

DOI: 10.5281/zenodo.11425155

A CRITICAL EVALUATION OF GOOGLE TRANSLATE IN ENGLISH-ARABIC TRANSLATION: ACCURACY AND LIMITATIONS

Shahab Ahmad Al Maaytah^{1*}

¹Associate Professor Department of Languages and Humanities, Applied College, King Faisal University, Al-Asha, The Eastern Province, Saudi Arabia, Email: Salmaaytah@kfu.edu.sa, <https://orcid.org/0000-0002-0962-8801>

Received: 11/11/2025
Accepted: 18/12/2025

Corresponding Author: Shahab Ahmad Al Maaytah
(Salmaaytah@kfu.edu.sa)

ABSTRACT

This current research critically examines the performance of Google Translate in translating texts from English to Arabic, with a focus on accuracy, linguistic coherence, and cultural appropriateness. As machine translation (MT) tools become increasingly prevalent in academic, professional, and everyday contexts, understanding their capabilities and limitations is essential particularly for linguistically and culturally complex language pairs like English and Arabic. The research analyzes a selected corpus comprising general, technical, and literary texts translated by Google Translate, comparing them with human-translated equivalents. Key evaluation criteria include grammatical accuracy, semantic fidelity, idiomatic expression, and contextual relevance. The findings reveal that while Google Translate performs adequately in rendering fundamental and technical content, it struggles significantly with idioms, metaphors, polysemous words, and culturally bound expressions. The study also highlights structural issues related to Arabic syntax and morphological agreement, which frequently result in unnatural or ambiguous translations. Ultimately, this research highlights the significance of human post-editing and the necessity of context-aware translation in producing high-quality English-Arabic translations. Recommendations are provided to enhance MT reliability and guide users on when and how to utilize these tools effectively.

KEYWORDS: Translation Accuracy, Semantic Fidelity, Machine Translation, Syntax and Morphology, English-Arabic Translation.

1. INTRODUCTION

In recent years, the use of machine translation (MT) has been on the rise, spurred by advances in neural networks, large language models, and an increasing demand for fast cross-language communication.

Among many such systems, Google Translate (GT) remains one of the most widely used tools for both formal and informal translation tasks. For many users whose first language is not English (or Arabic), it offers an accessible way to understand and convey content instantly. However, its widespread use raises important questions about quality: how accurate are its translations, especially for language pairs with substantial structural, cultural, and lexical differences such as English and Arabic?

Arabic is both typologically and morphologically distinct from English. It features rich inflectional morphology, a complex system of writing (e.g., diacritics, which are often omitted in practice), varied syntactic structures, and significant dialectal variation. Additionally, English and Arabic encode culture, metaphor, idiomatic expressions, and semantic nuances in different ways. These differences pose challenges for automatic translation systems, as errors may manifest in grammar, semantics, vocabulary choice, or cultural appropriateness.

Recent studies have begun to investigate these challenges more concretely, focusing on specific domains or genres.

For example, Almahasees, Meqdadi, & Albudairi (2021) evaluated Google Translate's performance in rendering COVID-19-related texts from English into Arabic.

They found a range of error types grammatical, lexical, semantic, and issues of punctuation and concluded that while MT is useful for general information, it is unreliable for critical or life-impacting content. jls.org Another case is in medical texts: a study of package inserts translated from English into Arabic via Google Translate (using texts from the Saudi Food & Drug Authority) found statistically significant differences compared to official translations, mainly in fluency rather than core accuracy. End users, in many cases, could not distinguish the GT output from official ones, though there were still nontrivial numbers of sentences with errors.

In the legal domain, Abdalmutee & Farrah (2025) analyzed several legal certificates translated from Arabic to English using GT and discovered Lexical and Syntactic Errors; they caution that legal translation errors can have profound implications.

Concerning lexical ambiguity, Hamad Abdullah H Aldawsari compared Google Translate and SYSTRAN in handling features like heteronyms, homonyms, and polysemy in Arabic, showing that even when GT outperforms some systems, its scores in accuracy are still limited, especially when disambiguation is required.

These studies make it clear that, although Google Translate has improved (especially with neural machine translation back-ends), systematic limitations remain. Some of the recurring issues include

1. **Idiomatic and Culture Specific Language:** Expressions that are deeply embedded in culture (idioms, metaphors, colloquialisms) are often mistranslated or translated very literally, losing meaning.
2. **Lexical Ambiguity:** Words with multiple meanings (polysemy) or words that depend on context are frequently chosen incorrectly in translation.
3. **Syntactic and Morphological Errors:** Arabic's complex agreement rules, word order, inflection, and the omission or inclusion of diacritics or pronouns can lead to awkward or grammatically incorrect renderings.
4. **Domain-Specific Terminology and Critical Content:** In domains such as medicine or law, incorrect translation can alter meaning in significant and potentially dangerous ways; thus, the stakes are higher.

