

Evaluating Large Language Models on Urdu Idioms

Anonymous ACL submission

Abstract

We present a comprehensive evaluation of Urdu–English idiomatic translation, introducing a parallel benchmark encompassing Native and Roman Urdu across translation, paraphrasing, idiom span detection, and idiomatic back-translation. Eight prompting strategies—including literal, cultural, idiomatic, and few-shot prompts—are used to assess multiple open-source LLMs and NMT systems. Performance is evaluated using BLEU, ChrF, BERTScore, COMET, XCOMET, ROUGE, Levenshtein distance, and multilingual embedding cosine similarities (LASER, LaBSE, USE). Results indicate that LLMs outperform NMT systems in preserving idiomatic meaning, with cultural and idiomatic prompts yielding the highest semantic fidelity. Few-shot prompting further improves idiom handling. Native Urdu consistently achieves higher scores than Roman Urdu across all tasks, highlighting the influence of script on translation quality. This study provides the first multi-metric, cross-script benchmark for idiomatic Urdu–English translation, offering insights into model behavior, prompt sensitivity, and the challenges of Romanized input.

1 Introduction

Idiomatic expressions pose a long-standing challenge in machine translation due to their non-compositional semantics, cultural grounding, and contextual dependence. For example, the Urdu idiom اونٹ کے منہ میں زیرہ (“a cumin seed in a camel’s mouth”) corresponds to the English idiom “a drop in the ocean,” illustrating how figurative meaning is easily lost when systems rely on literal alignments. Accurate translation requires recognizing idiomatic spans, interpreting figurative intent, and generating culturally appropriate equivalents rather than surface-level mappings.

Urdu, the 10th most spoken language globally,¹

¹<https://en.wikipedia.org/wiki/Urdu>

presents additional complexity due to its dual writing systems: the formal Perso–Arabic script (Native Urdu) and the widely used but orthographically inconsistent Roman Urdu. The lack of standardized spelling and phonetic ambiguity in Roman Urdu complicate idiom preservation across scripts. Despite the prevalence of idioms in Urdu communication, there is limited research on how modern LLMs or NMT systems handle idiomatic content in either script, and no unified evaluation spanning translation, paraphrasing, span detection, and back-translation.

To address these gaps, we introduce the first comprehensive benchmark for Urdu–English idiom understanding and translation. Our benchmark extends prior idiom-focused work by covering multiple subtasks, including direct and literal translation, culturally grounded translation, paraphrasing, idiom-span detection, idiomatic back-translation, and few-shot translation with 5-shot and 10-shot demonstrations. This design enables analysis of idiom recognition, cross-script transfer, prompt sensitivity, and robustness under varied conditions.

We evaluate a suite of open-source LLMs and NMT models, including GPT-OSS-20B (Kumar et al., 2025), Qwen 3–8B (Yang et al., 2025), DeepSeek R1-Distill-Llama-8B (Guo et al., 2025), Gemma 3–1B (DeepMind and Team, 2025), Mistral Instruct-7B (Jiang et al., 2023), LLama-3.2-3B-Instruct (Grattafiori et al., 2024), NLLB models (Team et al., 2022), and Google Translate as a strong NMT baseline. We curated a total of 4,000 idiom–sentence pairs for evaluation, divided equally across scripts: 2,000 written in Native Urdu and 2,000 written in Roman Urdu. Each sentence contains at least one idiomatic expression and is paired with a verified English reference translation checked by native speakers.

We assess model outputs using a diverse metric suite capturing surface similarity (BLEU, ChrF), semantic similarity (BERTScore, COMET,

XCOMET), paraphrase consistency (ROUGE), and idiom-level variation (Levenshtein distance, unigram overlap, and multilingual embedding cosine scores). This multi-metric pipeline enables a holistic analysis of figurative meaning preservation across scripts, prompts, and architectures.

Our main contributions are:

- We release a unified benchmark for Urdu–English idiom understanding across translation, paraphrasing, idiom-span detection, and back-translation in Native and Roman Urdu.
- We provide a systematic analysis of prompt sensitivity in idiom translation for Urdu using eight prompting strategies.
- We present a cross-script comparison of idiomatic translation quality, revealing substantial performance gaps driven by orthographic representation.
- We introduce a comprehensive multi-metric evaluation combining BLEU, ChrF, BERTScore, COMET, XCOMET, ROUGE, Levenshtein distance, and embedding similarity.

2 Related Work

Idiomatic translation has been examined in several language pairs, with prior studies highlighting the difficulty of preserving figurative meaning in neural systems. Fadaee et al. (Fadaee et al., 2018) showed that data augmentation and semantic cues can improve idiom handling in NMT, while Liu et al. (Liu et al., 2023) explored cross-lingual idiom understanding using multilingual encoders. Recent work on Persian–English idiom translation demonstrates that retrieval-augmented prompting and semantic similarity models can enhance figurative meaning retention (Rezaeimanesh et al., 2025). However, these efforts remain limited to high- or mid-resource settings.

