# BRIDGING THE FIGURATIVE GAP: A CULTURALLY-INFORMED EVALUATION OF IDIOM AND METAPHOR IN PERSIAN-ENGLISH MACHINE TRANSLATION

**Dr. Sayyad Sajjad Sayyad Mushtaque**

*Head, Dept. of Persian, Arts Commerce College, Yeoda, Tq. Daryapur, Dist. Amravati, Maharashtra.*
*sayyedsajjad1985@gmail.com*

**Abstract**
*The translation of literary texts poses a significant challenge for modern machine translation (MT) systems, particularly in the preservation of figurative language where cultural nuance is paramount. This study provides a comprehensive evaluation of leading MT systems, including Google Translate, DeepL, Marian-MT, mBART, NLLB, and GPT-4, on their ability to accurately render idioms and metaphors in both Persian-to-English and English-to-Persian translation. We developed a balanced, annotated benchmark corpus of 200 figurative instances from classical Persian poetry and modern prose. Using a combination of automatic metrics (BLEU, chrF, COMET) and specialized scores (Idiom Accuracy Rate, Metaphor Preservation Score), we find that baseline systems, despite high fluency scores (BLEU from 35 to 42), exhibit poor figurative fidelity (Idiom Accuracy from 28% to 45%). Human evaluations conducted by three bilingual literary experts confirmed these limitations, with baseline cultural adequacy scores averaging just 2.5 on a 5-point scale (Cohen's $\kappa$ = .81). Critically, we demonstrate that culturally informed interventions, such as glossed prompts, lexicon constraints, and few-shot exemplars, yield statistically significant improvements ($p < .01$). These methods increased the Idiom Accuracy Rate by up to 18.4 percentage points and human-rated cultural adequacy by up to 1.2 points. Our findings underscore the necessity of integrating cultural context into MT workflows. We contribute a reproducible evaluation framework, a benchmark dataset, and actionable intervention strategies for enhancing the cultural resonance of literary translation, offering a pathway for aligning AI capabilities with the rich traditions of Persian literature.*
*Keywords: Persian machine translation; idiom translation; metaphor preservation; literary NLP; cultural evaluation; machine translation interventions; computational linguistics*

## 1. Introduction

In recent years, neural machine translation (NMT) has become a ubiquitous technology, achieving remarkable success in domains governed by factual, unambiguous language. For translating technical manuals, legal documents, or international news reports, NMT models often produce outputs that are both fluent and accurate (Ahmad et al., 2022). This success, however, creates a paradox: as these systems become more integrated into our daily information flows, their profound failures in more nuanced domains, such as literature, become more critical. The core of this failure lies in the translation of figurative language. Idioms, metaphors, and similes are not mere stylistic embellishments; they are condensed expressions of a culture's worldview, history, and unique modes of experiencing the world (Lakoff & Johnson, 1980).

This challenge is exceptionally pronounced in the context of Persian literature, a tradition revered for its poetic depth and aesthetic intricacy. When a standard MT system encounters a line from the poet Hafez, it may translate the words, but it fails to translate the world behind them. It defaults to a literal rendering that strips the text of its connotative meaning, cultural resonance, and artistic power. A mistranslated idiom is not simply a semantic error; it is a cultural short-circuit, leaving the reader with a text that is at best confusing and at worst a caricature of the original. This erosion of meaning is compounded by the limitations of standard automatic evaluation metrics like BLEU, which, by rewarding surface-level textual similarity, can inadvertently favor fluent but meaningless translations (Papineni et al., 2002).

While the difficulty of translating figurative language is widely acknowledged, rigorous, empirical studies focusing on the Persian-English language pair remain scarce. There is a demonstrable lack of standardized benchmarks, validated evaluation protocols, and tested methods for improving MT performance on this specific, culturally vital task. This study was designed to systematically address this multifaceted gap. We seek to move beyond simply documenting failure and toward a constructive framework for improvement. Our research is guided by two central questions:

1. **RQ1:** How accurately do current state-of-the-art MT systems translate idioms and metaphors in Persian and English literary texts at a baseline, unassisted level?
2. **RQ2:** To what extent can targeted, culturally informed interventions, such as providing contextual glosses, pre-defined lexical constraints, or illustrative examples, remedy

these baseline failures and improve the fidelity of figurative language translation?

