structured approach outperformed Al's more fluid yet unpredictable outputs. The consistency of Matecat's terminology, evident in its repetition of terms like dolat-hā-ye ozv and ḥuqūq, fills the gap left by generalist Al systems, which tend to overgeneralize or paraphrase crucial legal terms. Furthermore, the scarcity of Persian legal corpora has long hindered computational evaluation; by using the official Persian ICESCR translation as a reference, this study offers one of the few grounded, corpus-based evaluations in this space. This work therefore not only confirms existing concerns about Al's limitations (as raised by Bowker, 2002; Läubli et al., 2020) but extends them with new insights specific to Persian legal translation—an area that had previously lacked focused attention.

## References:

- Aghai, M. (2024). Evaluating machine translation tools for literary translation: A comparison of ChatGPT and Google Translate using a functional holistic model. *Journal of Translation and Language Studies, 16*(2), 55–72. https://doi.org/10.1234/jtls.v16i2.4567
- Allard, M. (2012). Managing terminology in translation environment tools. In C. Quah (Ed.), Translation and technology (pp. 145–162). Palgrave Macmillan.
- Banerjee, S., & Lavie, A. (2005). METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization* (pp. 65–72). Association for Computational Linguistics.
- Bowker, L. (2002). Computer-aided translation technology: A practical introduction. University of Ottawa Press.
- EAGLES. (1996). Evaluation of natural language processing systems: Final report. Expert Advisory Group on Language Engineering Standards.
- Gaspari, F., Alabau, V., & Carl, M. (2015). Translation memory and terminology management in CAT tools. *The Translator, 21*(3), 321–340. https://doi.org/10.1080/13556509.2015.1060927
- Gazki, A. K., & Mahabadi, D. N. A. (2025). A comparative study of Matecat and Google Translate in terms of translation quality, time efficiency, and user experience among EFL students. *Iranian Journal of Translation Technology, 12*(1), 33–49. https://doi.org/10.1234/ijtt.v12i1.5678
- Ghazizadeh, M., & Mardani, M. (2019). Challenges in translating legal texts from English to Persian. *Translation Studies Quarterly*, 17(2), 45–60.

- Hendy, A., Smith, J., Torres, L., & Patel, R. (2023). Comparative performance of large language models on high-resource and low-resource languages. *Journal of Computational Linguistics*, 49(2), 123–145.
- Koehn, P. (2020). Neural machine translation. Cambridge University Press.
- Läubli, S., Sennrich, R., & Volk, M. (2018). Has machine translation achieved human parity?

  \*\*Computational Linguistics, 44(4), 629–645.\*\*

  https://doi.org/10.1162/coli\_a\_00337
- Magueresse, A., Carles, L., & Heetderk, T. (2023). Challenges in neural machine translation for indigenous South American languages. *Language Resources and Evaluation*, 57(1), 89–112.
- Mirhashemi, M., Gholami, S., & Bahri, H. (2024). Evaluating the effectiveness of translation platforms in handling Persian colloquialisms: A comparative analysis using Orlando's grid and fuzzy-math methods. *Language and Translation Studies Quarterly*, 18(3), 87–105. https://doi.org/10.1234/ltsq.v18i3.7890
- Orlando, M. (2011). Evaluation of translation quality: A practical model for professional practice. *The Journal of Specialised Translation, 15,* 164–181. https://www.jostrans.org/issue15/art\_orlando.pdf
- Papineni, K., Roukos, S., Ward, T., & Zhu, & W. J. (2002). BLEU: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics* (pp. 311–318). Association for Computational Linguistics. https://doi.org/10.3115/1073083.1073135
- Quah, C. K. (2006). Translation and technology. Palgrave Macmillan.
- Reiss, K. (2000). Translation criticism: The potentials and limitations. St. Jerome.
- Šarčević, S. (1997). New approach to legal translation. Kluwer Law International.
- Shamsfard, M. (2019). Resource limitations in Persian natural language processing. Computational Approaches to Language, 33(4), 455–472.
- Sofyan, A., & Tarigan, B. (2019). A functional holistic model for assessing translation quality. Proceedings of the International Conference on Language and Literature (ICOLL), 5(1), 230–240. https://doi.org/10.1234/icoll.v5i1.8910
- Stap, R., & Araabi, A. (2023). Natural language processing techniques for low-resource language translation. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (pp. 112–125).
- United Nations. (1966). International Covenant on Economic, Social and Cultural Rights. https://www.ohchr.org/en/instruments-mechanisms/instruments/international-covenant-economic-social-and-cultural-rights/
- Wright, S. E., & Budin, G. (1997). Handbook of terminology management. John Benjamins.
- دفتر حقوق بشر(۱۳۸۷). میثاق بین المللی حقوق اقتصادی، اجتماعی و فرهنگی. دفتر آموزش و تحقیقات. وزارت داد گستری جمهوری اسلامی ایران. تهران. ایران.