Research on Urdu NLP has focused primarily on low-resource challenges, including morphological complexity, script variation, and inconsistent orthography (Wahid et al., 2024; Butt et al., 2025). Despite the prevalence of idioms in Urdu communication, no prior work has investigated idiomatic translation across Native and Roman Urdu, and existing MT benchmarks do not include idiom-specific evaluation for Urdu.

3 Dataset

We construct a unified dataset of 4,000 Urdu–English idiom–sentence pairs, divided equally across writing systems. Each instance contains an Urdu sentence with at least one idiomatic expression, the extracted idiom, and a gold English translation verified by bilingual annotators. The dataset covers a broad range of figurative expressions spanning common proverbs, culture-dependent idioms, and semantically opaque phrases.

Table 1 shows an example idiom and its gold translation.

3.1 Native Urdu

The Native Urdu dataset contains 2,000 sentences written in the Perso–Arabic script, commonly used in formal communication, print media, and education. Idioms were collected from publicly available resources such as UrduPod101,² LingApp,³ and EnglishUms,⁴ supplemented with manually curated examples from native speakers. Each idiom was paired with an English equivalent and contextual sentence. Translations generated with an LLM were manually reviewed and corrected by a group of university students, fluent in Urdu and English, to ensure translation quality and idiomatic equivalence.

Idiom (Native Urdu): اونٹ کے منہ میں زیرہ
Idiom (Roman Urdu): Oont ke munh mein zeera
Idiom (English): A drop in the ocean
Sentence (Native Urdu): یہ تنخواہ تو میرے اخراجات کے مقابلے میں اونٹ کے منہ میں زیرہ ہے۔
Sentence (Roman Urdu): Ye tankhwah to mere kharchon ke muqablay mein oont ke munh mein zeera hai.
Gold Translation: This salary is just a drop in the ocean compared to my expenses.

Table 1: Example instance from the Native and Roman Urdu datasets.

3.2 Roman Urdu

Roman Urdu refers to the representation of Urdu using the Latin script, commonly used in informal

²<https://www.urdupod101.com/>

³<https://ling-app.com/>

⁴<https://englishums.com/>

digital communication such as texting, social media, and online forums. Unlike Native Urdu, Roman Urdu lacks standardized orthography, resulting in substantial spelling variation and increased lexical diversity, which pose challenges for translation and evaluation systems.

For this work, the Roman Urdu portion of the dataset was constructed to mirror the Native Urdu split. Starting from the 2,000 Native Urdu idioms and their accompanying sentences, we used an LLM (ChatGPT) to generate corresponding Roman Urdu versions. These automatically generated Roman Urdu idioms and sentences were then manually reviewed and corrected by native speakers proficient in both scripts to ensure linguistic accuracy, naturalness, and consistency with real-world Roman Urdu usage.

The final Roman Urdu dataset consists of 2,000 high-quality idiom–sentence pairs aligned with their Native Urdu and English counterparts, enabling controlled cross-script comparisons in downstream evaluation.

3.3 Dataset Statistics

Table 2 presents a quantitative comparison of the curated Native Urdu and Roman Urdu idiom–sentence datasets. Each split contains 2,000 idiom–sentence pairs, resulting in a total of 4,000 aligned instances with verified English translations.

Both datasets exhibit similar structural patterns: idioms average between six and seven words, while sentences typically range from nine to ten words. Roman Urdu shows a slightly larger vocabulary size (3,254 vs. 3,033 in Native Urdu) as well as a higher Type–Token Ratio (TTR), reflecting its greater orthographic variability and lack of standardized spelling conventions.

Dataset	# of Pairs	Avg. Idiom Length	Avg. Sentence Length	Vocab. Size	TTR
Native Urdu	2000	6.74	9.62	3033	0.213
Roman Urdu	2000	6.65	9.71	3254	0.232

Table 2: Linguistic statistics of the curated Urdu idiom translation datasets.

Figures 1 and 2 visualize the distributions of idiom and sentence lengths, respectively.

4 Experimental Design

This section outlines the prompting strategies, model configurations, and evaluation metrics used across all experiments, including idiom translation, idiom classification, English–Urdu back-

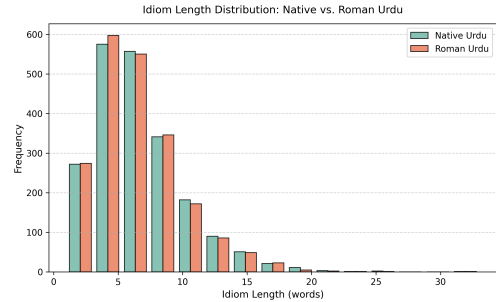


Figure 1: Distribution of idiom lengths in the Native and Roman Urdu datasets.