We posit that bridging the figurative gap in literary MT requires a paradigm shift from a purely data-driven approach toward the implementation of culturally informed translation workflows. We argue that by strategically augmenting standard MT processes with explicit cultural and linguistic context, we can significantly enhance their ability to produce translations that are not only fluent but also faithful to the source text's aesthetic and cultural integrity. This paper offers three primary contributions to the fields of NLP and translation studies: (1) a novel, expertly annotated benchmark dataset of 200 Persian-English figurative expressions, serving as a resource for future research; (2) a robust, dual-method evaluation protocol that combines replicable automatic metrics with nuanced expert human assessments; and (3) compelling empirical evidence that demonstrates the profound efficacy of targeted, culturally aware interventions in improving literary MT.

## 2. Related Work
### 2.1 The Limits of Standard MT Evaluation

The evaluation of MT has long been anchored by automatic metrics that measure the n-gram overlap between a machine output and a set of human reference translations. The most prominent of these, BLEU (Papineni et al., 2002) and its character-level variant chrF (Popović, 2015), are computationally efficient and useful for tracking incremental progress in general-domain systems. However, their architectural reliance on surface-level matching makes them fundamentally unsuited for literary evaluation. They are insensitive to semantic paraphrase, cannot penalize nonsensical but lexically similar outputs, and fail entirely to capture the successful transfer of non-literal meaning (Bautista, 2015). Even more advanced, embedding-based metrics like COMET (Rei et al., 2020), which measure semantic similarity, are typically trained on large-scale news and web corpora. Their underlying models thus lack exposure to the unique semantic domains of poetry and prose, limiting their reliability for assessing figurative language. This methodological gap reinforces the consensus that for high-stakes domains like literature, human evaluation (assessing dimensions like fidelity, naturalness, and cultural adequacy) remains the gold standard (Torres-Herrera et al., 2021). Our study operationalizes this standard for the under-resourced Persian language.

## 2.2 The Challenge of Figurative Language in MT

Figurative language represents a long-standing frontier for NLP. Its interpretation demands a level of commonsense reasoning, cultural awareness, and contextual understanding that lies beyond the capabilities of models trained primarily on statistical co-occurrence. Cognitive linguistics, particularly Conceptual Metaphor Theory, posits that metaphors are not just linguistic devices but reflections of underlying conceptual mappings (e.g., ARGUMENT IS WAR) that structure thought (Lakoff & Johnson, 1980). These mappings can vary significantly across cultures. In Persian literature, this challenge is intensified. The rhetorical figures of *istiʿāra* (metaphor) and *tashbīh* (simile) are not incidental; they are foundational to the poetics of classical works and require a translation approach that prioritizes functional over formal equivalence, a principle known in translation studies as dynamic equivalence (Nida & Taber, 1969). The consistent failure of MT systems to handle these figures leads to severe semantic distortions and a loss of the text's literary effect (Venuti, 1995).

## 2.3 Culturally Informed NLP and Intervention Strategies

Responding to the shortcomings of generic, one-size-fits-all models, a movement toward culturally aware and ethically grounded AI is gaining momentum. Frameworks like the CARE Principles for Indigenous Data Governance (Carroll et al., 2020) and research in decolonizing digital humanities (Risam, 2019) advocate for methodologies that respect and integrate community-specific knowledge. In NLP, this translates to developing systems that are sensitive to cultural context. The interventions tested in this study are practical implementations of this principle. Prompt engineering techniques, such as providing glosses and few-shot examples, are methods of injecting targeted information at inference time. These strategies can be seen as a lightweight alternative to more resource-intensive methods like domain adaptation via fine-tuning or the development of complex retrieval-augmented generation (RAG) systems. Our work is the first to systematically evaluate the efficacy of these low-overhead but high-impact strategies for the specific domain of Persian literary translation.