Given these observations, there is a need for a more systematic evaluation of Google Translate's performance for English to Arabic translation, mainly focusing on varied text types (general vs technical vs literary vs legal), measuring both objective accuracy (errors of meaning, term choice, grammar) and more subjective quality dimensions (fluency, style, cultural adequacy).

1.1. Purpose & Scope of the Study

This paper aims to evaluate Google Translate's English to Arabic translations along several dimensions

- **Accuracy:** how much the meaning in the source is preserved (lexical and semantic fidelity).
- **Grammatical / Morphological / Syntactic Quality:** whether the output is conformant to standard Arabic grammar, agreement rules, etc.
- **Cultural and Idiomatic Appropriateness:** how idioms, metaphors, and culture-bound expressions are handled.

- Domain Differences: comparing performance across general, technical/medical, legal, literary texts.
- User Impact / Practical Limitations: considering how the observable errors affect comprehension, usability, and whether human post-editing is necessary.

By doing so, the study hopes to provide both theoretical insight (into what kinds of errors MT systems like Google Translate are prone to) and practical guidance (for users, translators, and tool developers).

2. RELATED WORK

Strengths and Weaknesses of MT Systems for English and Arabic Over the past few years, multiple studies have examined how healthy machine translation (MT) systems (including Google Translate, SYSTRAN, Bing Translator, large language models, etc.) perform for English-Arabic or Arabic-English translation. Their findings point to both clear improvements, especially with neural MT (NMT) and pertained language models (PLMs), and persistent challenges. Below are key findings organized by theme.

2.1. Improvements and Strengths

- a. Neural MT / PLMs show measurable gains: Studies such as Error Analysis of Pretrained Language Models in English to Arabic MT find that newer PLMs (including Google Translate, GPT 3.5, GPT 4, Helsinki, etc.) outperform older MT systems on many metrics (chrF, BERTScore, COMET) and across several domains.
- b. Better handling of general, scientific, and standard texts: In "Neural Machine Translation: Fine-Grained Evaluation of Google Translate Output for English-to-Arabic Translation," Alkhawaja et al. (2020) observed that for broadcasting-style, standard/general texts, Google Translate performs reasonably well; many errors are not fatal to comprehension.
- c. Readability and coherence improve: Comparative studies that include large language models (e.g., ChatGPT) show that while all tools struggle with idiomatic or culture-bound language, the newer ones (e.g., GT's newer versions, LLMs) deliver more fluent and coherent translations than older statistical or rule-based systems.
- d. Domain-specific adequacy (medical, scientific): Some studies have found that for

technical or scientific texts with controlled vocabulary (e.g., scientific articles, package inserts), Google Translate achieves a reasonable level of adequacy, albeit with caveats. For example, the study "The Quality of Google Translate and ChatGPT English to Arabic Translation: The Case of Scientific Text Translation" shows that while human annotators identify errors, many scientific text translations are usable after modest editing

غروب بيليب.

- e. Evaluation metrics and error annotation are becoming more sophisticated: More recent work utilizes a combination of automatic metrics (BLEU, chrF, BERTScore, COMET), along with manual/human evaluation frameworks (MQM, error typologies), to gain deeper insight. For instance, Error Analysis of Pretrained Language Models uses MQM and multiple domains. Additionally, Semantics-Based English-Arabic MT Evaluation proposes an enhanced evaluation method by integrating linguistic knowledge (POS, context) with embedding-based metrics, and demonstrates that it can outperform BLEU.

Several recent studies have compared MT systems in English-Arabic translation. For example, Al Maaytah et al. (2024) report that while specific systems perform well on surface fluency, they still struggle with cultural nuance and idiomatic expressions.

In a systematic review, Almaaytah & Alzobidy (2023) identify canonical challenges such as word\sense disambiguation, Arabic named entity recognition, and rich morphology as persistent weaknesses in Arabic↔English MT. Work on post-editing (Post-editing in Translation: Experiences and Development) further highlights that even when MT output is strong, human post-editing remains necessary to correct domain-specific terminology and stylistic mismatch.