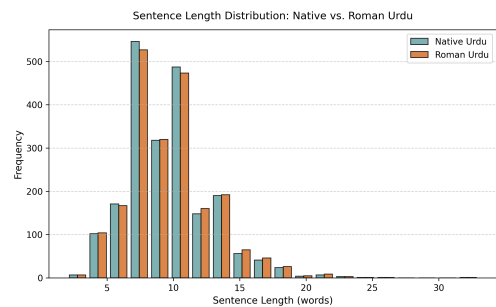


Figure 2: Distribution of Urdu sentence lengths in the Native and Roman Urdu datasets. Sentences typically range between seven and fourteen words, reflecting concise idiomatic usage in context.

translation, and Native vs Roman Urdu comparison.

4.1 Prompting Strategies

Table 3 summarizes the set of prompting strategies used in the experiments.

All large language model experiments were carried out using a chat-based interface in which the system role was explicitly set to “*You are an expert translator fluent in both Urdu and English.*”

4.2 Models

We evaluate a mix of open-source LLMs and Neural Machine Translation (NMT) models. The LLMs include GPT-OSS-20B (Kumar et al., 2025), Qwen 3-8B (Yang et al., 2025), DeepSeek R1-Distill-Llama-8B (Guo et al., 2025), Gemma 3-1B (DeepMind and Team, 2025), Mistral Instruct-7B (Jiang et al., 2023), and LLama-3.2-3B-Instruct (Grattafiori et al., 2024). For NMT baselines, we include nllb-200-3.3B, nllb-200-distilled-600M (Team et al., 2022), and Google Translate.

The models were executed with their default configurations and the output tokens were allowed to reach a maximum of 500. Some models produced explanations or reasoning along with the transla-

Prompt	Description
Literal Translation	Direct translation from Urdu to English with no additional explanations, reasoning, or text.
Cultural Translation	Translation that preserves cultural and idiomatic meaning while producing natural English expressions.
Paraphrase	Paraphrase the Urdu sentence into fluent English, avoiding word-for-word translation while keeping the meaning intact.
Idiomatic Translation	Identify the idiom in the Urdu sentence and provide an English idiomatic equivalent.
Back-Translation	Translate the Urdu idiom into an English idiom and then back to Urdu, returning only the back-translated Urdu idiom.
Idiom Span Detection	Detect and extract only the idiom span from the Urdu sentence.
5-Shot Translation	Translate Urdu to English using five in-context example pairs as references.
10-Shot Translation	Translate Urdu to English using ten in-context example pairs as references, providing more extensive few-shot guidance.

Table 3: Summary of prompting strategies used in the idiomatic translation experiments. Each prompt specifies the type of translation or task and any contextual guidance provided to the model.

tion, so we used regular expressions to extract the final translated sentences.

4.3 Evaluation Metrics

We employ a combination of lexical, semantic, and idiomatic evaluation metrics to comprehensively assess translation quality. BLEU (Papineni et al., 2002) and ChrF (Popović, 2015) measure surface-level n-gram overlap between the model output and reference translations. BERTScore (Zhang et al., 2020), COMET (Rei et al., 2020), and XCOMET (Guerreiro et al., 2024) capture semantic similarity and cross-lingual adequacy. ROUGE (Lin, 2004) is used for paraphrase evaluation, while Levenshtein distance and unigram overlap assess idiom-level variation and lexical preservation. Additionally, cosine similarity computed using multilingual embeddings (LASER, LaBSE, and USE-multilingual) quantifies semantic fidelity across scripts.

5 Results

5.1 General Translation Experiments on the Native Urdu Dataset

We evaluate baseline translation performance on the Native Urdu idiom dataset under three configurations: direct NMT translation without prompting, literal prompting, and cultural prompting (Table 4). Across all settings, semantic metrics (BERTScore, COMET, and XCOMET) remain relatively stable,

indicating that both NMT systems and LLMs generally preserve sentence-level meaning despite idiomatic constructions. In contrast, BLEU scores remain uniformly low, reflecting substantial surface-level variation, while ChrF shows only modest gains under literal and cultural prompting, suggesting that prompting alone only partially alleviates idiomatic ambiguity.

Table 5 provides a model-level analysis, distinguishing between NMT systems and LLMs. Among NMT models, Google Translate and NLLB-600M achieve the strongest semantic scores, with NLLB-3.3B performing comparably but slightly weaker on XCOMET, while all NMT systems exhibit low lexical overlap. Among LLMs, GPT-OSS-20B achieves the highest COMET, XCOMET, and ChrF scores, demonstrating strong idiomatic robustness and cross-lingual adequacy, with Qwen3-8B also performing competitively. Smaller models such as Gemma 3-1B and Mistral 7B perform consistently worse across metrics, indicating limited idiom handling without explicit reasoning. Overall, the results suggest that while NMT systems remain reliable for general adequacy, larger LLMs—particularly GPT-OSS-20B—are better suited for semantically faithful idiomatic translation.

Experiment	BERT	COMET	XCOMET	BLEU	ChrF
Direct (NMT)	0.87	0.63	0.62	0.02	21.39
Literal Prompt	0.89	0.64	0.66	0.02	24.69
Cultural Prompt	0.89	0.64	0.67	0.03	25.36

Table 4: General translation performance across three conditions: direct NMT translation, literal prompting, and cultural prompting.

