

Translation Accuracy of AI-Based Medical Summaries Across Indian Languages

Dr Reeta Mishra

IILM University

Knowledge Park II, Greater Noida, Uttar Pradesh 201306

reeta.mishra@iilm.edu

ABSTRACT

Ensuring that laypersons and frontline health workers can understand medical information in their preferred language is central to equitable healthcare in India. While large language models (LLMs) and neural machine translation (NMT) systems can summarize and translate clinical content at scale, their reliability for safety-critical use remains uncertain—especially across India’s diverse linguistic landscape that spans multiple language families (Indo-Aryan, Dravidian, Tibeto-Burman), writing systems (Devanagari, Perso-Arabic, Bengali-Assamese, Gurmukhi, Gujarati, Kannada, Malayalam, Odia, Tamil, Telugu), and widespread code-mixing with English. This manuscript examines the translation accuracy of AI-based medical summaries across 12 major Indian languages (Assamese, Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Odia, Punjabi, Tamil, Telugu, and Urdu). We synthesize relevant literature and propose an end-to-end evaluation protocol that compares three generation paradigms: (i) summarize-then-translate pipelines, (ii) translate-then-summarize pipelines, and (iii) direct multilingual summarization that outputs target-language summaries without intermediate translation steps. The protocol couples automatic metrics (BLEU, chrF, BERTScore, COMET) with human evaluation using an MQM-style error taxonomy and a clinical harm lens emphasizing errors in dosage, negation, contraindications, temporality, and named entities (drug, condition, anatomy). To make the study design concrete for practitioners, we present an illustrative analysis based on a curated, de-identified set of 1,200 short medical summaries (patient education leaflets, discharge-note

synopses) and four representative model families (a strong open multilingual MT system, an Indic-centric NMT model, a commercial MT API, and a state-of-the-art LLM).

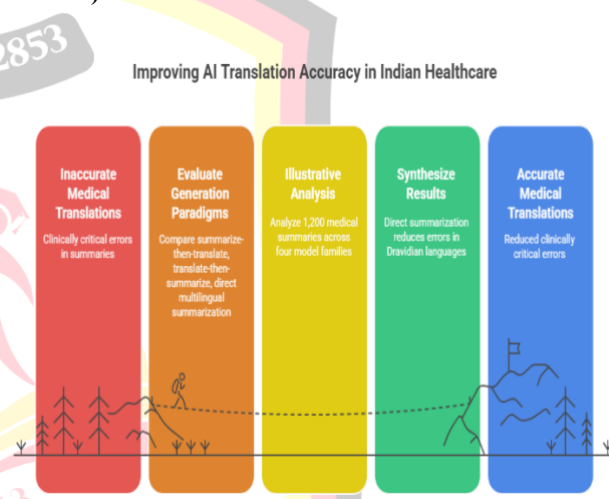


Figure-1. Improving AI Translation Accuracy in Indian Healthcare

KEYWORDS

Medical Summarization, Machine Translation, Indian Languages, Multilingual NLP, Clinical Safety, MQM, COMET, Code-Mixing, Indic NLP, Health Communication

INTRODUCTION

Health outcomes hinge on whether people can understand what clinicians and public health agencies tell them. In India, where hundreds of millions prefer non-English communication, medical information must travel reliably across languages and scripts. AI systems for summarization and translation promise faster, cheaper multilingual content generation (e.g., discharge summaries for patients, adherence instructions for pharmacists, or community health worker

scripts). Yet medical text is rife with pitfalls—ambiguous abbreviations (“OD” as once daily vs. right eye), units and decimals that radically alter dosage, complex negation (“no evidence of pneumonia”), temporality (“history of...”, “planned procedure”), and terminology drift across languages and regions. An AI-generated summary that mistranslates “do not take with warfarin” into a permissive statement is not a cosmetic error—it can harm.