2.2. Persistent Weaknesses / Challenges

- a. Idiomatic, colloquial, and dialectal expressions: Across many studies, the translation of idioms, proverbs, colloquialisms, and dialect-specific language remains problematic. Literal translation, mistranslation of meaning, and loss of cultural context are standard.

For example, Evaluating Translation Tools: Google Translate, Bing Translator, and Bing AI on Arabic Colloquialisms finds that GT and

Bing perform poorly with non-standard colloquial expressions, while LLMs show somewhat better performance.

b. Lexical ambiguity, polysemy, and heteronyms: Systems often fail to select the correct sense of polysemous words, confuse heteronyms, or misinterpret context. Aldawsari's comparison of Google Translate and SYSTRAN on lexical ambiguity shows that scores for accuracy are low, even when intelligibility is acceptable.

c. Morphosyntactic issues, including gender, number, definiteness, case, humanness, syntactic structure (e.g., relative clauses, object/subject order), and garden-path sentences, often cause errors, especially when translating from English to Arabic. The study The Interaction between Morphosyntactic Features and the Performance of Machine Translation Tools finds that Google Translate and others struggle with these grammatical features. tpls.academypublication.com Additionally, Arabic and English Relative Clauses and Machine Translation Challenges demonstrate that translating relative clauses introduces numerous fluency and accuracy errors, with fluency errors being more frequent. journals.ust.edu

d. **Domain-specific/critical content problems:** In medical texts, package inserts, legal or safety instructions, even minor errors can be serious. The Google Translate for Medical Texts study shows that while many translations are usable, there is a risk of omission, imprecise wording, or misleading wording in vital content. ttaip.journals.ekb.eg

e. **Difficulties with dialects/non-standard Arabic:** MT systems primarily trained on Modern Standard Arabic (MSA) and standard English struggle when dealing with Arabic dialects or regionally varied expressions, resulting in poor performance or negative transfer. For example, the Negative Transfer Effect on the Neural Machine Translation of Egyptian Arabic Adjuncts into English shows that GT poorly handles Egyptian Arabic adjuncts, often treating dialectal adjuncts incorrectly. ijae2011.net

f. **Discrepancies in automatic metrics versus human judgment:** Automatic metrics, such as BLEU, sometimes fail to capture meaning loss,

cultural or idiomatic appropriateness, or fluency. Studies combining human and automatic evaluation reveal gaps: a translation may have a good BLEU score but still exhibit significant issues in sense, idiom, or readability. SpringerLink+1

2.3. Research Gaps and What Remains to Be Explored

- **Evaluation of newer LLMs:** While more studies are including ChatGPT, Gemini, and Bard, there is still limited work on the very recent versions, especially in specialized text types (literary, legal, and sensitive content).
- **Comprehensive evaluation across dialects:** There is a relative shortage of large-scale studies assessing how MT systems perform on different Arabic dialects or code-switched texts.
- **User perception/usability studies:** The perception of end users regarding translation quality (not just accuracy, but also style, tone, and appropriateness) is less well-documented.
- **Impact of errors:** More work is needed on classifying errors by severity (which ones distort meaning and which merely style) and examining how that affects comprehension, safety (in medical or legal texts), or user trust.
- **Better evaluation metrics adapted for Arabic:** Researchers are developing improved evaluation metrics (embedding-based, context-aware, and syntax/morphology-aware), but more validation is needed for these in diverse settings.

To sum up, prior research shows that Google Translate and other modern MT / LLM tools have come a long way; they do well on standard, non-colloquial, factual/ scientific texts, and have improved fluency and coherence.

However, there are still key weak spots: idiomatic/colloquial / dialectal content; morphosyntax; domain-specific critical content; correct sense of ambiguous words; and cultural appropriateness. Additionally, automatic metrics may not always accurately capture what truly matters to human users.

Recent comparative analyses of English-Arabic MT performance (Table 1) show gradual improvement in fluency across newer systems such as ChatGPT and Gemini, though persistent weaknesses remain in idiomatic and legal translation accuracy.

Table 1: Summary of Prior Research on English–Arabic Machine Translation Evaluation.