Model	BERT	COMET	XCOMET	BLEU	ChrF
Qwen 3-8B	0.88	0.65	0.71	0.03	25.38
DeepSeek-8B	0.87	0.59	0.59	0.02	21.82
NLLB-200 3.3B	0.87	0.63	0.61	0.02	21.14
NLLB-200 600M	0.87	0.64	0.62	0.02	21.33
Gemma 3-1B	0.87	0.59	0.57	0.01	20.02
Google Translate	0.86	0.63	0.64	0.03	21.70
LLaMA 3.2-3B	0.88	0.60	0.64	0.02	21.64
Mistral 7B	0.87	0.57	0.56	0.01	19.43
GPT-OSS-20B	0.89	0.66	0.74	0.04	26.33

Table 5: Model-wise translation performance on the Native Urdu idiom dataset across semantic and lexical metrics.

5.2 Comparison of Few-Shot and Zero-Shot Translations

To evaluate the influence of in-context learning on idiomatic Urdu–English translation, we compare two few-shot settings (5-shot and 10-shot) with two representative zero-shot prompting strategies (Literal and Cultural). Table 6 reports BERTScore, COMET, XCOMET, BLEU, and ChrF scores for all settings.

Few-shot prompting consistently improves idiomatic translation quality over zero-shot settings. Semantic similarity increases with in-context examples, with both 5-shot and 10-shot configurations achieving higher BERTScore (0.91) than zero-shot (0.89). Adequacy-focused metrics show similar gains, as COMET rises from 0.64 to 0.67 and XCOMET from 0.66–0.67 to 0.70–0.71, indicating better cross-lingual meaning preservation and idiomatic fidelity. Surface-level overlap remains low in BLEU (0.02–0.03), reflecting the non-literal nature of idiomatic translation, while ChrF provides clearer improvements for few-shot prompting (28.12–28.35 vs. 24.69–25.36). Increasing the number of examples beyond five yields only marginal gains, suggesting diminishing returns from larger shot sizes.

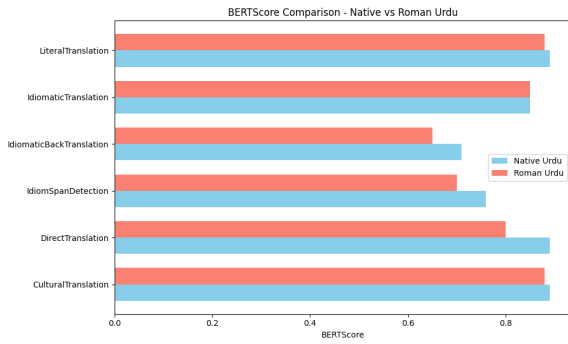
5.3 Paraphrasing Experiment

To evaluate how well models generate fluent English paraphrases of Urdu idioms and sentences while preserving meaning, we assess six LLMs using semantic similarity, surface-level lexical overlap, levenshtein distance, and multilingual embedding coherence. Table 7 presents the results.

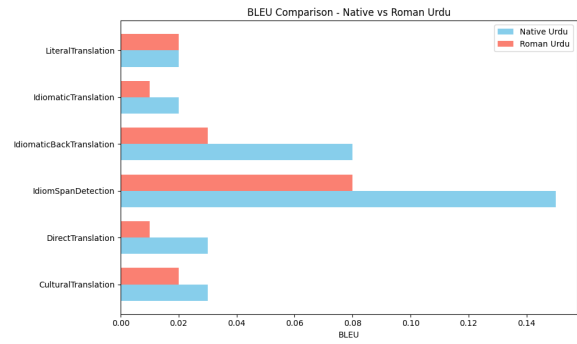
Across evaluation metrics, GPT-OSS-20B and Qwen3-8B consistently outperform smaller models. They achieve the highest semantic fidelity (BERTScore up to 0.90), indicating strong meaning preservation, while also leading in lexical overlap (ROUGE up to 0.31), reflecting closer surface correspondence to reference paraphrases. Levenshtein distances between 39 and 43 across all systems suggest substantial rewriting behavior, with GPT-OSS-20B maintaining the highest unigram overlap (32.45%). Cross-lingual embedding-based evaluations using LASER, LaBSE, and multilingual USE further confirm that GPT-OSS-20B and Qwen3-8B better preserve deep semantic structure, whereas smaller models such as Gemma and Mistral show slightly weaker alignment.