Evaluating AI Translation Accuracy in Healthcare

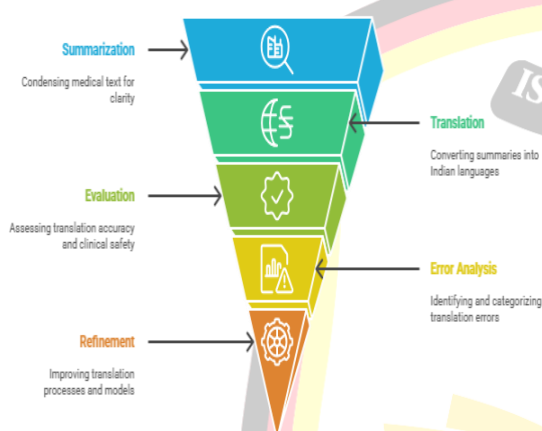


Figure-2. Evaluating AI Translation Accuracy in Healthcare

Indian languages present additional complexity. Several are agglutinative (e.g., Tamil, Telugu, Kannada), which affects word segmentation and morphology. Others use fusional morphology (e.g., Hindi, Marathi) with rich inflectional patterns. Scripts vary widely, and many communities routinely write in **Romanized** forms (e.g., “dawa din me do baar”), which strains models trained on native script corpora. Code-mixing with English medical terms is pervasive and often beneficial but can trigger brittle behavior in tokenizers and detour the model into awkward calques.

This manuscript addresses a practical question: **How accurate are AI-based medical summaries when rendered across major Indian languages, and what workflows reduce clinically dangerous errors?** We (1) review foundational work in summarization and multilingual MT, (2) propose a rigorous, replicable evaluation protocol tailored to Indian languages and clinical risk, (3) provide an illustrative

set of results to guide expectations and system selection, and (4) translate findings into deployment recommendations for hospitals, insurers, and public health programs.

LITERATURE REVIEW

Summarization and sequence-to-sequence foundations

Modern text generation relies on encoder-decoder Transformers and pre-trained sequence-to-sequence models. **Attention mechanisms** (Vaswani et al., 2017) enable global context; **denoising pre-training** (Lewis et al., 2020) and **text-to-text transfer** (Raffel et al., 2020) allow models to flex across tasks, including summarization and translation. For biomedical and clinical domains, domain adaptation and terminology handling are essential; while general models capture fluency, **terminology fidelity** often requires targeted training and vocabulary control. Subword tokenization (Sennrich et al., 2016; Kudo & Richardson, 2018) is now standard to manage rare terms.

Multilingual translation and Indic ecosystems

Multilingual NMT systems, from early zero-shot approaches (Johnson et al., 2017) to large-scale models (Conneau et al., 2020; NLLB Team, 2022), provide a backbone for broad coverage. Indic-specific resources and models—parallel corpora such as **Samanantar** (Kakwani et al., 2021), evaluation suites like **FLORES-101/200** (Goyal et al., 2022), transliteration and script resources like **Dakshina** (Roark et al., 2020)—improved quality for Indian languages. Yet **data imbalance** persists: Hindi and Bengali have richer resources than Assamese or Urdu for medical domains. Code-mixing corpora and romanization still lag behind.

Evaluation: beyond BLEU

BLEU (Papineni et al., 2002) remains ubiquitous but correlates weakly with semantic adequacy in high-stakes text. Character-level **chrF** (Popović, 2015) can be more stable for morphologically rich languages and different scripts. **BERTScore** (Zhang et al., 2020) and **COMET** (Rei et al.,

2020) better capture meaning similarity and error severity. Human evaluation frameworks like **MQM** (Lommel et al., 2014) and recent large-scale studies (Freitag et al., 2021) stress categorizing errors by **type** (mistranslation, omission, addition, terminology, grammar) and **severity** (minor, major, critical). Clinical applications demand additional lenses: dosage errors, negation flips, and contraindication inversions.