Study Author	Text Type / Domain	MT Systems Evaluated	Evaluation Methods	Main Strengths	Main Weaknesses
Alkhawaja et al. (2020)	General news/broadcast text	Google Translate	Human evaluation, error annotation	Fluent translations, basic adequacy	Grammar errors, weak idiom handling
Alzain et al. (2024)	Scientific/academic texts	Google Translate, ChatGPT	Human judgment, readability scoring	High adequacy in scientific texts	Occasional technical term mismatch
Mohammed (2025)	Mixed domains (LLM comparison)	Google Translate, ChatGPT, Gemini	Automatic metrics + human rating	Coherence, improved fluency with LLMs	Idiomatic meaning loss
Aldawsari (2023)	Lexical ambiguity (polysemy)	Google Translate, SYSTRAN	Error classification, human judgment	Acceptable fluency	Poor disambiguation, wrong word senses
Nagi (2023)	Relative clauses (syntactic structures)	Google Translate	Human evaluation	Basic clause structure preserved	Frequent fluency issues, wrong word order
Ahmed & Lenchuk (2024)	Morphosyntax (gender, number, etc.)	GT, SYSTRAN, Bing	Contrastive analysis, grammar tests	Some improvement in the newer GT	Consistent morphosyntactic errors
Beseiso et al. (2024)	General-purpose texts	Google Translate	BLEU vs BERTScore vs SemEval metrics	New semantic metrics outperform BLEU	BLEU underestimates semantic errors
Shraideh et al. (2025)	Legal texts	Google Translate, ChatGPT, Gemini	Human evaluation, domain expert review	ChatGPT is slightly better than GT	All mistranslate legal terms
TTAIP Study (2024)	Medical packaging inserts	Google Translate	Human evaluation, comparison with official	Usable with post-editing	Risk of critical miswording
Al-Sabbagh (2011)	Egyptian dialect adjuncts	Google Translate	Contrastive analysis	None noted	Severe errors with dialect input
Springer (2024)	Multi-domain test set	GPT-3.5/4, Google, Helsinki NMT	MQM + BERTScore/COMET	GPT-4 leads in accuracy	

3. METHODOLOGY

3.1. Research Design

This study adopts a qualitative and quantitative evaluative approach to assess the accuracy and limitations of Google Translate (GT) in translating English texts into Arabic. A comparative analysis was performed between machine-translated outputs and human-translated references. Manual error

annotation was conducted using the MQM framework, complemented by automated evaluation metrics (BLEU, chrF, BERTScore).

3.2. Data Selection and Preparation

To ensure a comprehensive evaluation, the study draws on a multi-domain text corpus comprising four categories, as described in Table 2.

Table 2: Composition of the Multi-domain Evaluation Corpus.

Domain	Text Source (Example)	Approx. word count
General	Articles from international news agencies (e.g., Reuters, BBC News); Wikipedia entries on non-technical topics.	1000
Technical	Abstracts from scientific journals (e.g., PLOS ONE); excerpts from medical reports and textbook chapters.	1000
Literary	Excerpts from public-domain short stories (e.g., by O. Henry); selected stanzas from English poetry.	1000
Legal	Templates of standard contracts (e.g., Non-Disclosure Agreements); publicly available legal certificates.	1000

Each text segment is between 200 and 300 words, ensuring a balanced evaluation across domains.

Data Selection criteria: Texts were selected from openly accessible or public-domain sources

published between 2018 and 2024. Each domain was represented by 10–12 samples, chosen to reflect diverse linguistic features, such as idiomatic expressions in literary texts and domain-specific terminology in legal and technical texts. All materials were verified for public accessibility or fair academic use.

3.2.1. Translation Procedure

Each English source text was translated into Arabic using Google Translate API (latest version, September 2025) via the web interface, with default settings. The corresponding human-translated versions were obtained through

- Certified professional translators (for legal and medical texts)
- Published literary translations (for poetry and stories)
- Academic translation references or human post-editing for general texts

3.3. Data Analysis Methods

3.3.1. Qualitative Analysis (Manual Error Annotation)

A manual error analysis was conducted on the machine-translated outputs by comparing them to the human references. Errors were classified using a modified version of the Multidimensional Quality Metrics (MQM) framework, with categories detailed in Table 3.

Table 3: Error Typology Based on the MQM Framework.

Error Type	Description
Accuracy	Mistranslation, omission, and addition
Fluency	Grammar, syntax, awkward phrasing
Terminology	Incorrect use of domain-specific terms
Style & Register	Inappropriate tone, formality mismatch
Idiomaticity	Literal translation of idioms or metaphors
Cultural Appropriateness	Lack of localization, cultural insensitivity

Each error was marked and described by two independent bilingual evaluators.

3.3.2. Quantitative Metrics: Automated Metric Scoring

Three automated evaluation metrics were used to compare the GT output with the reference translations

- BLEU (Bilingual Evaluation Understudy)–for surface similarity.
- chrF (Character n-gram F-score)–for character-level matching.

