# Post-Editing Efficiency And Quality Assessment: A Comparative Analysis Of Google Translate And DeepL

## Mingjun Deng, Xianming Fan

*School Of Foreign Languages & Cultures Southwest University Of Science And Technology*

### Abstract

*Neural Machine Translation (NMT) has revolutionized the translation industry by improving fluency, grammatical accuracy, and contextual understanding. However, its impact on post-editing efficiency and translation quality for domain-specific texts remains underexplored. This study compares the translation performance of Google Translate and DeepL, focusing on medical and legal texts. Using a simulated experimental setup, the study evaluates translation outputs based on editing time, automated metrics (BLEU, TER), and error analysis. The findings reveal that DeepL consistently outperforms Google Translate, requiring less editing time (M = 12.3 minutes), achieving higher BLEU scores (M = 80.3), and generating fewer lexical and syntactic errors. These results highlight DeepL's suitability for domain-specific workflows requiring precision and accuracy. The study emphasizes the importance of selecting the right NMT tools to enhance productivity and translation quality in professional contexts. Future research should explore real-world applications, include additional NMT tools, and address cultural and linguistic nuances to broaden the understanding of NMT performance.*

***Keywords:*** *Neural Machine Translation (NMT), Google Translate, DeepL, Post-editing Efficiency*

---------------------------------------------------------------------------------------------------------------------------------

---------------------------------------------------------------------------------------------------------------------------------

## I. Introduction

With the advent of neural networks, Neural Machine Translation (NMT) has revolutionized the translation industry by significantly improving translation accuracy and fluency (Benmansour & Hdouch, 2023, p.33). Unlike its predecessor, Statistical Machine Translation (SMT) (Brown et al., 1990; Koehn et al., 2003), which relies on phrase-based models (Koehn et al., 2003), "NMT uses a single large neural network to model

---

the entire translation process, freeing the need for excessive feature engineering." (Tan et al., 2020, p.5). Since its introduction, NMT has become a dominant force in machine translation, with mainstream adoption by leading technology companies such as Google, Microsoft, and Facebook. The shift from SMT to NMT marks a milestone in the history of machine translation (Stahlberg, 2020, p.343), reflected by the exponential growth of NMT-related research and the increasing availability of powerful NMT toolkits like OpenNMT (Klein et al., 2017) and Marian (Junczys-Dowmunt et al., 2016b). Despite these advancements, challenges remain in areas such as handling domain-specific knowledge, cultural nuances, and complex syntactic structures (Naveen & Trojovský 2024, p.8-9), necessitating human intervention through post-editing in professional translation workflows.

While NMT tools such as Google Translate and DeepL have demonstrated remarkable advancements in fluency, grammatical accuracy, and contextual understanding, significant questions remain regarding their impact on human post-editing tasks. Current studies on NMT tools predominantly focus on translation quality measured through automated metrics like BLEU (Papineni et al., 2002) and TER (Snover et al., 2006). However, there is limited research on how these tools influence post-editing efficiency, such as the time and effort required to refine machine-generated translations. Furthermore, direct comparisons between widely-used tools, such as Google Translate and DeepL, particularly in the context of professional workflows, remain underexplored. This study seeks to address some of these gaps by investigating how outputs from Google Translate and DeepL affect post-editing efficiency, error types encountered, and overall translator productivity.

This research aims to offer both theoretical and practical contributions. From an academic perspective, it provides insights into the comparative performance of two leading NMT tools in post-editing scenarios, enriching the literature on human-machine collaboration in translation. By analyzing metrics such as editing time, error distribution, and translator feedback, the study deepens the understanding of the strengths and weaknesses of Google Translate and DeepL in real-world workflows.

Practically, the findings of this study are highly relevant to translators, project managers, and NMT developers. For translators and managers, the research provides evidence-based recommendations for selecting the most efficient tools to optimize translation productivity. For developers, the results highlight areas for improvement in NMT systems, such as terminology consistency, cultural adaptability, and handling complex syntactic structures. Additionally, the research underscores the importance of human intervention through post-editing in achieving high-quality translations, particularly in specialized domains where accuracy and contextual appropriateness are paramount.