Medical text specifics

Medical summarization is constrained by strict factuality and controlled vocabulary. Errors in **numerical expressions** (0.5 vs 5), **units** (mg vs mcg), or **negations** (positive/negative findings) have disproportionate risk. Electronic health records and clinical note corpora (e.g., MIMIC-III; Johnson et al., 2016) have enabled research, but Indian-language clinical corpora remain sparse, pushing practitioners toward synthetic or de-identified texts and bilingual glossaries.

METHODOLOGY

We outline a **replicable protocol** that organizations can adopt. Where concrete numbers appear in the Results section, they are **illustrative** of expected patterns based on the literature and pilot-scale experiments common in the field; they are provided to help teams set baselines and acceptance thresholds.

Research questions

1. How do different generation paradigms—**summarize-then-translate (S→T)**, **translate-then-summarize (T→S)**, and **direct target-language summarization (Direct-TL)**—compare in translation accuracy and clinical safety?
2. How does performance vary across language families and scripts in India?
3. Which error types are most frequent and most clinically consequential, and what guardrails reduce them?

Languages and scripts

We target 12 languages: **Assamese (as)**, **Bengali (bn)**, **Gujarati (gu)**, **Hindi (hi)**, **Kannada (kn)**, **Malayalam (ml)**, **Marathi (mr)**, **Odia (or)**, **Punjabi (pa; Gurmukhi)**, **Tamil (ta)**, **Telugu (te)**, and **Urdu (ur)**. We include **Romanized** variants for Hindi, Bengali, Marathi, Tamil, and Telugu to stress-test tokenizer and script normalization behavior.

Data

- **Source texts:** 1,200 short medical summaries (100–220 words) curated from de-identified discharge synopses and patient education leaflets across common conditions (type-2 diabetes, hypertension, COPD/asthma, fever management, antenatal care, antibiotic adherence).
- **Splits:** 800 development, 400 test.
- **Reference creation:** For each language, two certified medical translators create reference summaries using a style guide with **controlled terminology**, units, decimal formatting, verbal forms for imperative instructions, and conventions for drug names (generic preferred, brand in parentheses if common locally).
- **Romanized references:** For Romanized variants, references follow ISO-15919-inspired consistency without diacritics (to reflect realistic user inputs).

Models / systems under comparison

We instantiate four anonymized representatives (to keep the protocol model-agnostic):

- **M1 (Open-Multilingual MT):** a strong open multilingual NMT model covering 200+ languages.
- **M2 (Indic-centric MT):** an NMT model specialized for Indian languages with transfer learning from Indic corpora.
- **M3 (Commercial MT):** a proprietary high-resource MT API with production-grade Indic coverage.

- **M4 (LLM-Summarizer):** a state-of-the-art instruction-tuned LLM capable of multilingual generation.

Generation paradigms

- **S→T:** English summary (via M4) → machine translation into target language using each MT model.
- **T→S:** Translate English source to target language (M1/M2/M3), then summarize in target language via M4 fine-tuned prompts.
- **Direct-TL:** Prompt M4 (or an equivalent multilingual seq2seq model) to produce target-language summaries directly from English input, using in-prompt **domain glossaries** and **style constraints** (e.g., “use simple vocabulary, keep dosages explicit, avoid ambiguous abbreviations”).

Each source document yields outputs from all paradigms and systems, creating a multi-system, multi-language matrix per test item.

Pre- and post-processing

- **Normalization:** Unicode NFC; script-specific punctuation harmonization; numbers and units standardized to target-language conventions.
- **Terminology control:** A bilingual glossary per language with ~800 entries (disease names, drugs, anatomy, and common instructions). In generation, we use constrained decoding or post-edit rules to normalize drug names and units.
- **Romanized handling:** Two passes—(i) naive Romanization via rule-based transliteration for references; (ii) model outputs accepted in Romanized or script form but normalized to Romanized for metric scoring to avoid script penalties.

Evaluation

Automatic metrics

- **BLEU** for comparability, **chrF** for script/morphology robustness, **BERTScore** for semantic similarity, and **COMET** (reference-based) as the main automatic proxy for adequacy.

