ARTICLE

# Cross-Lingual Evidence-Based Strategies for Identifying Fabrications in Neural Translation Systems

**Phm Quc Huy**[1]

[1] Can Tho University, Department of Computer Science, Street, Ninh Kieu District, Can Tho, Vietnam.

## Abstract

This paper investigates cross-lingual evidence-based strategies for uncovering fabrications that emerge in neural translation systems, focusing on the phenomenon known as hallucination in machine-generated texts. Although neural architectures have achieved remarkable performance in many translation tasks, they remain prone to generating content that is factually ungrounded or entirely fabricated. Such fabrications compromise reliability, especially in domains where accuracy is mandatory. The proposed framework integrates linguistic alignment, bilingual term matching, and extrinsic verification checks across multiple language pairs. Emphasis is placed on constructing domain-specific corpora that underscore knowledge-intensive expressions, enabling robust identification of fabricated segments. Empirical analysis covers a spectrum of neural machine translation models, investigating how each architecture handles rare terms and nuanced syntactic constructs. The paper also develops a robust suite of metrics designed to quantify hallucination severity, leveraging lexical similarity measures and cross-entropy differentials. Results demonstrate that employing external knowledge sources and semantic aligners can reduce fabrication rates across a variety of languages, thereby enhancing translation integrity. The implications of this research extend to areas such as cross-border communications, international legal proceedings, and medical translations. By synthesizing empirical findings, this study offers a nuanced roadmap for future explorations in cross-lingual integrity verification, highlighting the evolving interplay between data-driven models and linguistic fidelity.

## 1 Introduction

Neural Machine Translation (NMT) architectures have redefined the scope of automated language conversion, benefiting from large-scale parallel corpora and sophisticated deep learning frameworks. Statistical models once held primacy, relying on phrase-based and word-based alignments that demanded explicit feature engineering. Newer models rely on self-attention mechanisms that capture long-range dependencies, yielding improved fluency in both high-resource and low-resource languages. Despite these successes, growing evidence indicates that even the most advanced systems exhibit vulnerabilities wherein words, phrases, or entire segments become fabricated. Researchers often employ the term hallucination for these manifestations, but the stakes are elevated when disseminating misinformation, as readers and end-users expect accurate representations of the source content. Such hallucinatory outputs compromise user trust and raise broader ethical implications, ranging from misplaced reliance in critical sectors to potential legal ramifications when official documents are mistranslated.

A central difficulty arises from the neural network's reliance on learned representations that might

| Architecture | Core Mechanism | Advantages | Limitations |
|---|---|---|---|
| RNN-based NMT | Recurrent connections | Captures sequential dependencies | Struggles with long-range dependencies |
| CNN-based NMT | Convolutions | Parallelizable, faster training | Limited context window |
| Transformer | Self-attention | Handles long-range dependencies well | Requires extensive training data |
| Hybrid Models | Combination of architectures | Leverages strengths of multiple methods | Increased complexity |

**Table 1.** Comparison of different Neural Machine Translation (NMT) architectures.

not fully capture the nuances of context. When encountering rarely seen terminology, archaic expressions, or insufficient training data, the model can resort to invention [1], [2]. In cross-lingual contexts, the phenomenon becomes even more opaque, since there may be fewer bilingual experts able to detect inaccuracies. Moreover, with the proliferation of commercial translation services, detecting these fabrications has become a global priority. Significant research and industry efforts have emerged to mitigate inaccuracies, but the complexity of language phenomena and the global diversity of linguistic structures continue to challenge model reliability. Traditional evaluation metrics such as BLEU, METEOR, or TER focus on surface-level overlaps with reference texts, potentially overlooking deeper semantic inconsistencies [3].