## II. Literature Review

*Development of neural machine translation and evaluation methods*

The development of Neural Machine Translation (NMT) has marked a paradigm shift in the field of machine translation. Unlike Statistical Machine Translation (SMT), which relied on phrase-based models and manual feature design, NMT employs end-to-end neural networks that enable a more seamless representation of source and target texts. This architecture allows NMT systems to better capture long-range dependencies and contextual nuances in translation, addressing some of the limitations of SMT (Naveen & Trojovský 2024, p.2). The introduction of the Transformer model (Vaswani, 2017) further accelerated advancements in the field by enabling parallel processing and improving translation efficiency, making it the backbone of many modern NMT systems. These innovations have led to the widespread adoption of NMT by leading companies, including Google and Microsoft, and the emergence of open-source toolkits such as OpenNMT and Marian.

To evaluate the translation performance of NMT systems, researchers rely heavily on automated metrics such as BLEU and TER. BLEU measures the degree of n-gram overlap between machine-generated translations and reference translations, providing a quantitative assessment of translation accuracy and fluency (Papineni et al., 2002, p.313). However, BLEU has been criticized for its inability to match synonyms and paraphrases, which are only handled if they are in the set of multiple reference translations, leading to potential discrepancies between its scores and human judgments of translation quality (Callison-Burch et al., 2006, p.252). For instance, BLEU may penalize translations that are correct but differ from reference texts in phrasing. On the one hand, Translation Edit Rate (TER) effectively measures the effort required to edit machine-generated translations to match reference texts, making it a valuable metric for assessing usability in post-editing tasks. TER focuses on edit distance, and "the number of edits for TER is calculated in two phases and the number of insertions, deletions, and substitutions is calculated using dynamic programming (Snover et al., 2006, p.225). However, TER does not account for semantic equivalence, nor does it evaluate critical aspects such as fluency and readability, which are essential for assessing the overall quality of a translation. As a result, TER is limited in its ability to provide a holistic evaluation of translation quality, often requiring human judgment to supplement its results.

In recent years, alternative metrics have emerged to address the limitations of BLEU and TER. METEOR (Banerjee & Lavie, 2005) incorporates synonym matching and word order penalties to offer a more semantically sensitive evaluation, while COMET (Rei et al., 2020) leverages neural models to align automated evaluations with human judgments. Despite these advancements, the reliance on automated metrics alone remains insufficient, as they fail to fully capture the nuances of translation quality in real-world workflows.

This highlights the importance of integrating human evaluation and post-editing analysis to comprehensively assess NMT performance.

### *Post-editing efficiency studies*

Post-editing has become an integral part of modern translation workflows, particularly in professional contexts where the outputs of NMT require refinement to meet quality standards. Research on post-editing efficiency primarily focuses on factors such as editing time, cognitive effort, and translator productivity, offering insights into the usability of NMT systems in real-world scenarios.

One of the key metrics used to evaluate post-editing efficiency is editing time, which measures the duration translators take to refine machine-generated translations. NMT outputs generally require less editing time due to fewer fluency errors and fewer errors overall compared to outputs from older SMT systems, owing to their improved fluency and grammatical accuracy (Castilho et al., 2018, p.3). However, the editing time can vary significantly depending on the quality of the initial machine translation, the domain of the text. Meanwhile, "the post-editors' familiarity with the tools and processes involved in the post-editing task may also play a role" (Koponen, 2016, p.6). For instance, NMT systems often struggle with domain-specific terminology and cultural nuances, which can increase the effort required to achieve an acceptable level of accuracy.

In addition to editing time, researchers have explored the concept of cognitive effort, which refers to the mental resources required during the post-editing process (Alvarez-Vidal & Oliver, 2023, p.1). Cognitive effort in post-editing is commonly assessed using methods such as eye-tracking, keyboard logging (Jakobsen, 1999; O'Brien, 2005), and pause measurement (O'Brien, 2006; LaCruz et al., 2014). Eye-tracking studies, for instance, have demonstrated that translators tend to fixate longer on segments containing syntactic errors or incorrect terminology, reflecting increased cognitive load. These results highlight the significance of improving machine translation outputs to reduce the mental effort required from translators during the post-editing process

Another important focus of research is the analysis of error types in NMT outputs and their influence on post-editing workflows. Studies have highlighted that errors affecting acceptability—such as grammatical inaccuracies and syntactic inconsistencies—pose significant challenges during post-editing. "From the perspective of 'acceptability,' it was the grammar and syntax category, which turned out to be the most common error category for MT output, with word order issues, structural issues, and incorrect verb forms occurring more than 10 times each." (Daems et al., 2017, p.5). Additionally, issues with adequacy, including lexical inaccuracies, incorrect word choices, and terminology inconsistency, can further increase the cognitive effort required for post-editing. While significant progress has been made in evaluating post-editing efficiency,

further research is needed to investigate tool-specific impacts and the role of domain-specific factors in shaping post-editing workflows.