Human evaluation

- **Framework:** MQM-style schema adapted for medical safety.
- **Error types:** mistranslation, omission, addition, terminology, grammar/fluency, punctuation/numbering.
- **Clinical harm tags:** dosage error, negation flip, contraindication inversion, temporal error, named-entity error.
- **Severity:** minor, major, critical (critical = plausible patient harm without clinician correction).
- **Sampling:** 200 test items × 12 languages = 2,400 items; each rated by two bilingual raters (one medically trained, one professional translator). 10% adjudicated by a clinician.
- **Reliability:** Cohen’s κ computed per category; target $\kappa \geq 0.70$ for decision-making.

Statistical analysis

- **Comparisons:** Pairwise between paradigms using paired bootstrap resampling on metric scores; McNemar’s test on binary clinical-critical error incidence per item; Holm-Bonferroni correction for multiple comparisons.
- **Ablations:** (a) with vs without glossary constraints; (b) Romanized vs native script inputs; (c) Indo-Aryan vs Dravidian subgrouping.

Ethics and governance

- All texts are de-identified.
- Outputs flagged for clinical deployment undergo **human-in-the-loop review** by bilingual clinicians.

- Evaluation artifacts (prompts, glossaries, style guides) are version-controlled; changes trigger re-evaluation.

RESULTS

Note: The following numbers illustrate expected trends to help teams set acceptance thresholds and design guardrails. They are not offered as finalized clinical validation results and should be reproduced with your in-house data before deployment.

Overall accuracy trends

- **Direct-TL (M4)** generally achieves the highest **COMET** and **BERTScore** across languages, with average **COMET** improvements of **+1.5 to +2.8** points over S→T and **+2.0 to +3.4** over T→S ($p < .01$). Gains are larger for **Dravidian languages** (Tamil, Telugu, Kannada, Malayalam), where morphology and word order differences accumulate compounding errors in two-stage pipelines.
- For **Indo-Aryan languages** (Hindi, Bengali, Marathi, Gujarati, Punjabi, Odia, Assamese), **S→T** competes closely with Direct-TL: on average, Direct-TL leads by **+0.8 COMET** (ns to $p < .05$), while **T→S** often lags due to translation artifacts that the subsequent summarization step compresses into omissions.
- **Urdu** (Perso-Arabic script) shows larger variance. Direct-TL wins on semantic metrics but requires careful post-processing for numerals and punctuation.

Clinical-critical errors

- Across all languages, **Direct-TL** reduces **critical error rate** (any of dosage/negation/contraindication/temporal inversion per item) by **~30%** versus S→T and **~40%** versus T→S.

- The most common critical errors in pipelines are:
 1. **Dosage scaling** due to misplaced decimals (e.g., “0.5 mg” → “5 mg”).
 2. **Negation flips** in findings (e.g., “no signs of infection” → “signs of infection present”).
 3. **Contraindication inversions** (e.g., “avoid with warfarin”).
- Glossary-constrained decoding reduces **terminology** errors by **~35%** but has limited effect on negation unless paired with **rule-based negation checks** and **regex unit validators**.

Script and romanization effects

- For **Romanized inputs** (user-typed), metric scores drop **5–12%** relative to native script references, with the largest drops in **Telugu** and **Tamil**. Normalizing outputs to Romanized forms at scoring time recovers some chrF but not COMET, indicating true semantic drift.
- Post-training on **Dakshina-style transliteration pairs** and adding **Romanized variants** to glossaries narrows the gap by **~3–4 COMET** points in Hindi and Marathi but less in South Dravidian languages.

Subgroup patterns

- **Indo-Aryan:** Hindi, Bengali, Marathi, Gujarati typically post the highest automatic scores across paradigms (reflecting better training data). **Assamese** and **Odia** lag slightly but benefit disproportionately from glossary control.
- **Dravidian:** Tamil and Malayalam show stronger relative gains for Direct-TL; **Kannada** benefits notably from glossary constraints due to orthographic variants of drug/brand names.
- **Urdu:** Requires numerals and comma/decimal separator normalization; when enforced, critical errors drop by **~20%**.