Identifying fabrications in neural translations remains an essential endeavor if high-stakes fields like healthcare, legal documentation, and scientific publishing are to employ machine translation reliably. Subtle errors, such as altered drug dosages or changes in chemical compound names, can produce catastrophic outcomes. Even less critical mistranslations can hamper cross-cultural communications, limit the reach of vital information, and erode user confidence. This prompts an urgent need for a framework that emphasizes evidence-based validation strategies. Instead of relying solely on standard reference-based evaluation, researchers can harness external knowledge repositories, curated bilingual dictionaries, or domain-specific lexicons to confirm the veracity of translated segments [4], [5].

Multiple strategies have been proposed to address

these issues. Some rely on decoder modifications that encourage faithful generation by augmenting the training dataset with adversarial examples [6]. Others employ retrieval-based paradigms, wherein the translation model consults external data sources to validate or refute generated hypotheses. Yet, these approaches may require specialized infrastructure or extensive domain-specific knowledge. A modular, cross-lingual framework that harmonizes evidence-based validation with the model's inherent capacity for contextual inference offers a more robust pathway. In this study, the goal is to propose and systematically evaluate such an integrated approach.

The subsequent sections are structured to explore the multifaceted nature of cross-lingual hallucinations, elucidate evidence-based validation mechanisms, and propose quantifiable metrics for fabrication detection. This investigation covers multiple neural architectures, including standard sequence-to-sequence, Transformer-based, and convolutional variants, comparing how each is susceptible to generating falsified data under resource constraints. Empirical findings highlight the capacity of evidence-augmented methods to substantially reduce the frequency of these errors. Ultimately, this research provides a foundation for better understanding and mitigating neural hallucinations, offering a roadmap for future explorations that blend computational linguistic strategies with domain-specific knowledge integration.

## 2 Overview of Neural Translation Fabrications

Hallucination in neural translation emerges when models generate linguistically plausible sequences

| Type | Description | Causes | Impact |
|---|---|---|---|
| Intrinsic Hallucination | Fluent but incorrect output | Training data biases | Misleading translations |
| Extrinsic Hallucination | Unrelated content generation | Lack of contextual grounding | Loss of meaning |
| Omission Errors | Missing crucial details | Model underconfidence | Incomplete information |
| Insertion Errors | Addition of non-existent elements | Overgeneralization | Distorted message |

**Table 2.** Categorization of hallucinations in Neural Machine Translation (NMT).

| Metric | Type | Strengths | Weaknesses |
|---|---|---|---|
| BLEU | Reference-based | Measures n-gram overlap | Ignores semantic correctness |
| METEOR | Reference-based | Considers synonymy | Computationally expensive |
| TER | Edit-distance based | Evaluates effort needed for correction | Doesn't capture fluency well |
| BERTScore | Embedding-based | Captures semantic similarity | Requires deep models |

**Table 3.** Common evaluation metrics for detecting hallucinations in machine translation.

ungrounded in the source text. These anomalies can take various forms, from isolated word substitutions to entire passages that deviate from the intended meaning. Scholars have traced this behavior to factors including inherent network overconfidence, insufficient domain coverage in training sets, and an over-reliance on statistical patterns rather than genuine semantic alignment. In extreme cases, fabrications can appear so coherent that they escape rudimentary detection methods, leading to a more insidious problem of misinformation. Recognizing the magnitude of risk, researchers seek theoretical underpinnings that explain why neural translation systems are prone to this defect.

Empirical examinations of hallucinations have revealed that the phenomenon often arises in under-resourced language pairs where parallel corpora are minimal or imbalanced. For instance, widely spoken languages such as English, Spanish, or French typically have robust training data, which constrains the scope of hallucination to rare domains or obscure expressions. In contrast, languages with fewer digital resources see more pronounced inaccuracies,

because the model attempts to interpolate from patterns learned in dissimilar language contexts. Additionally, code-switching scenarios can exacerbate the problem, as the presence of multiple languages within a single sentence complicates the distributional cues upon which the translation engine depends.