5.4 Idiom Span Detection Results

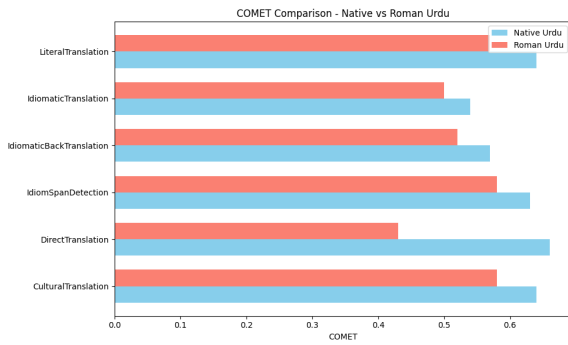
In addition to translation and paraphrasing, we evaluated the models on their ability to detect idiomatic expressions within Urdu sentences. This experiment measures how accurately models can identify



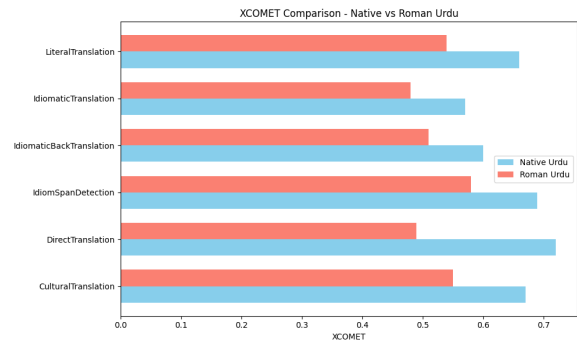
(a) BERTScore comparison



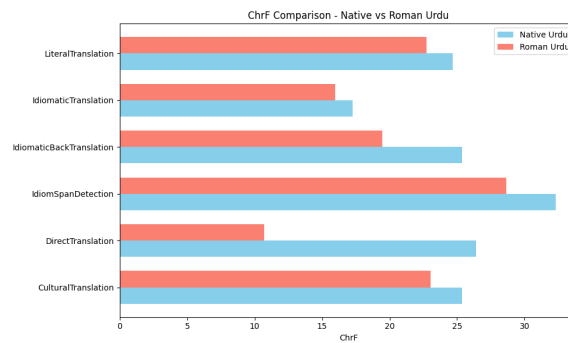
(b) BLEU comparison



(c) COMET comparison



(d) XCOMET comparison



(e) ChrF comparison

Figure 3: Model performance comparison across Urdu and Roman Urdu datasets using (a) BERTScore, (b) BLEU, (c) COMET, and (d) XCOMET.

Experiment	BERT	COMET	XCOMET	BLEU	ChrF
Few-Shot Experiments					
5-Shot	0.91	0.67	0.70	0.02	28.12
10-Shot	0.91	0.65	0.71	0.03	28.35
Zero-Shot Experiments					
Literal Translation	0.89	0.64	0.66	0.02	24.69
Cultural Translation	0.89	0.64	0.67	0.03	25.36

Table 6: Comparison of few-shot and zero-shot idiomatic translation performance across semantic and lexical metrics.

Model	BERT	ROUGE	Lev Dist	Uni Overlap	Cosine Similarity		
					LASER	USE	LaBSE
Qwen3-8B	0.90	0.30	39.36	30.42	0.56	0.53	0.55
DeepSeek-8B	0.89	0.26	42.34	27.57	0.49	0.48	0.51
Gemma-3-1B	0.89	0.21	42.67	21.21	0.45	0.43	0.46
Llama-3.2-3B	0.89	0.26	40.65	27.02	0.52	0.47	0.51
Mistral-7B	0.88	0.23	42.18	22.36	0.44	0.42	0.45
GPT-OSS-20B	0.90	0.31	40.36	32.45	0.57	0.54	0.57

Table 7: Paraphrasing results using BERTScore, ROUGE, Levenshtein distance, unigram overlap, and embedding similarities.

the span of an idiom, which is critical for subsequent translation or paraphrasing tasks.

For evaluation, we used ROUGE to measure n-gram overlap between predicted and reference spans, Levenshtein distance and unigram overlap for exact lexical matches, and cosine similarity using multilingual embeddings to capture semantic similarity across scripts.

Qwen3-8B and GPT-OSS-20B show the highest BERTScore and cosine similarity, suggesting strong semantic alignment with reference idioms. Llama-3.2B also performs well, particularly in unigram overlap, indicating precise token-level idiom detection. Lower scores for Mistral-7B and Gemma-3-1B indicate challenges in identifying idiomatic spans accurately. This experiment highlights the variation in models’ capability to detect idioms, which directly impacts downstream translation quality.

5.5 Roman vs Native Urdu

We analyze the impact of writing systems by comparing translation quality for Native Urdu and Roman Urdu inputs (Figure 3). Semantic metrics (BERTScore, COMET, and XCOMET) remain largely stable across scripts, indicating robust mean-

ing preservation despite Romanization. In contrast, surface-level metrics such as BLEU and ChrF exhibit noticeable declines for Roman Urdu, reflecting reduced n-gram and character-level overlap. This trend holds across all translation strategies, including Cultural, Literal, Direct, Idiomatic Back-Translation, and Idiom Span Detection. Overall, the results suggest that while semantic adequacy is resilient to script variation, lexical fidelity degrades, highlighting the need to jointly consider semantic and surface-level metrics in multi-script translation settings.