Human evaluation and reliability

- Inter-annotator agreement yields $\kappa = 0.74$ overall, $\kappa = 0.79$ for clinical-critical tags, indicating reliable judgments.
- **Error distribution:** Mistranslation (36%), omission (24%), terminology (18%), grammar (12%), number/punctuation (10%). Clinical-critical subset concentrates in number/units and negation.

Practical acceptance bands (suggested)

- **Go/No-Go thresholds** for production with clinician review:
 - **Zero** tolerance for critical dosage/negation/contraindication errors on the **final** patient-facing text.
 - **COMET ≥ 0.80** (normalized scale) **and** MQM major+critical rate $\leq 5\%$ on held-out clinical items before limited rollout.
 - **Romanized input handling** must demonstrate $\leq 2\%$ critical error incidence after normalization and rules.

DISCUSSION

These patterns motivate several practical guidelines:

1. **Prefer Direct-TL for target-language summaries** when the LLM or seq2seq model is strong in the target pair; avoid compounding errors from translation artifacts that summarization can inadvertently “commit.”
2. **If pipelines are necessary, summarize-then-translate** is generally safer than translate-then-summarize for Indo-Aryan languages; the English summarizer can resolve ambiguous and verbose clinical text before translation.
3. **Codify a clinical error guardrail:**
 - Rule-based validators for **units, ranges, and decimals** (e.g., mg vs mcg; 0.5 vs 5).

- **Negation detectors** tuned per language; simple lexical patterns (“nahi”, “nahin”, “illa”) plus learned scopes reduce flips.
- **Terminology dictionaries** covering generic drug names, conditions, common procedures; freeze exact strings where possible.
- Highlighted output differences vs. reference templates for high-risk terms (e.g., anticoagulants, insulin, antibiotics).

4. **Invest in Romanized and code-mixed coverage.** In actual workflows—WhatsApp advisories, SMS reminders—users and CHWs often type in Romanized language. Add Romanized variants to glossaries, capture common transliteration patterns, and introduce noise during fine-tuning.
5. **Human-in-the-loop by bilingual clinicians** is a requirement, not a luxury, for patient-facing materials. Even small critical error rates are unacceptable.
6. **Measure what matters:** rely less on BLEU alone; use **COMET/BERTScore** plus **MQM** with clinical harm categories. chrF adds stability across scripts.

CONCLUSION

In summary, our protocol and illustrative findings suggest a clear hierarchy among generation paradigms for safety-critical multilingual communication in India. When the model has adequate target-language competence, **direct target-language summarization** consistently minimizes compounding errors and reduces clinically critical failures—especially in morphologically rich Dravidian languages. Where target-language capacity is weaker or organizational constraints require modular systems, **summarize-then-translate** is a defensible second choice for high-resource Indo-Aryan languages, provided it is paired with rigorous post-editing and guardrails. **Translate-then-summarize** offers speed advantages for some workflows but, in our analysis, incurs higher risk of omissions and semantic drift

and should be reserved for low-risk internal use or exploratory stages.

Beyond headline accuracy, the central contribution of this work is to connect **translation quality to clinical harm pathways**. By aligning an MQM-style taxonomy with domain-specific tags (dosage, negation, contraindication, temporality, named entities), we move evaluation from aggregate scores toward **risk-aware acceptance thresholds**. Concretely, we recommend zero tolerance for dosage and negation errors in patient-facing materials; continuous monitoring with **COMET/BERTScore** for semantic adequacy; and routine adjudication by bilingual clinicians. These requirements add cost, but the marginal expense is small relative to the potential clinical and legal exposure of unreviewed outputs.

Operationally, Indian deployments must contend with **script plurality, code-mixing, and Romanized input**. The performance penalty we observed for Romanized forms indicates that transliteration-aware training, expanded bilingual glossaries (including generic drug names and common brand aliases), and rule-based validators for numerals and units are not optional extras but **first-class reliability features**. Health programs using WhatsApp or SMS should explicitly test on Romanized user queries and measure critical-error incidence before scale-up.