The architecture of modern NMT systems contributes to how such fabrications manifest. Early sequence-to-sequence models employed recurrent neural networks (RNNs) with attention, enabling the decoder to attend selectively to relevant regions of the source text. Although this mechanism improved translation accuracy compared to older phrase-based models, it does not entirely mitigate the risk of fabrications. Transformers, relying on multi-headed self-attention, have further amplified gains in translation quality and speed. Nonetheless, the global attention mechanism can still spread errors over different layers, thereby embedding certain inaccuracies deep within the network. Moreover, subword tokenization schemes, while beneficial for handling out-of-vocabulary terms, can occasionally amplify confusion around morphological variants,

| Cause | Description | Impact | Example |
|---|---|---|---|
| Overconfidence | Model assigns high probability to incorrect outputs | Leads to fluent but false translations | Fabricated medical terms |
| Low-Resource Training | Insufficient parallel data for certain languages | Higher hallucination rates in rare languages | Errors in indigenous language translation |
| Code-Switching | Mixing of languages in input text | Confuses tokenization and alignment | Incorrect handling of multilingual phrases |
| Subword Tokenization | Breakdown of words into smaller units | Can generate plausible but incorrect terms | Incorrect morpheme combinations |

Table 4. Key factors contributing to hallucinations in neural translation models.

leading to invented forms that appear plausible on the surface.

Large-scale pretraining also contributes to this dynamic. NMT systems increasingly benefit from pretrained language models like BERT, RoBERTa, or GPT variants. Although these models offer a powerful linguistic prior, they are trained predominantly on monolingual corpora. The cross-lingual adaptation often relies on additional fine-tuning or bridging techniques, but certain lexical or syntactic irregularities remain unaddressed. When the model encounters novel terms, domain-specific jargon, or idiomatic expressions, it may revert to generating approximate translations that risk deviating from the factual source content. Subtle distortions can accumulate, especially if the system's decoding mechanism selects high-probability tokens that maintain syntactic fluency without guaranteeing semantic fidelity.

The variability of hallucinations underscores the need to analyze them from multiple perspectives. Lexical deviations may occur through synonyms or near-synonyms that misrepresent domain-specific terms. Syntactic fabrications might reorder clauses in a way that alters the original meaning. There are also rhetorical-level hallucinations in which entire conceptual elements appear, diverging significantly from the source intent. Researchers must therefore develop frameworks that can detect anomalies at each level. Advances in interpretability shed light on hidden representations, showing how attention weights, activation patterns, and intermediate embeddings correlate with the eventual output.

Yet even interpretability methods sometimes fail to definitively isolate the root cause of hallucinations, given the complexity of deep architectures.

Achieving a comprehensive understanding of the neural underpinnings of fabrication informs the development of robust countermeasures. Evidence-based strategies, which rely on external verifiers, linguistic constraints, or cross-lingual alignments, represent a significant leap toward mitigating these errors. By systematically reviewing the current landscape, it becomes apparent that a unified approach that combines multiple signals—ranging from domain-specific dictionaries to alignment heuristics—can offer stronger protection against covert translation errors. The following sections propose methods that incorporate bilingual lexicon checks, semantic alignment tools, and knowledge graphs, all integrated within a pipeline designed to ensure factual accuracy. Such measures assume heightened importance in high-stakes applications, where even a minor distortion can have serious consequences.

## 3 Cross-Lingual Evidence-Based Approaches

Recent research underscores the utility of coupling internal model representations with external evidence sources to identify and correct hallucinatory output. One approach leverages bilingual dictionaries or terminological databases that list valid word correspondences for specialized domains. When the translation engine produces a segment containing terms absent in these references, a

| Fabrication Type | Definition | Example | Potential Consequence |
|---|---|---|---|
| Lexical Hallucination | Incorrect word substitution | "Aspirin" translated as "Ibuprofen" | Misinterpretation in healthcare settings |
| Syntactic Hallucination | Reordered clauses affecting meaning | "The patient was given medicine" → "The medicine was given to a patient" | Ambiguity in legal texts |
| Rhetorical Hallucination | Unfounded conceptual additions | Adding a reason for an event not present in source text | Misinformation in news translation |
| Structural Hallucination | Entirely new sentence formation | Inserting a fabricated summary in a translated document | Distortion of original content meaning |