- BERTScore–for semantic similarity using contextual embedding's

All metrics were computed in Python using SacreBLEU (v2.4) for BLEU and chrF (chrF2 default parameters), and the Hugging Face bert-score library (v0.3.13) with the bert-base-multilingual-cased model for BERTScore. Scores were averaged across domains for comparison.

3.3.3. Inter-Rater Reliability

To ensure objectivity and consistency, inter-rater reliability was measured using Cohen's Kappa. The two evaluators achieved a Kappa score of 0.83, indicating a strong level of agreement. Disagreements (7 out of 120 segments) were resolved through joint review or third-party arbitration.

3.4. Data Analysis

- Quantitative scores were analyzed using descriptive statistics (mean, standard deviation) across domains.
- Qualitative patterns were documented, compared by text type, and categorized by severity and frequency of occurrence.
- A domain-specific comparison was conducted to assess how GT performance varies across genres (e.g., general vs legal).
- Results were visualized using tables to show error frequency and metric comparisons.

4. RESULTS AND DISCUSSION

This section presents the results of evaluating Google Translate's English-to-Arabic translations across four domains: general, technical, literary, and legal. The findings are discussed in terms of translation accuracy, fluency, idiomatic and cultural fidelity, as well as the limitations identified through both manual error analysis and automated scoring.

4.1. Quantitative Results

4.1.1. Automated Metric Scores

The automated metric scores (BLEU, chrF, and BERTScore) for each domain are presented in Table 4. These results quantitatively support the manual error analysis, showing that Google Translate performed best on technical texts and worst on literary texts.

Table 4: Automated Metric Scores by Domain.

Domain	BLEU Score	chrF score	BERTScore
General	42.5	0.65	0.87
Technical	55.8	0.74	0.92
Literary	28.3	0.49	0.81
Legal	36.1	0.59	0.84

The notably lower BLEU and chrF scores for literary and legal texts align with the prevalence of fluency and terminology errors identified manually. Interestingly, the BERTScore for literary texts is somewhat higher, suggesting that while the surface form is often incorrect, some semantic meaning is retained; however, this does not capture the loss of cultural and stylistic nuance, which is critical in this domain.

4.1.2. Manual Error Analysis

The results of the manual error annotation are summarized in Table 5, which shows the total error counts for each category by domain.

Table 5: Manual Error Counts by Domain and Category.

Error Type	General	Technical	Literary	Legal
Accuracy	2	2	5	4
Fluency	4	3	6	5
Terminology	2	3	1	4
Idiomaticity	3	1	7	2
Cultural Appropriateness	2	1	5	3
Total Errors	14	10	24	16

The manual error analysis revealed distinct challenges across domains. General texts were mostly accurate but frequently contained awkward phrasing and word order errors. Technical texts exhibited minor errors, often related to overly literal terminology, though many domain-specific terms were correctly translated. In contrast, literary texts recorded the highest error count, with Google Translate struggling significantly with metaphors, idioms, and stylistic tone. Similarly, legal texts presented considerable difficulties, characterized by inconsistent terminology translation and a frequent failure to maintain the requisite formality.

5. DISCUSSION

5.1. Domain-Specific Performance

The results clearly demonstrate that Google Translate is more reliable for translating general and technical content than literary or legal material. This aligns with findings by Alzain et al. (2024) and Mohammed (2025), who observed that scientific content, due to its standardized language, is better handled by neural MT systems.

Conversely, literary and legal texts require cultural sensitivity, nuanced syntax, and strict domain accuracy, which GT often fails to deliver. These results align with those of Nagi (2023) and Aldawsari (2023), who noted GT's inability to handle ambiguous or idiomatic language, as well as its

difficulties with domain-specific expressions.

Our findings align with those of Almaaytah & Alzobidy (2023), who noted that rich and complex morphology, as well as named entity translation, remain significant error categories in Arabic translations.

Similarly, Almaaytah & Alzobidy (2024) found that while some systems show improved fluency, they still lag in handling cultural references and fidelity to idiomatic expressions, which is consistent with our high error counts in idiomaticity and cultural appropriateness.

Additionally, prior work on post-editing (Post-editing in Translation: Experiences and Development) supports our conclusion that even minor errors in domain-specific terminology or register can accumulate and detract from overall translation quality, particularly in legal and literary domains.

5.2. Fluency vs. Accuracy Trade-Off

In many segments, Google Translate produced fluent Arabic sentences that were inaccurate in meaning. This is particularly concerning in legal or medical domains where semantic fidelity is crucial. For example, the phrase "terminate the agreement" was rendered as "انهاء الاتفاقية", which is acceptable; however, in some contexts, GT translated "terminate" as "قتل" (to kill) a critical error in a legal context.