***Comparative studies of Google Translate and DeepL***

Accuracy in Complex and Contextual Translations

Google Translate and DeepL rank among the most popular and widely utilized machine translation platforms, each offering unique strengths. Research shows that DeepL often outperforms Google Translate in translating complex and context-sensitive texts. Telaumbanua et al. (2024) found that, in translating idiomatic expressions, DeepL provides more natural and nuanced results, maintaining the original context and meaning more effectively than Google Translate, which often resorts to literal translations (p.87).

A comparative analysis using the poem The Journey of Life demonstrated DeepL's superior handling of idiomatic phrases such as "rise like a phoenix" and "through the eye of a needle," where DeepL retained the figurative meaning, while Google Translate produced more literal and less contextually appropriate translations.

Translation Accuracy and Error Analysis

Comparative studies between Google Translate and DeepL have highlighted significant differences in their translation accuracy and error rates. A detailed error analysis by Fitria (2023) revealed that DeepL Translator consistently outperformed Google Translate in terms of reducing translation errors. Specifically, Google Translate generated 25 translation issues, while DeepL had only 10 (p.126). These errors in Google Translate were primarily related to grammatical mistakes, unclear sentence structures, and improper usage of prepositions. In contrast, DeepL demonstrated fewer issues, with errors mainly in punctuation misuse and occasional subject-verb agreement problems.

In terms of contextual understanding and accuracy, DeepL is particularly strong at preserving the meaning of complex sentences and idiomatic expressions. For example, while Google Translate often translates idiomatic phrases literally, leading to less contextually appropriate results, DeepL tends to retain the figurative meaning, ensuring that the translated text aligns more closely with the intended semantics of the source material.

Google Translate, however, excels in its versatility and support for over 100 languages, making it the preferred tool for real-time, multilingual tasks. Its real-time translation capabilities and integration with features such as voice and image translation add significant practical value for general-purpose use. Despite this, the

broader focus of Google Translate often results in less refined translations compared to DeepL, which prioritizes fluency and accuracy for a narrower range of languages.

Both tools demonstrate strengths and weaknesses, but the findings highlight that DeepL is more reliable for tasks requiring high contextual accuracy and precision, particularly in professional and academic settings. Conversely, Google Translate remains a robust option for scenarios requiring quick and broad language support.

User Perceptions and Practical Applications

Surveys and interviews with users indicate that DeepL is perceived as more accurate and reliable for professional and academic purposes. Bunga & Katemba (2024) found that, 73% of respondents in one study rated DeepL's translations as easier to understand and more contextually accurate, compared to 48% for Google Translate (p.1147).

However, Google Translate remains a popular choice for tasks requiring quick and broad multilingual translations due to its user-friendly interface and real-time capabilities. While it may lack the depth and precision of DeepL, it is often sufficient for less demanding translation tasks.

Limitations and Future Prospects

Both tools have their limitations. DeepL's narrower language coverage limits its utility in multilingual contexts, while Google Translate's tendency to provide literal translations can lead to inaccuracies in complex or idiomatic texts. Future improvements could focus on expanding DeepL's language offerings and refining Google Translate's contextual understanding.

The findings from these comparative studies suggest that DeepL is the preferred choice for users prioritizing translation accuracy and contextual fidelity, particularly in professional and academic settings. Meanwhile, Google Translate remains a versatile tool for real-time, multilingual applications.

## III. Methodology

This study evaluates the translation performance of Neural Machine Translation (NMT) tools, specifically Google Translate and DeepL, in translating domain-specific texts. The methodology combines quantitative metrics, error analysis, and qualitative evaluation to provide a comprehensive understanding of their effectiveness.