**Table 5.** Classification of hallucinations based on linguistic structure.

| Strategy | Method | Advantages | Challenges |
|---|---|---|---|
| Evidence-based Validation | Uses external knowledge bases | Improves factual accuracy | Requires structured data sources |
| Bilingual Lexicon Checks | Cross-references translation with curated dictionaries | Ensures term consistency | Limited coverage for rare languages |
| Semantic Alignment Tools | Measures contextual coherence across languages | Reduces meaning distortion | Computationally expensive |
| Knowledge Graph Integration | Links translations to verified entity relationships | Helps prevent misinformation | Needs extensive preprocessing |

**Table 6.** Methods to detect and mitigate hallucinations in neural translation models.

flag is raised, prompting either a revision or an alert for human post-editing. Although seemingly straightforward, dictionary-based checks must balance comprehensiveness with the risk of false positives. Overly strict dictionaries may reject legitimate neologisms or updated nomenclature, thus stifling the evolution of language usage.

A more nuanced method employs cross-lingual semantic alignment. This technique projects source and target sentences into shared embedding spaces, ensuring that words or phrases sharing similar meanings cluster together. Tools such as multilingual BERT or XLM-R excel in creating embedding representations that capture semantic equivalences across languages. By examining alignment patterns, the system can detect segments that drift away from expected distributions, highlighting potential fabrications. This alignment-based detection can be further refined by introducing constraints derived from domain knowledge. For example, in a medical context, a term referencing a specific procedure in the source should map closely to the recognized target-language equivalent. If the alignment vector deviates significantly, a deeper investigation is triggered.

In parallel, knowledge graphs offer a structured repository of cross-lingual links between entities, concepts, and relationships. These graphs, which might contain factual data curated from encyclopedic sources, can operate as a robust verification layer. After the neural translator generates a candidate segment, the system checks whether the entities match known cross-lingual links in the graph. In cases where an entity is misrepresented or replaced with an unrelated item, the evidence-based pipeline flags that portion. Such checks are indispensable in specialized areas like legal or scientific translations, where an incorrect entity reference can transform the meaning of an entire passage. Coupling knowledge graphs with contextual embeddings refines the verification process further, since the system not only identifies entity mismatches but also evaluates the plausibility of relationships in the translated text.

A crucial aspect of implementing evidence-based strategies involves real-time or near-real-time operations, especially in live translation systems. Incorporating alignment checks, dictionary lookups, or knowledge graph queries during inference demands computational efficiency. Caching frequently used terms and indexing knowledge graph subcomponents aligned to specific domains can mitigate latency concerns. Some frameworks adopt a hybrid approach, performing partial checks for common terms on-the-fly while deferring more complex validations for a post-processing phase. Such trade-offs balance accuracy with throughput requirements, accommodating large-scale deployments. Importantly, the synergy between evidence-based checks and neural inference should be configured to preserve the fluidity of generation, preventing the system from becoming overly rigid or reliant on a limited reference set.

Human-in-the-loop strategies can further refine cross-lingual evidence-based mechanisms. Expert translators can review flagged segments, accepting or rejecting suggested corrections, thereby creating iterative feedback loops. This feedback enhances the system's dictionary entries and alignment thresholds, progressively lowering the incidence of false positives. In specialized domains where language evolves rapidly—such as technology or bioscience—expert oversight ensures that any expansions in domain vocabulary are integrated swiftly. Machine translation developers are increasingly moving toward collaborative solutions, acknowledging that purely automated verification may not suffice for highly complex or sensitive texts. By designing user-friendly interfaces where domain experts can annotate or comment on flagged segments, the entire pipeline gains from continuous refinement grounded in real-world usage patterns.

## 4 Implementation and Experimental Setup

Empirical validation of the proposed cross-lingual evidence-based strategies demands careful experimental design, covering diverse language pairs and text domains. This section outlines the setup adopted in this study, detailing data collection, model architectures, and implementation specifics for dictionary checks, semantic alignment modules, and knowledge graph integrations. A core objective is to ensure that the methodology is transparent and reproducible, enabling comparisons with baseline systems that lack evidence-based mechanisms.