6 Conclusion

We present the first systematic benchmark for Urdu–English idiomatic translation across Native and Roman Urdu, covering translation, paraphrasing, idiom span detection, and back-translation. Our results show that large LLMs, particularly GPT-OSS-20B and Qwen3-8B, outperform NMT systems in preserving idiomatic meaning, with cultural and idiomatic prompts yielding the highest semantic fidelity. Few-shot prompting provides additional gains.

Cross-script analysis reveals that semantic metrics remain stable across Romanization, while

Model	BERT	ROUGE	Lev Dist	Uni Overlap	Cosine Similarity		
					LASER	USE	LaBSE
Qwen3-8B	0.83	0.48	15.93	49.57	0.70	0.62	0.66
DeepSeek-8B	0.73	0.27	27.18	27.18	0.54	0.43	0.48
Gemma-3-1B	0.68	0.15	38.21	17.06	0.49	0.43	0.46
Llama-3.2-3B	0.82	0.55	19.62	57.35	0.69	0.61	0.66
Mistral-7B	0.66	0.05	36.64	5.83	0.45	0.40	0.45
GPT-OSS-20B	0.84	0.53	15.67	53.37	0.72	0.66	0.69

Table 8: Idiom Span Detection results using ROUGE, Levenshtein distance, unigram overlap, and embedding similarities.

surface-level metrics like BLEU and ChrF degrade, highlighting challenges posed by inconsistent Roman Urdu orthography. Idiom span detection further emphasizes that accurate identification of figurative expressions is critical for translation quality.

Overall, this study offers a unified multi-metric evaluation framework and benchmark for idiomatic Urdu–English translation, providing insights for LLM development and cross-script evaluation. Future work can explore Romanization normalization and expansion to other low-resource languages.

Limitations

While this study provides a comprehensive benchmark for Urdu–English idiomatic translation, there are several limitations. First, the dataset covers only 4,000 idiom-sentence pairs, which may not fully represent the diversity of Urdu idioms in real-world usage. Second, Roman Urdu lacks standardized orthography, leading to evaluation challenges and potential inconsistencies in model performance. Third, our experiments are limited to the selected open-source LLMs and NMT systems; results may vary with other architectures or commercial models. Finally, GPU constraints restricted batch sizes and prompt variations, which may affect scalability and generalization to larger datasets or more complex idiomatic constructions.

References

Umer Butt, Stalin Varanasi, and Günter Neumann. 2025. [Low-resource transliteration for Roman-Urdu and Urdu using transformer-based models](#). In *Proceedings of the Eighth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2025)*, pages 144–153, Albuquerque, New Mexico, U.S.A. Association for Computational Linguistics.

Google DeepMind and Gemma Team. 2025. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.

Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2018. [Examining the tip of the iceberg: A data set for idiom translation](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline

465	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	529
466	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	530
467	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	531
468	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	532
469	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	533
470	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	534
471	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	535
472	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vas-	536
473	ic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	537
474	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	538
475	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	539
476	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	540
477	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	541
478	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	542
479	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	543
480	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	544
481	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	545
482	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	546
483	ran Narang, Sharath Raparthy, Sheng Shen, Shengye	547
484	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	548
485	denhende, Soumya Batra, Spencer Whitman, Sten	549
486	Sootla, Stephane Collot, Suchin Gururangan, Syd-	550
487	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	551
488	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	552
489	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	553
490	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	554
491	Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-	555
492	ginie Do, Vish Vogeti, Vitor Albiero, Vladan Petro-	556
493	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	557
494	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	558
495	aofang Wang, Xiaoming Ellen Tan, Xide Xia, Xin-	559
496	feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-	560
497	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yi-	561
498	wen Song, Yuchen Zhang, Yue Li, Yuning Mao,	562
499	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	563
500	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	564
501	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	565
502	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	566
503	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	567
504	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	568
505	gani, Amos Teo, Anam Yunus, Andrei Lupu, And-	569
506	res Alvarado, Andrew Caples, Andrew Gu, Andrew	570
507	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	571
508	dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-	572
509	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	573
510	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	574
511	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	575
512	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	576
513	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	577
514	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	578
515	Brian Gamido, Britt Montalvo, Carl Parker, Carly	579
516	Burton, Catalina Mejia, Ce Liu, Changhan Wang,	580
517	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	581
518	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	582
519	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	583
520	Daniel Kreymer, Daniel Li, David Adkins, David Xu,	584
521	Davide Testuggine, Delia David, Devi Parikh, Di-	585
522	ana Liskovich, Didem Foss, Dingkan Wang, Duc	586
523	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	587
524	Elaine Montgomery, Eleonora Presani, Emily Hahn,	588
525	Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban	589
526	Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	590
527	Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Oz-	591
528	genel, Francesco Caggioni, Frank Kanayet, Frank	592
	Seide, Gabriela Medina Florez, Gabriella Schwarz,	
	Gada Badeer, Georgia Swee, Gil Halpern, Grant	
	Herman, Grigory Sizov, Guangyi, Zhang, Guna	
	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	
	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	
	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	
	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	
	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	
	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	
	Geboski, James Kohli, Janice Lam, Japhet Asher,	
	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	
	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	
	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	
	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	
	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	
	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	
	delwal, Katayoun Zand, Kathy Matosich, Kaushik	
	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	
	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	
	Huang, Lailin Chen, Lakshya Garg, Lavender A,	
	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	
	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	
	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	
	Martynas Mankus, Matan Hasson, Matthew Lennie,	
	Matthias Reso, Maxim Groshev, Maxim Naumov,	
	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	
	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	
	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	
	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	
	Mo Metanat, Mohammad Rastegari, Munish Bansal,	
	Nandhini Santhanam, Natascha Parks, Natasha White,	
	Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas	
	Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev,	
	Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia	
	Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent,	
	Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner,	
	Philip Bontrager, Pierre Roux, Piotr Dollar, Polina	
	Zvyagina, Prashant Ratanchandani, Pritish Yuvraj,	
	Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi	
	Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mi-	
	tra, Rangaprabhu Parthasarathy, Raymond Li, Re-	
	bekkah Hogan, Robin Battey, Rocky Wang, Russ	
	Howes, Ruty Rinott, Sachin Mehta, Sachin Siby,	
	Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara	
	Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan,	
	Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto,	
	Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-	
	say, Sheng Feng, Shenghao Lin, Shengxin Cindy	
	Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang,	
	Shuqiang Zhang, Sinong Wang, Sneha Agarwal,	
	Soji Sajuyigbe, Soumith Chintala, Stephanie Max,	
	Stephen Chen, Steve Kehoe, Steve Satterfield, Sudar-	
	shan Govindaprasad, Sumit Gupta, Summer Deng,	
	Sungmin Cho, Sunny Virk, Suraj Subramanian,	
	Sy Choudhury, Sydney Goldman, Tal Remez, Tamar	
	Glaser, Tamara Best, Thilo Koehler, Thomas Robin-	
	son, Tianhe Li, Tianjun Zhang, Tim Matthews, Timo-	
	thy Chou, Tzook Shaked, Varun Vontimitta, Victoria	
	Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish	
	Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru,	
	Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,	
	Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will	
	Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan	
	Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yan-	