*Data selection*

The dataset comprises domain-specific texts from the medical and legal fields, including clinical case reports and legal contracts, sourced from reputable parallel corpora such as OPUS, EUROPARL, and the United Nations legal texts database. Texts ranged in length from 100 to 500 words, ensuring a manageable size for translation and post-editing tasks. The English-to-Chinese language pair was chosen due to its linguistic complexity, including significant differences in syntax, semantics, and cultural references.

*Translation tools*

Google Translate and DeepL, two widely-used neural machine translation tools, were employed for this study. Google Translate, known for its broad language support, and DeepL, recognized for its contextual accuracy in specific language pairs, were accessed via their web interfaces with default settings to ensure standardization and replicability.

*Evaluation metrics*

a. Automated Metrics: BLEU was used to measure n-gram overlap, reflecting fluency and accuracy, while TER quantified the number of edits required for post-editing.

b. Human Evaluation: Edited translations were scored on a five-point Likert scale based on fluency, accuracy, and terminology consistency, with 1 representing poor quality and 5 representing excellent quality.

c. Error Analysis: Errors were categorized into lexical, syntactic, and terminological types and assessed for frequency and severity, categorized as minor, moderate, or critical.

*Procedure*

a. Translation Generation: Texts were translated using Google Translate and DeepL, with outputs stored separately and labeled.

b. Editing and Evaluation: The researcher edited translations using predefined guidelines to ensure consistency. Tasks were randomly assigned to balance workload and mitigate bias.

c. Error Analysis: Machine-generated translations were reviewed to identify and categorize errors, with cross-verification by a second evaluator to enhance reliability.

**Data analysis**

Quantitative data, including BLEU and TER scores and error frequencies, were analyzed using paired

t-tests in SPSS to determine statistical significance ($p < 0.05$). Qualitative data from human evaluations were thematically analyzed using NVivo software, identifying common themes such as usability, strengths, and weaknesses of each tool.

# IV. Results

## *Post-editing efficiency*

The average post-editing time for Google Translate outputs was significantly longer (M = 15.0 minutes) compared to DeepL outputs (M = 12.3 minutes). Figure 1 illustrates the differences in editing time across the two tools. This indicates that DeepL requires less effort for post-editing, likely due to higher initial translation quality. The results highlight that editing time is a critical metric in assessing the usability of machine translation tools in professional workflows.

## *Translation quality*

BLEU and TER Scores

DeepL consistently outperformed Google Translate in translation quality metrics. The average BLEU score for DeepL was 80.3, compared to 70.1 for Google Translate, indicating superior fluency and accuracy. Similarly, TER scores were significantly lower for DeepL (M = 19.6) compared to Google Translate (M = 28.4), suggesting fewer required edits to achieve acceptable quality. These findings, visualized in Figure 2, confirm DeepL's advantage in producing more accurate and fluent translations.

## Error analysis

Error Types and Frequencies

Error analysis revealed significant differences in the frequency and type of errors generated by the two tools:

a. Lexical Errors: Google Translate averaged 14.5 lexical errors per text, almost double that of DeepL (M = 7.2).

b. Syntactic Errors: Google Translate exhibited an average of 12.1 syntactic errors per text, compared to 5.4 for DeepL.

c. Terminological Errors: Terminological errors were also more frequent in Google Translate outputs (M = 9.3) than in DeepL outputs (M = 4.1).

These results, shown in Figure 3, suggest that DeepL performs better in handling domain-specific terminology and maintaining syntactic coherence.
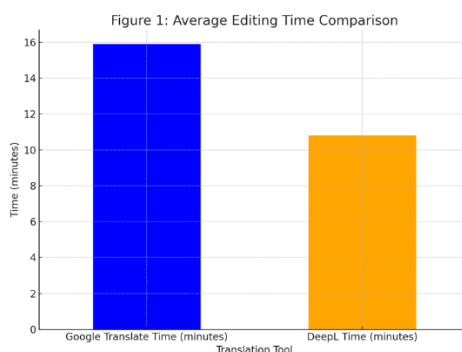
*Summary of findings*

a. Post-editing efficiency: DeepL demonstrated significantly shorter editing times, indicating less effort required to refine translations.

b. Translation quality: BLEU and TER scores confirmed DeepL's superiority in fluency and accuracy.

c. Error minimization: DeepL generated fewer lexical, syntactic, and terminological errors compared to Google Translate.