Data curation begins by selecting parallel corpora for four language pairs: English–Spanish, English–German, English–Chinese, and English–Swahili. These pairs capture different degrees of resource availability, from widely studied European languages to relatively under-resourced African languages. For each language pair, the corpora include both general-domain texts (such as Wikipedia articles) and specialized texts from scientific, legal, and medical domains. This partition allows the study to assess how well the proposed methods handle both commonly encountered phrases and highly domain-specific expressions. An additional monolingual corpus is included for fine-tuning the semantic alignment module, ensuring that context vectors capture broader usage patterns beyond direct translation pairs.

Model selection encompasses three principal NMT architectures: a Transformer-based model, a recurrent sequence-to-sequence model with attention, and a convolutional sequence-to-sequence model. Each is trained using standard hyperparameters derived from published benchmarks, with minor modifications to integrate external evidence modules. The baseline configuration operates without any external evidence, producing translations in a purely data-driven manner. The experimental configuration incorporates dictionary-based validation, semantic alignment checks, and knowledge graph lookups. For the dictionary-based component, specialized lists derived from domain-specific glossaries supplement a general-purpose bilingual dictionary. A threshold-based mechanism determines when to flag

a translation for containing lexical items that do not appear in the external references.

The semantic alignment module employs a multilingual pretrained model, such as XLM-R, fine-tuned on parallel sentence pairs to maximize cross-lingual consistency. Sentences are projected into a shared embedding space, and an alignment score is computed by comparing average or max-pooled embeddings across segments. If the alignment score falls below a certain threshold, indicating semantic drift, the translation is flagged for further examination. Knowledge graph integration relies on linking extracted named entities from both source and target to a cross-lingual knowledge base that contains entity IDs and interlingual links. The system checks whether the entity mapping is valid and whether the relationships between entities align with those stored in the graph. This step is carried out in a pipeline fashion, after the translator generates a preliminary output and the alignment module signals an acceptable global match.

To ensure that the evaluation is robust, the experimental suite features both automatic metrics and human judgments. For the automatic metrics, BLEU, METEOR, and TER are used to quantify surface-level accuracy relative to reference translations. A specialized hallucination metric is introduced to capture instances of fabricated text, assigning a penalty score when flagged translations deviate significantly from the source meaning. Human evaluators, each with domain expertise in the relevant subject matter, conduct blind assessments on subsets of flagged segments. They categorize issues based on severity (minor lexical deviation versus major factual error) and domain risk level (low-stakes general text versus high-stakes medical passage). This combination of automated and manual evaluation offers a clearer picture of how well the evidence-based strategies perform, compared with baselines that rely solely on large-scale neural inference.

Implementation details are facilitated using a modular framework that divides the translation pipeline into distinct components. Python-based scripts orchestrate data preprocessing and post-processing, while specialized modules written in C++ or CUDA handle tokenization, alignment calculations, and knowledge graph queries for speed. The integration points between modules are standardized, allowing multiple NMT backends to be interfaced with the same evidence-based layer. This design further promotes extensibility, as new dictionaries, additional domain corpora, or updated knowledge graphs can be incorporated without rewriting the core pipeline. Logs capturing flagged translations, alignment scores, and final acceptance or rejection decisions are stored in a structured database, enabling in-depth analysis and reproducibility.

## 5 Evaluation Metrics and Analysis

A robust evaluation framework is critical for measuring the effectiveness of evidence-based strategies. Traditional reference-based metrics such as BLEU, METEOR, and TER quantify the overlap between the machine-generated translation and one or more human-generated references. While these metrics can highlight gross deviations in lexical choice or word order, they offer limited insights into deeper semantic or factual consistency. The phenomenon of hallucination requires metrics that capture the alignment of named entities, domain-specific vocabulary, and conceptual fidelity. Consequently, this study introduces a dedicated hallucination index (HI), which weights detected fabrications based on their severity and potential impact.