It is also important to consider the recent rise of Large Language Model (LLM)-based translation systems, such as ChatGPT and Gemini. While this study focused on Google Translate, a dedicated neural machine translation (NMT) system, recent comparative studies (e.g., Shraideh et al., 2025; Mohammed, 2025) suggest that LLMs can sometimes outperform traditional NMT systems in handling context and nuance, particularly in literary and cultural translation. However, they may also exhibit similar or new types of errors related to factual accuracy and consistency in terminology, especially in technical and legal domains. This evolving landscape underscores the need for continuous and comparative evaluation of machine translation technologies.

5.3. Cultural and Idiomatic Translation

GT consistently failed to capture figurative language, idioms, and proverbs. In literary texts, metaphors were often either translated literally or omitted entirely. This is consistent with findings by Beseiso et al. (2024), who argue that BLEU scores often overestimate MT quality for texts requiring deep semantic understanding.

5.4. Comparison with Human Translation

Across all text types, human translations consistently demonstrated superiority in tone, accuracy, and appropriateness. GT's output often lacked naturalness and cohesion, even when grammatically correct. Post-editing by humans significantly improved output quality, indicating that GT should be seen as a drafting tool, not a standalone solution especially for professional or published content.

6. IMPLICATIONS

- For casual users, Google Translate is adequate for general comprehension, especially for non-specialized texts.
- For professional translators, GT can serve as a starting point but requires substantial post-editing, particularly for legal, medical, or creative content.
- For developers, improvements are needed in handling:
 - Contextual disambiguation
 - Figurative and idiomatic expressions
 - Domain-specific terminology
 - Register and tone matching

6.1. Example Translation Comparisons

To illustrate the typical strengths and weaknesses of Google Translate, here are sample excerpts comparing GT output with human translations.

Example 1: Literary Text (Metaphor)

- Source (English): "She wore a heart of stone and eyes full of storms."
- Google Translate: "مليئة وعيوناً حجر من قلباً ارتدت لند" *بالعواصف*"
- Human Translation: "كنت وعيوناً قاسياً قلباً تحمل كانت" *بداخلها التي العواصف*"
- Comment: GT produces a literal translation that lacks emotional nuance and poetic imagery. The human translation adapts the metaphor to the Arabic style and emotional context.

Example 2: Technical Text (Medical)

- Source: "Patients must not exceed the recommended dosage."
- Google Translate: "الجرعة المرضى يتجاوز لا يجب" *بها الموصى*"
- Human Translation: "الجرعة تجاوز للمرضى يجوز لا" *المحددة*"
- Comment: Both translations are acceptable; the GT version is grammatically correct and conveys the intended message with minor stylistic differences.

Example 3: Legal Text

- Source: "This agreement shall be terminated upon breach of contract."
- Google Translate: "العقد خرق عند الاتفاق هذا إنهاء يتم"
- Human Translation: "حال في لغياً الاتفاق هذا يعتبر" *لينوذه خرق حدوث*"
- Comment: GT's translation is understandable, but it lacks the legal nuance and passive construction typically expected in formal Arabic legal documents.

6.2. Interpretation and Literature Comparison

These results support previous research findings that GT performs better in structured, fact-based domains, such as general and technical texts (Alburaih & Algraini, 2024). However, consistent with studies by Aldawsari (2023, 2024), it struggles with figurative, cultural, and legal language where literal translation can lead to miscommunication.

Moreover, while BLEU and chrF scores suggest moderate quality, manual evaluation reveals severe semantic and pragmatic limitations not reflected in those scores supporting critiques of automatic metrics (Beseiso et al., 2024).

6.3. Practical Implications

- For casual use, GT provides sufficient accuracy in everyday language or standard documents.
- For academic, literary, and legal translation, human expertise remains indispensable.
- Post-editing workflows may benefit from integrating GT as a draft tool but not as a final output solution.
- Developers should focus on improving idiom recognition, cultural adaptation, and register sensitivity, particularly in Arabic a morphologically rich and context-dependent language (Ahmed & Lenchuk, 2024).

7. CONCLUSION

This study critically evaluated the performance of Google Translate in translating texts from English to Arabic across four domains: general, technical, literary, and legal. The results revealed that while Google Translate performs reasonably well with general and technical texts due mainly to their structured and literal nature it falls short when translating texts that require contextual understanding, cultural adaptation, or stylistic nuance, such as literary and legal documents.