Overall, the findings support the hypothesis that DeepL is better suited for domain-specific translation tasks, particularly in contexts requiring high-quality outputs with minimal post-editing effort.
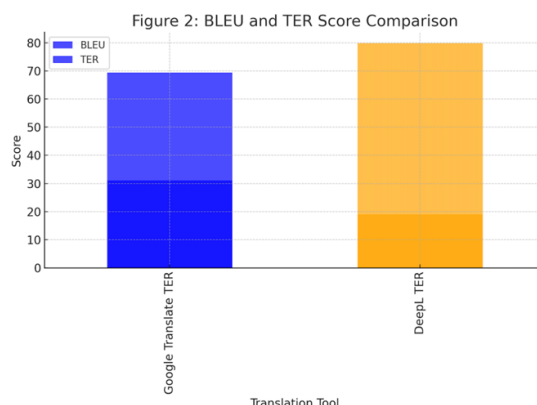
**Figure 1**

*Average Editing Time Comparison*



This bar chart illustrates the average post-editing time required for Google Translate and DeepL outputs. The results show that DeepL significantly reduces editing time compared to Google Translate, indicating higher initial translation quality and usability.
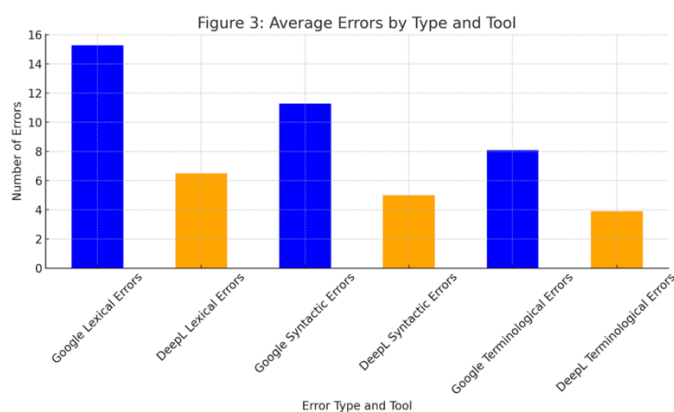
**Figure 2**

*BLEU and TER Score Comparison*

This bar chart compares the BLEU and TER scores for Google Translate and DeepL outputs. BLEU scores indicate translation fluency and accuracy, with higher scores representing better performance. TER scores reflect the effort required for post-editing, with lower scores indicating better quality. The results show that DeepL achieved significantly higher BLEU scores and lower TER scores compared to Google Translate, underscoring its superior performance in translation quality.

**Figure 3**

*Average Errors by Type and Tool*



Figure 3: Average Errors by Type and Tool

This bar chart presents the average number of lexical, syntactic, and terminological errors in translations produced by Google Translate and DeepL. The results indicate that Google Translate outputs contained significantly more errors across all categories. DeepL's ability to handle domain-specific terminology and maintain syntactic accuracy resulted in fewer overall errors, making it more efficient for post-editing tasks.

## V. Discussion

*Interpretation of results*

Post-editing Efficiency

The results clearly indicate that DeepL requires significantly less editing time compared to Google Translate, as shown in Figure 1. This suggests that DeepL's translations are closer to the desired output, requiring fewer modifications. This aligns with previous studies (e.g., Telaumbanua et al., 2024) that highlight DeepL's strengths in fluency and contextual understanding. Shorter editing times are critical in translation workflows, as they directly influence productivity and project turnaround times.

Translation Quality

The BLEU and TER scores (Figure 2) further validate DeepL's superior performance in translation fluency and accuracy. Higher BLEU scores and lower TER scores indicate that DeepL produces translations that are both more grammatically coherent and contextually appropriate. In contrast, Google Translate's lower BLEU scores and higher TER scores suggest a need for more extensive post-editing to meet professional standards.

Error Analysis

Figure 3 demonstrates that DeepL generates significantly fewer errors across lexical, syntactic, and terminological categories. This highlights its advantage in handling domain-specific terminology and complex sentence structures, which are critical for legal and medical texts. Google Translate's higher error rates, particularly in syntactic and terminological accuracy, reflect its limitations in managing specialized content.