The HI begins by categorizing errors flagged by the evidence-based pipeline. Each flagged instance is mapped to one of four categories: (1) minor lexical mismatch, (2) moderate lexical mismatch with partial semantic drift, (3) major factual error, and (4) complete hallucination. Category weights reflect the potential harm or confusion resulting from each error. For example, a minor lexical mismatch, such as an incorrect preposition, carries a low weight. A major factual error, such as translating a drug name incorrectly, carries a high weight. Summing weighted errors and normalizing by the number of words in the translation yields the HI. Lower HI values indicate translations that remain faithful to the source, whereas higher HI values point to significant hallucination. This formulation allows for a granular assessment of how each component of the evidence-based system contributes to reducing different classes of errors.

Beyond numerical metrics, a qualitative analysis illuminates the specific conditions under which fabrications arise and how external validation curtails them. Detailed error typology tables reveal that dictionary-based checks excel at catching specialized vocabulary issues, while semantic alignment modules are adept at spotting more subtle divergences. Knowledge graph validations, on the other hand, primarily target named entities and relational integrity.

In some samples, the pipeline flags an error that turns out to be a legitimate expression unknown to the dictionary but semantically aligned with the source context. These instances lead to false positives, underscoring the balancing act required when applying rigid evidence rules to evolving language usage. Human evaluators help refine these thresholds, suggesting that domain-specific expansions or dynamic updates to the external references can mitigate false alarms.

Statistical analysis of the experiment's results across the four language pairs reveals interesting trends. For high-resource language pairs, like English–Spanish and English–German, baseline hallucinations are less frequent, and the introduction of evidence-based checks yields modest but tangible improvements. For lower-resource settings, like English–Swahili, baseline hallucination rates are substantially higher. In these contexts, the pipeline leads to a more pronounced reduction in false content. Domain-specific texts, particularly those dealing with technical jargon, benefit from dictionary-based checks, which sharply reduce lexical fabrications. Conversely, knowledge graph integration proves invaluable in legal or medical texts abundant in named entities. The synergy of these components often yields the most significant performance gains, a finding consistent with the hypothesis that multiple evidence streams are necessary to capture the breadth of possible fabrication modes.

Further introspection into model-specific behaviors uncovers how different architectures respond to evidence augmentation. Transformer-based models, while strong in capturing long-distance dependencies, show occasional vulnerabilities in translating very rare terms under domain constraints. In these instances, dictionary validation curbs the model's tendency to invent plausible-sounding but incorrect terms. Recurrent models, which rely heavily on contextual cues from hidden states, gain more from alignment checks, suggesting that explicit cross-lingual embeddings help rectify their comparatively lower capacity for global context modeling. Meanwhile, convolutional models display a mix of behaviors, sometimes excelling in short sentence translations but faltering on longer sequences involving multiple clauses. Across all architectures, the introduction of evidence-based modules consistently lowers the HI, confirming the utility of combining internal neural inferences with external validation signals.

Human judgment remains the ultimate arbiter of translation fidelity. Qualitative feedback from experts highlights that while the pipeline avoids many glaring errors, it may still struggle with nuanced idiomatic expressions or newly coined terms absent from dictionaries and knowledge graphs. In addition, experts emphasize the desirability of interactive or incremental validation, enabling them to selectively override flagged segments when context justifies it. The concluding analyses underscore the pivotal role of robust external evidence in elevating translation reliability and provide direction for refining these strategies to accommodate evolving language patterns [7].

# 6 Discussion of Challenges and Future Directions

Despite the advancements reported, the implementation of cross-lingual evidence-based strategies for detecting fabrications in neural translation systems faces several unresolved challenges. First, the ongoing evolution of languages complicates dictionary-based and knowledge graph-based validations. Many languages grow through the adoption of loanwords, the creation of neologisms, or the reclamation of regional dialects that were previously under-documented. These constant updates require continuous curation of external references, or else the system might flag legitimate terms as hallucinations. Although crowdsourcing and community-driven initiatives have proven helpful in updating lexical repositories, the incorporation of these changes into production-grade pipelines remains a non-trivial undertaking [8].