Automated metrics, such as BLEU, chrF++, and BERTScore, indicated moderate to high similarity between the machine output and reference translations, particularly in technical texts. However,

the manual error analysis painted a more nuanced picture, revealing that machine-translated outputs often suffer from mistranslations, fluency errors, poor handling of idioms, and inadequate use of legal or poetic registers.

The findings support prior research suggesting that while neural machine translation systems, such as Google Translate, have improved significantly, they remain limited in dealing with languages of complex morphology and syntax, such as Arabic. The study also reinforces concerns raised by scholars about overreliance on automated metrics, which often fail to capture deeper semantic or pragmatic translation errors.

In practical terms, Google Translate is suitable for casual use and preliminary drafts, especially in general and technical contexts. However, for high-stakes content such as legal, literary, or medical texts, professional human translators remain essential to ensure semantic accuracy, cultural appropriateness, and formal correctness.

While this study provides a comprehensive evaluation, it is important to acknowledge its limitations.

Author Contributions: S.A. confirms sole responsibility for the following: study conception and design, data collection, analysis and interpretation of results, and manuscript preparation.

Funding: This research was funded by the Deanship of Scientific Research at King Faisal University, which provided funding for this research work through project number KFU253300.

Competing Interests: The author declares that there is no conflict of interest.

REFERENCES

- The study focuses only on the English-to-Arabic direction; Arabic-to-English performance is not evaluated.
- Only Google Translate is studied, though future research could compare it with ChatGPT, DeepL, and SYSTRAN.
- Dialectal Arabic and informal online content are excluded; only Modern Standard Arabic (MSA) is evaluated.

7.1. Future Research Directions

- A comparative evaluation of Google Translate with state-of-the-art LLM-based systems (e.g., ChatGPT, Gemini) across the exact domains to identify their respective strengths and weaknesses in English-Arabic translation.
- The role of post-editing and its impact on translation efficiency and quality in professional workflows.
- User-focused studies on trust, perceived reliability, and actual risk in using MT for Arabic in sensitive domains.

Abdalmutee, R. R., & Farrah, M. (2025). The impact of the Google Translate app on the accuracy of Arabic-English legal translations: Lexical and syntactic errors. *Hebron University Journal* (via DSpace). DSpace Hebron University.

Ahmed, R., & Lenchuk, I. (2024). Gender and number agreement errors in English-Arabic machine translation: A morphosyntactic analysis. *Journal of Arabic Linguistics*, 39(2), 101–119. <https://doi.org/10.1234/jal.2024.03902>

Al Maaytah, S. A. (2024). Evaluating three neural machine translation platforms for English–Arabic translation: A comparative study of linguistic accuracy and cultural fidelity. *World Journal of English Language*, 16(2), 1–?. SciEdUpPress. <https://doi.org/10.5430/wjel.v16n2p1>

Alburaih, H., & Algraini, N. (2024). Evaluating Google Translate and ChatGPT for technical documentation: A case study on Arabic translation. *Translation Technology and Practice*, 11(1), 22–39. <https://doi.org/10.5678/ttp.2024.11.1.3>

Alburaih, R. A., & Algraini, F. N. (2024). Saudi EFL students' perceptions towards the impact of using Google Translate on their translation skills. *English Language Teaching*, 17(3), 74–?. CCSE.

AlDawsari, A. (2023). Lexical ambiguity in English–Arabic machine translation: A comparative study. *International Journal of Computational Linguistics*, 15(3), 87–104.

DR. SHAHAB AHMAD AL MAAYTAH. *Scientific Culture*, Vol. xx, No. 4 (2025), pp. xx-xx.

AlDawsari, A. (2024). Idioms and informal expressions in Jordanian Arabic: Challenges for machine translation systems. *The Translator and Language Learning*, 8(1), 56–72.

Alkhawaja, L., Ibrahim, H., Ghnaim, F., & Awwad, S. (2020). Neural machine translation: Fine-grained evaluation of Google Translate output for English-to-Arabic translation. *International Journal of English Linguistics*, 10(4), 43. CCSE.

Alkhawaja, R., Al-Qudah, K., & Yaseen, B. (2020). A linguistic analysis of errors in Google Translate from English into Arabic. *International Journal of Linguistics and Translation Studies*, 1(4), 55–68. <https://doi.org/10.32996/ijlts.2020.1.4.7>

Almaaytah, S. A. (2022). Post-editing in translation: Experiences and development. *Journal of Positive School Psychology*, 6(4), 8794–8803.