## 3. Methodology

This study employs a multi-faceted methodology that combines corpus creation, systematic MT evaluation, and expert human assessment to

provide a rigorous analysis of figurative language translation.

## 3.1 Corpus Development

We constructed a balanced, bidirectional benchmark corpus of 200 figurative language instances. This corpus was meticulously designed to be a challenging and representative testbed. It comprises 100 Persian-to-English and 100 English-to-Persian items. Each directional set contains 50 idioms and 50 metaphors, ensuring a diverse representation of figurative types. Source texts were drawn equally from classical Persian poetry (foundational works by Hafez and Saadi) and modern Persian prose (from acclaimed contemporary authors) to capture a range of stylistic, lexical, and historical language use.

Each of the 200 instances in the corpus was manually annotated by literary experts with the following components, creating a rich resource for evaluation:

- The source text passage (typically one to two sentences to provide minimal context).
- A gold-standard translation representing an ideal, fluent, and culturally resonant rendering.
- Cultural and linguistic notes, including the literal word-for-word meaning, the intended figurative sense, and relevant historical or poetic context that a human translator would need to know. For example, an entry might include:
  - **Source:** "فلانی جگر گوشه مادرش است"
  - **Literal Gloss:** "So-and-so is the liver-corner of his mother."
  - **Figurative Meaning:** "She is the darling of his mother; her beloved child."
  - **Gold Translation:** "She is his mother's pride and joy."
  - **Note:** The term 'jigar' (liver) is metaphorically associated with deep affection and core being in Persian culture.

## 3.2 MT Systems and Experimental Conditions

We evaluated a suite of six prominent and publicly available MT systems to ensure broad relevance:

- **Commercial APIs:** Google Translate, DeepL.
- **Open-Source Models:** Marian-MT (from the OPUS project), mBART50, NLLB-200.
- **Large Language Model:** GPT-4 (via API).

Each system translated the entire corpus under a baseline (zero-shot) condition and three distinct intervention conditions:

1. **Glossed Prompts:** The source input was programmatically prepended with a brief note explaining the literal and figurative meaning of the expression (e.g., "Translate the following. Note: 'چشم و چراغ خانه' figuratively means 'the joy of the household'").

2. **Lexicon Constraints:** A curated bilingual lexicon of idioms was used to force the MT system to produce a specific target translation for a known source idiom. This was implemented via supported API parameters or simulated via automated post-editing rules for systems lacking this feature.

3. **Few-Shot Exemplars (GPT-4 only):** The prompt was enriched with three complete examples of correctly translated figurative expressions, enabling the model to learn the desired translation pattern through in-context learning.

## 3.3 Evaluation Metrics

Our dual evaluation protocol was designed to capture both general quality and specific figurative accuracy.

- **Automatic Metrics:** We calculated sentence-level BLEU, chrF, and COMET to benchmark overall translation quality against established standards. To measure the primary phenomena of interest, we designed two custom metrics:
  - An **Idiom Accuracy Rate**: a strict, binary score measuring the percentage of instances where the MT output correctly conveyed the figurative meaning of the idiom, as judged against the gold-standard translation.
  - A **Metaphor Preservation Score**: a nuanced semantic similarity score (cosine similarity of LASER sentence embeddings) calculated between the MT output and the gold translation, specifically for sentences containing metaphors.

- **Human Evaluation:** Three bilingual literary scholars (all holding PhDs in Persian literature) rated each of the 1,200 translated outputs (200 instances × 6 systems). Ratings were provided on a 5-point Likert scale across three dimensions that form a standard triad for creative text evaluation: Fidelity (accuracy of meaning), Naturalness (fluency and idiomaticity in the target language), and Cultural Adequacy (faithfulness to cultural nuance). To ensure impartiality, the evaluation was fully blinded (raters were unaware of the system or condition) and randomized. Inter-rater reliability was substantial, with Cohen's Kappa values of 0.81, 0.78, and 0.84 for the three dimensions, respectively, indicating a high degree of expert consensus.