A second challenge lies in the inherent domain-dependence of many verification methods. A dictionary or knowledge graph tailored for legal texts may have limited utility in medical or technical domains. Conversely, expanding a knowledge base to cover multiple domains can lead to an unwieldy repository that increases query times and introduces potential conflicts in term usage. This domain specificity underscores the need for modular designs that allow easy swapping or customization of external resources. Another approach might incorporate machine-learning techniques that learn to weigh the relevance of domain-specific references on a per-document or per-segment basis, thus mitigating performance degradations when switching contexts [5], [9].

Third, potential biases embedded in the training data

and the external evidence sources can undermine the reliability of these strategies [10]. If a dictionary or knowledge graph predominantly represents one cultural or regional perspective, the system may inadvertently propagate those biases, flagging terms from marginalized dialects as errors. Conversely, it may fail to detect certain types of fabrications that arise from underrepresented viewpoints [11], [12]. Addressing such biases requires not only technical adjustments but also broader sociolinguistic awareness in dataset selection [13], curation, and annotation processes. Ethics boards or cross-disciplinary committees may need to be consulted when deploying translation systems in sensitive or culturally diverse environments.

The computational overhead of evidence-based validation is another significant issue. Integrating dictionary lookups, alignment checks, and knowledge graph queries into a real-time inference pipeline can strain memory and processing resources, especially in edge computing scenarios or large-scale web services that handle millions of translation requests daily. Research into more efficient alignment algorithms, caching strategies, or approximate matching techniques may help alleviate these overheads. Likewise, hardware accelerations—such as GPUs specialized for matrix operations or domain-specific accelerators—can potentially be harnessed to streamline external evidence checks without compromising throughput.

In addition to these challenges, future directions include expanding the scope of evidence-based strategies to cover morphological and phonological dimensions of language. Many current methods focus on lexical and named-entity consistency but overlook subtler aspects of linguistic variation. Morphological mismatches can lead to inaccuracies in languages with complex inflectional systems. Extending cross-lingual embedding models to incorporate morphological features, or developing specialized subword dictionaries that account for inflectional patterns, may offer improvements in capturing nuanced fabrications. Researchers might also examine the intersection of prosody and semantics in spoken language translation systems, an area that remains under-explored in large-scale neural frameworks [14], [15].

Finally, collaborative pipelines that combine generative models with retrieval modules or symbolic reasoning engines hold promise. Retrieval modules can supply candidate phrases or factual statements from large corpora, while symbolic reasoning can ascertain logical consistency. By merging these paradigms, the next generation of machine translation systems may achieve a level of factual grounding far surpassing current capabilities. Although these developments demand significant engineering and conceptual innovations, they align with the emerging vision of machine intelligence that integrates symbolic and sub-symbolic approaches to better handle the intricate tapestry of human language. These trajectories underscore that the effort to eradicate fabrications in neural translation is not solely a matter of incremental improvement but an ongoing reimagining of how data-driven and knowledge-based processes can be fused [13].

## 7 Conclusion

The analysis presented here demonstrates that cross-lingual evidence-based strategies offer a substantial layer of protection against fabrications in neural translation systems [16], [17]. Through the integration of dictionaries, semantic alignment tools, and knowledge graph validation, the incidence of hallucinatory outputs is reduced across multiple language pairs and domains. These methods mitigate both minor lexical mismatches and major factual distortions, underscoring the versatility of hybrid approaches that merge data-driven inference with curated external references. Results from automated metrics and human evaluations converge on the conclusion that such evidence-based pipelines substantially enhance translation fidelity, especially for specialized or high-stakes applications. Nonetheless, this study has highlighted the numerous challenges that must be addressed before fully robust deployment becomes possible. Continuous updates to external resources, careful handling of domain dependencies, and vigilance regarding data biases remain critical. Future work will likely explore the integration of morphological and phonological features, further optimization of computational overhead, and deeper fusions of retrieval-based and symbolic reasoning paradigms. These directions promise a new generation of translation systems that align more closely with the complex and rapidly evolving tapestry of human language, bringing us closer to a scenario in which mistranslations born from hallucination become rare exceptions rather than routine risks.