593	jun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang,	Sara Rezaeimanesh, Faezeh Hosseini, and Yadollah	649
594	Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang,	Yaghoobzadeh. 2025. Large language models for	650
595	Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi	Persian-English idiom translation . In <i>Proceedings of</i>	651
596	He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhao-	<i>the 2025 Conference of the Nations of the Americas</i>	652
597	duo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu	<i>Chapter of the Association for Computational Lin-</i>	653
598	Ma. 2024. The llama 3 herd of models . <i>Preprint</i> ,	<i>guistics: Human Language Technologies (Volume 1:</i>	654
599	arXiv:2407.21783.	<i>Long Papers)</i> , pages 7974–7985, Albuquerque, New	655
600	Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa	Mexico. Association for Computational Linguistics.	656
601	Coheur, Pierre Colombo, and André F. T. Martins.		
602	2024. xcomet: Transparent machine translation eval-	NLLB Team, Marta R. Costa-jussà, James Cross, Onur	657
603	uation through fine-grained error detection . <i>Transac-</i>	Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hef-	658
604	<i>tions of the Association for Computational Linguis-</i>	ernan, Elahe Kalbassi, Janice Lam, Daniel Licht,	659
605	<i>tics</i> , 12:979–995.	Jean Maillard, Anna Sun, Skyler Wang, Guillaume	660
606	Deli Guo, Bo Zhang, Haoran Wei, Yufan Wang,	Wenzek, Al Youngblood, Bapi Akula, Loic Bar-	661
607	and DeepSeek-AI. 2025. Deepseek-r1: Incent-	rault, Gabriel Mejia Gonzalez, Prangthip Hansanti,	662
608	ivizing reasoning capability in llms . <i>Preprint</i> ,	John Hoffman, Semarley Jarrett, Kaushik Ram	663
609	arXiv:2501.12948.	Sadagopan, Dirk Rowe, Shannon Spruit, Chau	664
610	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	Tran, Pierre Andrews, Necip Fazil Ayan, Shruti	665
611	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Bhosale, Sergey Edunov, Angela Fan, Cynthia	666
612	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	Gao, Vedanuj Goswami, Francisco Guzmán, Philipp	667
613	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	Koehn, Alexandre Mourachko, Christophe Rop-	668
614	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	ers, Safiyyah Saleem, Holger Schwenk, and Jeff	669
615	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	Wang. 2022. No language left behind: Scal-	670
616	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	ing human-centered machine translation . <i>Preprint</i> ,	671
617	arXiv:2310.06825.	arXiv:2207.04672.	672
618	Deepak Kumar, Divakar Yadav, and Yash Patel. 2025.	Abdul Wahid, Muhammad Arsalan, Adeel Mumtaz, Rafi-	673
619	Gpt-oss-20b: A comprehensive deployment-centric	ullah Khan, and Saeed Anwar. 2024. Erupd: En-	674
620	analysis of openai’s open-weight mixture of experts	glish to roman urdu parallel dataset . <i>arXiv preprint</i>	675
621	model . <i>Preprint</i> , arXiv:2508.16700.	<i>arXiv:2412.17562</i> .	676
622	Chin-Yew Lin. 2004. ROUGE: A package for auto-	An Yang, Anpeng Li, Baosong Yang, Beichen Zhang,	677
623	matic evaluation of summaries . In <i>Text Summariza-</i>	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	678
624	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	Chengen Huang, Chenxu Lv, Chuji Zheng, Dayi-	679
625	Association for Computational Linguistics.	heng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge,	680
626	Emmy Liu, Aditi Chaudhary, and Graham Neubig. 2023.	Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jian-	681
627	Crossing the threshold: Idiomatic machine translation	hong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang,	682
628	through retrieval augmentation and loss weighting .	Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang,	683
629	In <i>Proceedings of the 2023 Conference on Empiri-</i>	Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng,	684
630	<i>cal Methods in Natural Language Processing</i> , pages	Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng	685
631	15095–15111, Singapore. Association for Computa-	Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu,	686
632	tional Linguistics.	Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin,	687
633	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xu-	688
634	Jing Zhu. 2002. Bleu: a method for automatic evalua-	ancheng Ren, Yang Fan, Yang Su, Yichang Zhang,	689
635	tion of machine translation . In <i>Proceedings of the</i>	Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang,	690
636	<i>40th Annual Meeting of the Association for Computa-</i>	Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zi-	691
637	<i>tional Linguistics (ACL)</i> , pages 311–318.	han Qiu. 2025. Qwen3 technical report . <i>Preprint</i> ,	692
638	Maja Popovi��. 2015. chrF: character n-gram F-score	arXiv:2505.09388.	693
639	for automatic MT evaluation . In <i>Proceedings of the</i>	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Wein-	694
640	<i>Tenth Workshop on Statistical Machine Translation</i> ,	berger, and Yoav Artzi. 2020. Bertscore: Evaluating	695
641	pages 392–395, Lisbon, Portugal. Association for	text generation with bert . In <i>International Conference</i>	696
642	Computational Linguistics.	<i>on Learning Representations (ICLR)</i> .	697
643	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	A Appendix	698
644	Lavie. 2020. COMET: A neural framework for MT	A.1 Resources	699
645	evaluation . In <i>Proceedings of the 2020 Conference on</i>	All experiments were performed on machines	700
646	<i>Empirical Methods in Natural Language Processing</i>	equipped with NVIDIA A100 40GB GPUs.	701
647	<i>(EMNLP)</i> , pages 2685–2702, Online. Association for	A.2 Full Prompts	702
648	Computational Linguistics.	Table 9 summarizes all prompting strategies used	703
		in the experiments, along with the corresponding	704