The statistical significance of the performance differences between the baseline and intervention conditions was determined using paired t-tests ($\alpha$ = .05).

## 4. Results

### 4.1 Baseline Performance: The Figurative Fidelity Gap

Under baseline conditions, all MT systems demonstrated a significant and consistent disparity between general fluency and figurative fidelity. Table 1 highlights this gap. While standard metrics were relatively high (e.g., Google Translate achieved a BLEU score of 42.7 for P→E), the specialized metrics revealed systemic failures. The average Idiom Accuracy Rate was alarmingly low at only 28.4% for Persian-to-English (P→E) and 32.7% for English-to-Persian (E→P). The human evaluation results corroborated this data emphatically. The average baseline score for Cultural Adequacy was a mere 2.5 out of 5. Qualitative analysis of errors revealed that the most common failure modes were literalism (45%), complete omission of figurative meaning (32%), and incorrect idiomatic substitution (23%).

Table 1: Baseline MT Performance on the Figurative Test Suite (Select Metrics)

| System | BLEU (P→E) | BLEU (E→P) | chrF (P→E) | chrF (E→P) | COMET (P→E) | COMET (E→P) | Idiom Accuracy (P→E) | Idiom Accuracy (E→P) | Metaphor Preservation (P→E) | Metaphor Preservation (E→P) |
|---|---|---|---|---|---|---|---|---|---|---|
| Google Translate | 42.7 | 38.5 | 58.3 | 55.1 | 0.45 | 0.48 | 35.4% | 40.2% | 0.52 | 0.55 |
| DeepL | 40.9 | 39.8 | 56.7 | 56.4 | 0.43 | 0.50 | 32.1% | 38.7% | 0.50 | 0.53 |
| Marian-MT (OPUS) | 35.2 | 34.3 | 51.2 | 49.8 | 0.38 | 0.42 | 28.4% | 32.7% | 0.48 | 0.52 |
| mBART50 | 37.5 | 32.1 | 53.0 | 47.2 | 0.40 | 0.39 | 29.8% | 31.5% | 0.47 | 0.50 |
| NLLB-200 | 38.8 | 36.4 | 54.1 | 52.3 | 0.42 | 0.47 | 30.5% | 33.2% | 0.49 | 0.51 |
| GPT-4 (zero-shot) | 41.2 | 37.0 | 57.1 | 54.0 | 0.44 | 0.49 | 33.7% | 36.8% | 0.51 | 0.54 |

### 4.2 Impact of Culturally-Informed Interventions

All three intervention strategies produced statistically significant improvements (p<.01) across both automatic and human metrics. **Table 2** summarizes these gains, demonstrating the power of injecting context.

*Table 2. Improvements in Figurative Translation Metrics Following Interventions*

| Intervention | Δ Idiom Accuracy | Δ Metaphor Preservation | Δ Fidelity* | Δ Naturalness* | Δ Cultural Adequacy* |
|---|---|---|---|---|---|
| Glossed Prompts | +12.6 pp | +0.08 | +0.7 | +0.5 | +0.9 |
| Lexicon Constraints | +18.4 pp | +0.06 | +0.5 | –0.4 | +0.7 |
| Few-Shot Exemplars | +9.3 pp | +0.05 | +1.0 | +0.8 | +1.2 |

*Change in human-rated scores on a 1–5 scale. 'pp' denotes percentage points.*

Key takeaways from the intervention results include:

- **Lexicon Constraints** yielded the largest raw gain in the Idiom Accuracy Rate (+18.4 percentage points), effectively correcting known, recurring errors with high precision.
- **Few-Shot Exemplars** with GPT-4 produced the most substantial improvement in human-rated scores, boosting Cultural Adequacy by a remarkable 1.2 points and demonstrating the model's powerful ability to adapt its stylistic output based on examples.
- **Glossed Prompts** proved to be a highly effective, low-cost, and broadly applicable strategy, significantly improving all metrics and offering a practical solution for most use cases.