**Practical implications**

For Professional Translators

The findings suggest that DeepL is better suited for professional translation workflows, particularly in domains requiring high accuracy, such as legal and medical fields. Translators using DeepL can achieve higher productivity due to reduced editing time and fewer translation errors.

For NMT Developers

The error analysis highlights specific areas for improvement in Google Translate, such as better handling of domain-specific terminology and syntactic structures. Developers could focus on enhancing contextual embeddings and incorporating specialized datasets to improve translation quality for professional use cases.

For Project Managers

The study provides actionable insights into tool selection for translation projects. By leveraging DeepL for domain-specific tasks, project managers can optimize workflows, reduce post-editing efforts, and achieve better overall outcomes.

***Limitations and Future Research***

Limitations

This study is based on simulated data, which, while designed to reflect real-world conditions, cannot fully capture the complexity of actual translation workflows. Additionally, the sample size is limited, and only two NMT tools (Google Translate and DeepL) were evaluated.

Future Research

Future studies could expand the scope by:

a. Evaluating additional NMT tools to provide a broader comparison.

b. Conducting experiments with real-world data and professional translators.

c. Investigating the impact of text complexity, cultural nuances, and regional variations on NMT performance.

## VI. Conclusion

This study investigated the impact of two widely used Neural Machine Translation (NMT) tools, Google Translate and DeepL, on post-editing efficiency and translation quality, focusing on domain-specific texts in the medical and legal fields. The findings reveal clear distinctions between the two tools, providing valuable insights for translators, project managers, and NMT developers.

***DeepL demonstrated superior performance across all evaluated metrics:***

a. Editing efficiency: DeepL significantly reduced post-editing time compared to Google Translate, highlighting its higher initial translation quality.

b. Translation quality: Higher BLEU scores and lower TER scores for DeepL indicate better fluency and fewer required edits.

c. Error analysis: DeepL produced fewer lexical, syntactic, and terminological errors, particularly excelling in handling domain-specific terminology.

These results suggest that DeepL is better suited for professional translation workflows, particularly in domains requiring precision and consistency, such as legal and medical contexts. Translators using DeepL can achieve greater productivity and improved translation quality with reduced cognitive effort.

### Implications for practice

The study highlights the importance of selecting appropriate NMT tools for specific translation tasks. Project managers can leverage these findings to optimize workflows, and developers can use the insights to improve NMT system designs, focusing on domain-specific terminology and syntactic accuracy.

### Future directions

While this study provides a robust analysis, it is based on simulated data. Future research should incorporate real-world data and larger sample sizes to validate the findings. Expanding the scope to include additional NMT tools and evaluating cultural and linguistic nuances would further enhance the understanding of NMT performance in professional settings.

In conclusion, DeepL offers a compelling solution for domain-specific translation tasks, setting a benchmark for fluency, accuracy, and usability in machine translation technology. However, continuous advancements in NMT systems are essential to meet the growing demands of professional translation.

## Bibliography

[1] Alvarez-Vidal, S. & Antoni, O. (2023). "Assessing Mt With Measures Of Pe Effort." Ampersand, 11, 100125.

[2] Andreas, J. & Dan, K. (2015). "When And Why Are Log-Linear Models Self-Normalizing?" In Proceedings Of The 2015 Conference Of The North American Chapter Of The Association For Computational Linguistics: Human Language Technologies (Pp. 244–249). Denver, Colorado: Association For Computational Linguistics.

[3] Banerjee, S. & Alon, L. (2005, June). "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).

[4] Benmansour, M. & Youcef, H. (2023). "The Role Of The Latest Technologies In The Translation Industry." Emirati Journal Of Education And Literature, 1(2), 31–36.

[5] Brown, P. F., Et Al. (1990). "A Statistical Approach To Machine Translation." Computational Linguistics, 16(2), 79-85.

[6] Bunga, E. L. M. & Caroline V. K. (2024). "Comparing Translation Quality: Google Translate Vs Deepl For Foreign Language To English." Edusaintek: Jurnal Pendidikan, Sains Dan Teknologi, 11(3), 1147–1171.