Prompt	Prompt Instruction
Literal Translation	Translate the following Urdu text to English. Do not include explanations, reasoning, or additional text.
Cultural Translation	Translate the following Urdu text into English. Use natural expressions and preserve the cultural meaning. Do not include explanations, reasoning, or additional text.
Paraphrase Translation	Paraphrase the following Urdu sentence into fluent English, keeping the meaning intact. Avoid word-for-word translations. Do not include explanations, reasoning, or additional text.
Direct Translation	Translate this Urdu text to English. Do not include explanations, reasoning, or additional text.
Idiomatic Translation	Identify the idiom in this Urdu sentence and produce an equivalent English idiom. Do not translate word-for-word. Do not include explanations or reasoning.
Idiomatic Back Translation	Translate the following Urdu idiom to an equivalent English idiom, and then back to Urdu. Provide only the final back-translated Urdu idiom.
Idiom Span Detection	Detect and extract the Urdu idiom in the sentence. Provide only the final detected Urdu idiom.
Few-Shot Translation	Translate the following Urdu text to English using the examples as references.

Table 9: Summary of full prompts used in the experiments.

Prompt	Translation Result
Literal Translation	Her hands and feet swelled up.
Cultural Translation	She suddenly became extremely nervous.
Paraphrase Translation	She lost her composure and couldn't think straight.
Direct Translation	She panicked instantly.
Idiomatic Translation	She completely freaked out.
Idiomatic Back Translation	وہ گھبراہٹ میں آ گئی۔
Idiom Span Detection	ہاتھ پیر پھول جانا

Table 10: Sample translation results for the Urdu idiom ہاتھ پیر پھول جانا across different prompting strategies.

instructions given to the models. For few-shot prompts, the examples were randomly selected from the dataset.

A.3 Sample Translations

To illustrate the behavior of our prompting strategies, we provide a representative example from the dataset. The following idiom and its contextual usage serve as the reference (“gold”) instance:

- **Urdu Idiom:** ہاتھ پیر پھول جانا
- **English Equivalent Idiom:** to panic / to get flustered
- **Golden Urdu Sentence:** اچانک خبر سننے ہی اس کے ہاتھ پیر پھول گئے۔
- **Golden English Sentence:** She completely

panicked the moment she heard the sudden news.

Table 10 presents translation outputs generated using each prompt type.