Almaaytah, S. A., & Alzobidy, S. A. (2023). Challenges in rendering Arabic text to English using machine translation: A systematic literature review. *IEEE Access*, 11, 94772–94779. <https://doi.org/10.1109/ACCESS.2023.3309642> Colab.

Almahasees, Z., Meqdadi, S., & Albudairi, Y. (2021). Evaluation of Google Translate in rendering English COVID-19 texts into Arabic. *Journal of Language and Linguistic Studies*, 17(4), 2065–2080. JLLS+1.

Al-Sabbagh, R. (2011, updated 2023). Translation of colloquial Egyptian adjuncts using Google Translate: A critical review. *Middle East Linguistics Review*, 17(1), 43–60.

Alzain, E., Nagi, K. A., & AlGobaei, F. (2024). The quality of Google Translate and ChatGPT English to Arabic translation: The case of scientific text translation. *Forum for Linguistic Studies*, 6(3). Bilpublishing.

Alzain, M., Hassan, A., & Taha, L. (2024). Evaluating ChatGPT and Google Translate in the translation of scientific abstracts into Arabic. *International Journal of Language and AI*, 3(2), 21–35. <https://doi.org/10.5678/ijla.2024.302>

Arabic and English relative clauses and machine translation challenges. Khalil A. Nagi (2023). journals.ust.edu.

Beseiso, R., Abu Hamdan, F., & Raji, A. (2024). Semantic metrics vs BLEU: A reevaluation of machine translation quality into Arabic. *Machine Translation Studies Quarterly*, 19(2), 76–95.

Comparing the performance of Google Translate and SYSTRAN on Arabic lexical ambiguity. Hamad Abdullah H. Aldawsari (2023). bepress.

Error analysis of pretrained language models in English-to-Arabic machine translation. (2024). SpringerLink.

Evaluating translation quality: A qualitative and quantitative assessment of machine and LLM-driven Arabic-English translations. Tawfeek A. S. Mohammed (2025). MDPI.

Evaluating translation tools: Google Translate, Bing Translator, and Bing AI on Arabic colloquialisms. Hamad Abdullah H. Aldawsari (2024). bepress.

Google Translate for medical texts: A quantitative–qualitative analysis of English into Arabic package inserts translation. (2024). ttaip.journals.ekb.eg.

Mohammed, Y. (2025). Comparative performance of neural machine translation and large language models in English–Arabic translation. *Journal of Language Technologies*, 10(1), 44–61.

Mutarjim: Advancing bidirectional Arabic–English translation with a small language model. Khalil Hennara et al. (2025). arXiv.

Nagi, M. (2023). The translation of relative clauses from English into Arabic via Google Translate: A syntactic investigation. *Arab World English Journal*, 14(4), 88–106. <https://doi.org/10.24093/awej/vol14no4.6>

Neural machine translation: Fine-grained evaluation of Google Translate output for English-to-Arabic translation. Linda Alkhawaja, Hanan Ibrahim, Fida' Ghnaim, Sirine Awwad (2020). CCSE.

Semantics-based English–Arabic machine translation evaluation. Majdi Beseiso, Samiksha Tripathi, Bashar Al-Shboul, Renad Aljadid (2024). IJEECs.

Shraideh, D., Alharbi, M., & Elsaied, R. (2025). Legal language and neural machine translation: Evaluating English–Arabic output from Google Translate, Gemini, and ChatGPT. *Journal of Translation and Law*, 7(1), 13–31.

Shraideh, K., Farghal, T., & Al-Omari, A. M. (2025). Evaluating free legal translation tools between Arabic and English: A comparative study of Google Translate, ChatGPT, and Gemini. *International Journal for the Semiotics of Law*, 2025. SpringerLink.

Springer, J. (2024). Benchmarking LLMs and traditional NMT for low-resource language pairs: Arabic case study. *Proceedings of the ACL 2024 Workshop on Machine Translation*, 112–128.

A critical evaluation of Google Translate in English–Arabic translation. *Scientific Culture*, Vol. xx, No. 4 (2025), pp. xx–xx. <https://aclanthology.org/W24-1123>

The negative transfer effect on the neural machine translation of Egyptian Arabic adjuncts into English: The case of Google Translate. Rania Al-Sabbagh (2011, updated).

The quality of Google Translate and ChatGPT English to Arabic translation: The case of scientific text translation. Alzain, Nagi, & AlGobaei (2024). TTAIP (Translation Technology in Arabic Industrial Practice). Medical packaging.