## Conflicts of Interest

## Acknowledgement

## References

[1] L. Ranaldi and G. Pucci, "Does the english matter? elicit cross-lingual abilities of large language models," in *Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)*, 2023, pp. 173–183.

[2] Y. Qiu, Y. Ziser, A. Korhonen, E. M. Ponti, and S. B. Cohen, "Detecting and mitigating hallucinations in multilingual summarisation," *arXiv preprint arXiv:2305.13632*, 2023.

[3] R. Mehta, A. Hoblitzell, J. O'keefe, H. Jang, and V. Varma, "Halu-nlp at semeval-2024 task 6: Metacheckgpt-a multi-task hallucination detection using llm uncertainty and meta-models," in *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, 2024, pp. 342–348.

[4] W. Zhu, Y. Lv, Q. Dong, *et al.*, "Extrapolating large language models to non-english by aligning languages," *arXiv preprint arXiv:2308.04948*, 2023.

[5] J. Wang, Y. Liang, F. Meng, *et al.*, "Zero-shot cross-lingual summarization via large language models," *arXiv preprint arXiv:2302.14229*, 2023.

[6] S. V. Bhaskaran, "Tracing coarse-grained and fine-grained data lineage in data lakes: Automated capture, modeling, storage, and visualization," *International Journal of Applied Machine Learning and Computational Intelligence*, vol. 11, no. 12, pp. 56–77, 2021.

[7] H. Ye, T. Liu, A. Zhang, W. Hua, and W. Jia, "Cognitive mirage: A review of hallucinations in large language models," *arXiv preprint arXiv:2309.06794*, 2023.

[8] R. Mehta, A. Hoblitzell, J. O'Keefe, H. Jang, and V. Varma, "Metacheckgpt–a multi-task hallucination detection using llm uncertainty and meta-models," *arXiv preprint arXiv:2404.06948*, 2024.

[9] H. Huang, T. Tang, D. Zhang, *et al.*, "Not all languages are created equal in llms: Improving multilingual capability by cross-lingual-thought prompting," *arXiv preprint arXiv:2305.07004*, 2023.

[10] S. V. Bhaskaran, "Enterprise data architectures into a unified and secure platform: Strategies for redundancy mitigation and optimized access governance," *International Journal of Advanced Cybersecurity Systems, Technologies, and Applications*, vol. 3, no. 10, pp. 1–15, 2019.

[11] N. M. Guerreiro, D. M. Alves, J. Waldendorf, *et al.*, "Hallucinations in large multilingual translation models," *Transactions of the Association for Computational Linguistics*, vol. 11, pp. 1500–1517, 2023.

[12] A. Bruno, P. L. Mazzeo, A. Chetouani, M. Tliba, and M. A. Kerkouri, "Insights into classifying and mitigating llms' hallucinations," *arXiv preprint arXiv:2311.08117*, 2023.

[13] S. V. Bhaskaran, "Integrating data quality services (dqs) in big data ecosystems: Challenges, best practices, and opportunities for decision-making," *Journal of Applied Big Data Analytics, Decision-Making, and Predictive Modelling Systems*, vol. 4, no. 11, pp. 1–12, 2020.

[14] C. M. Bishop and H. Bishop, *Deep learning: Foundations and concepts*. Springer Nature, 2023.

[15] K. Saitoh, *Deep learning from the basics: Python and deep learning: Theory and implementation*. Packt Publishing Ltd, 2021.

[16] M. Mahrishi, K. K. Hiran, G. Meena, and P. Sharma, *Machine learning and deep learning in real-time applications*. IGI global, 2020.

[17] P. Grohs and G. Kutyniok, *Mathematical aspects of deep learning*. Cambridge University Press, 2022.