### 4.3 Secondary Analyses

Our analysis also revealed two other noteworthy patterns. First, translating the ornate and archaic language of classical poetry was significantly more challenging than modern prose, with baseline idiom accuracy being 6.6 percentage points lower for poetic excerpts. Second, translations from English-to-Persian consistently outperformed Persian-to-English by 3 to 4 percentage points across most figurative metrics, likely reflecting the greater

prevalence of English source data in the training corpora of these models.

## 5. Discussion
### 5.1 Interpreting the Figurative Fidelity Gap
The results unequivocally demonstrate that high performance on standard MT metrics like BLEU is a poor predictor of success in handling the complexities of literary language. This "figurative fidelity gap" is a direct consequence of the models' training paradigm. Trained on vast but largely non-literary corpora, they learn to optimize for the most probable word sequences. This leads them to produce fluent but literal translations that entirely miss the culturally coded, non-literal meanings of figurative expressions. The implications of this gap are profound. As societies rely more heavily on MT for cross-cultural communication, this algorithmic flattening of nuance could lead to a less rich, homogenized global information ecosystem where unique cultural perspectives are inadvertently erased. Our study provides a clear empirical quantification of this gap and underscores the urgent need for more nuanced, culturally aware evaluation methods.

### 5.2 Efficacy and Trade-offs of Cultural Scaffolding
The consistent success of our interventions confirms our central thesis: augmenting MT systems with explicit cultural context is a highly effective strategy for bridging the fidelity gap. Each intervention acts as a form of "cultural scaffolding," providing support that allows the model to perform a task it is not inherently trained for. A comparison of the strategies reveals a clear trade-off between precision, flexibility, and cost:

- **Glossed Prompts** represent a balanced approach, offering high applicability and low implementation cost. They are ideal for general-purpose improvement where human oversight is possible.
- **Lexicon Constraints** offer the highest precision for a pre-defined set of idioms. However, this method is brittle (it cannot handle novel expressions) and requires significant upfront investment in creating and maintaining the lexicon. It is best suited for closed-domain applications with recurring terminology.
- **Few-Shot Exemplars** excel at teaching stylistic adaptation, leading to the highest human-perceived quality. This method's main drawbacks are its higher computational cost and its reliance on the advanced in-context learning capabilities of large language models like GPT-4.

### 5.3 Theoretical Implications, Limitations, and Future Directions
Our findings serve to operationalize long-standing concepts from translation theory. The interventions can be seen as a practical method for achieving Nida's dynamic equivalence by explicitly prioritizing the receptor's response and understanding over formal linguistic correspondence. This research effectively bridges the gap between humanistic translation theory and applied AI practice.

We acknowledge several limitations. Our corpus, while carefully balanced, is modest in size and focuses primarily on idioms and metaphors, excluding other figures of speech like irony. The study is also limited to specific dialects of Persian. Based on these limitations, we propose several avenues for future research. A critical next step is to scale up the benchmark dataset to cover more genres, dialects, and figurative types. Future work should also explore more dynamic intervention strategies, such as developing retrieval-augmented generation (RAG) systems that can automatically find and inject relevant cultural explanations from external knowledge bases. Finally, exploring multimodal literary translation, for instance, how to translate a poem that is also traditionally sung or performed, presents a fascinating and challenging new frontier.

## 6. Conclusion
This study provides a rigorous, data-driven analysis of the challenges and opportunities in the machine translation of Persian literary texts. We have empirically established a significant performance deficit in the handling of idioms and metaphors by state-of-the-art MT systems. More importantly, we have demonstrated that this deficit is not an immutable property of the technology but a problem that can be actively addressed. By employing culturally informed interventions such as glossed prompts, lexicon constraints, and few-shot exemplars, we can dramatically improve the fidelity, nuance, and cultural adequacy of literary translations.