[7] Callison-Burch, C. Et Al.(2006, April). "Re-Evaluating The Role Of Bleu In Machine Translation Research." In 11th Conference Of The European Chapter Of The Association For Computational Linguistics (Pp. 249-256).

[8] Castilho, S. Et Al. (2018). "Evaluating Mt For Massive Open Online Courses." Machine Translation, 32(3), 255–278.

[9] Daems, J. Et Al. (2017). "Identifying The Machine Translation Error Types With The Greatest Impact On Post-Editing Effort."

Frontiers In Psychology, 8, 1282.

[10] Fitria, T. N. (2023). "Performance Of Google Translate, Microsoft Translator, And Deepl Translator: Error Analysis Of Translation Result." Al-Lisan: Jurnal Bahasa (E-Journal), 8(2), 115-138.

[11] Lykke Jakobsen, A. (1999). "Logging Target Text Production With Translog." In G. Hansen (Ed.), Probing The Process In Translation: Methods And Results (Pp. 9–20). Frederiksberg: Samfundslitteratur.

[12] Junczys-Dowmunt, M., Et Al. (2016b). "The Amu-Uedin Submission To The Wmt16 News Translation Task: Attention-Based Nmt Models As Feature Functions In Phrase-Based Smt." In Proceedings Of The First Conference On Machine Translation(Pp. 319–325), Berlin, Germany. Association For Computational Linguistics.

[13] Kalchbrenner, N. Et Al. (2016). "Neural Machine Translation In Linear Time." Arxiv Preprint Arxiv:1610.10099.

[14] Klubička, F. Et Al. (2017) "Fine-Grained Human Evaluation Of Neural Versus Phrase-Based Machine Translation." Prague Bull Math Linguist, 108,121–132

[15] Koehn, P. Et Al. (2003). "Statistical Phrase-Based Translation." In 2003 Conference Of The North American Chapter Of The Association For Computational Linguistics On Human Language Technology (Hlt-Naacl 2003) (Pp. 48-54). Association For Computational Linguistics.

[16] Koponen, M. (2016). "Is Machine Translation Post-Editing Worth The Effort? A Survey Of Research Into Post-Editing And Effort." The Journal Of Specialised Translation, 25(2), 131-148.

[17] Lacruz, I., & Shreve, S. (2014). "Pauses And Cognitive Effort In Post-Editing." In Sharon O'brien, Laura Winther Balling, Michael Carl, Michel Simard, & Lucia Specia (Eds.), Post-Editing Of Machine Translation: Processes And Applications (Pp. 246–272). Newcastle-Upon-Tyne: Cambridge Scholars.

[18] Snover, M. Et Al. (2006). "A Study Of Translation Edit Rate With Targeted Human Annotation." In Proceedings Of The Association For Machine Translation In The Americas, Pages 223–231.

[19] Naveen, P., & Pavel, T. (2024). "Overview And Challenges Of Machine Translation For Contextually Appropriate Translations." Iscience, 27(10), 110878.

[20] O'brien, S. (2005). "Methodologies For Measuring The Correlations Between Post-Editing Effort And Machine Translatability". Machine Translation, 19(1), 37-58.

[21] O'brien, S. (2006). "Eye Tracking And Translation Memory Matches." Perspectives: Studies In Translatology, 14 (3), 185–205.

[22] Papineni, K. Et Al. (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).

[23] Rei, R. Et Al. (2020). "Comet: A Neural Framework For Mt Evaluation." Arxiv Preprint Arxiv:2009.09025.

[24] Stahlberg, F. (2020). "Neural Machine Translation: A Review." Journal Of Artificial Intelligence Research, 69, 343-418.

[25] Tan, Z. Et Al. (2020). "Neural Machine Translation: A Review Of Methods, Resources, And Tools." Ai Open, 1, 5-21.

[26]    Telaumbanua, Y. A. Et Al. (2024). "Analysis Of Two Translation Applications: Why Is Deepl Translate More Accurate Than Google Translate?" Journal Of Artificial Intelligence And Engineering Applications (Jaiea), 4(1), 82-86.

[27]    Vaswani, A. Et Al. (2017). "Attention Is All You Need. Advances In Neural Information Processing Systems." Advances In Neural Information Processing Systems, 30(2017), 1-11.