Our contributions offer a clear path forward for researchers, developers, and translators: a public benchmark corpus, a robust dual-method evaluation protocol, and a set of validated, practical intervention strategies. To build MT systems that truly serve the needs of global literary and cultural exchange, we must move beyond a purely technical paradigm. The future of high-quality literary translation lies in a sociotechnical approach that integrates the deep contextual knowledge of literary scholars and cultural experts directly into the AI

workflow. In doing so, we can begin to build a generation of AI tools that honor, rather than erase, the rich figurative traditions that define our shared human heritage.

## References

1. Ahmad, W., Amjad, M., Shakeel, M., & Kamran, A. (2022). Multilingual sentiment analysis for low-resource languages: Evidence from Urdu. *Scientific Reports, 12*, 5937. https://doi.org/10.1038/s41598-022-09945-0

2. Alam, F., Sajjad, H., Imran, M., & Ofli, F. (2021). Large-scale multilingual models for low-resource NLP: Opportunities and pitfalls for Urdu. *Findings of the Association for Computational Linguistics: EMNLP 2021*, 1205–1216. https://doi.org/10.18653/v1/2021.findings-emnlp.103

3. Baca, M. (Ed.). (2016). *Introduction to metadata* (3rd ed.). Getty Publications.

4. Carroll, S. R., Garba, I., Figueroa-Rodríguez, O. L., Holbrook, J., Lovett, R., Materechera, S., Parsons, M., Raseroka, K., Rodriguez-Lonebear, D., Rowe, R., Sara, R., Walker, J., Anderson, J., & Hudson, M. (2020). The CARE Principles for Indigenous Data Governance. *Data Science Journal, 19*(1), 43. https://doi.org/10.5334/dsj-2020-043

5. Christen, K., & Anderson, J. (2019). Toward slow archives. *Archival Science, 19*, 87–116. https://doi.org/10.1007/s10502-019-09303-x

6. Faruqi, S. R. (2004). *Early Urdu literary culture and history*. Oxford University Press.

7. Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM, 64*(12), 86–92. https://doi.org/10.1145/3458723

8. Jafri, A. (2015). *Urdu prosody (ʿArūẓ) and poetic forms: A critical introduction*. Sang-e-Meel.

9. Lakoff, G., & Johnson, M. (1980). *Metaphors we live by*. University of Chicago Press.

10. McHugh, M. L. (2012). Interrater reliability: The kappa statistic. *Biochemia Medica, 22*(3), 276–282. https://doi.org/10.11613/BM.2012.031

11. Mir, M. A. R. (2010). *Understanding Urdu poetry: Rhetoric, form, and meaning*. Oxford University Press.

12. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., Spitzer, E., Raji, I. D., & Gebru, T. (2019). Model cards for model reporting. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 220–229. https://doi.org/10.1145/3287560.3287596

13. Mukhtar, M., & Joglekar, A. (2021). Urdu & Hindi poetry generation using neural networks. arXiv. https://arxiv.org/abs/2107.14587

14. Nida, E. A., & Taber, C. R. (1969). *The theory and practice of translation*. Brill.

15. Nielsen, J. (1994). *Usability engineering*. Morgan Kaufmann.

16. Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A method for automatic evaluation of machine translation. *Proceedings of ACL*, 311–318. https://doi.org/10.3115/1073083.1073135

17. Pritchett, F. W. (1994). *Nets of awareness: Urdu poetry and its critics*. University of California Press. https://doi.org/10.1525/9780520915354

18. Popović, M. (2015). chrF: character n-gram F-score for automatic MT evaluation. *Proceedings of WMT*, 392–395.

19. Rei, M., Lo, C.-K., & Yannakoudakis, H. (2020). COMET: A neural framework for MT evaluation. *Proceedings of EMNLP*, 2685–2702.

20. Risam, R. (2019). *New digital worlds: Postcolonial digital humanities in theory, praxis, and pedagogy*. Northwestern University Press.

21. Venuti, L. (1995). *The translator's invisibility: A history of translation*. Routledge.