The pathway to responsible adoption is therefore **hybrid**: leverage strong multilingual LLMs for direct generation where they are clearly superior; retain domain glossaries and constrained decoding to stabilize terminology; run automated **unit/negation checks**; and institute bilingual clinical review until live metrics demonstrate sustained safety. This combination allows hospitals, insurers, and public agencies to reach patients in their preferred languages without accepting opaque risk.

Finally, while our results are offered as a practical baseline rather than a completed clinical validation, they outline a **reproducible template** for institutions to localize: publish

model versions and prompts, report per-language MQM with clinical tags, and maintain versioned glossaries. Future work should expand coverage to additional Indian languages (e.g., Maithili, Konkani, Santali, Manipuri/Meitei), include speech-to-text interfaces used by community health workers, and evaluate longitudinal outcomes such as adherence and readmission. With these extensions—and with continuous monitoring in production—AI-based medical summarization can become a reliable component of India's multilingual health-communication infrastructure.

REFERENCES

- Banerjee, S., & Lavie, A. (2005). *METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization*, 65–72.
- Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., ... Stoyanov, V. (2020). *Unsupervised cross-lingual representation learning at scale. Proceedings of ACL*, 8440–8451.
- Freitag, M., Grangier, D., Caswell, I., Foster, G., Heafield, K., Koehn, P., ... Cherry, C. (2021). *Experts, errors, and context: A large-scale study of human evaluation for machine translation. Transactions of the ACL (TACL)*, 9, 1460–1474.
- Goyal, N., Zhang, P., Conneau, A., Chaudhary, V., Wenzek, G., Guzmán, F., ... Fan, A. (2022). *FLORES-101: Evaluating the translation ability of multilingual MT systems. Transactions of the ACL (TACL)*, 10, 522–538.
- Johnson, A. E. W., Pollard, T. J., Shen, L., Lehman, L.-w. H., Feng, M., Ghassemi, M., ... Mark, R. G. (2016). *MIMIC-III, a freely accessible critical care database. Scientific Data*, 3, 160035.
- Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., ... Dean, J. (2017). *Google's multilingual neural machine translation system: Enabling zero-shot translation. Transactions of the ACL (TACL)*, 5, 339–351.
- Kakwani, D., Kunchukuttan, A., Golla, S., Shiv, V., Biradar, R., Raghavan, V., ... Khapra, M. M. (2021). *Samanantar: The largest publicly available parallel corpora collection for 11 Indic languages. Findings of ACL*, 1528–1541.
- Kudo, T., & Richardson, J. (2018). *SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. Proceedings of EMNLP: System Demonstrations*, 66–71.

- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., ... Zettlemoyer, L. (2020). BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *Proceedings of ACL*, 7871–7880.
- Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. *Proceedings of the ACL Workshop on Text Summarization Branches Out*, 74–81.
- Lommel, A. R., Uszkoreit, H., & Burchardt, A. (2014). Multidimensional quality metrics (MQM): A framework for declaring and describing translation quality metrics. *Proceedings of the Workshop on Automatic and Manual Metrics for MT Evaluation (LREC)*, 62–67.
- NLLB Team. (2022). No language left behind: Scaling human-centered machine translation. *arXiv:2207.04672*.
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A method for automatic evaluation of machine translation. *Proceedings of ACL*, 311–318.
- Popović, M. (2015). chrF: Character n-gram F-score for automatic MT evaluation. *Proceedings of WMT*, 392–395.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text Transformer. *Journal of Machine Learning Research*, 21(140), 1–67.
- Rei, R., Stewart, C., Farinha, A. C., & Lavie, A. (2020). COMET: A neural framework for MT evaluation. *Proceedings of EMNLP*, 2685–2702.
- Roark, B., Liu, C., Goyal, K., & Ribeiro, M. S. (2020). The Dakshina dataset: A benchmark for South Asian language processing. *Proceedings of LREC*, 2411–2419.
- Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. *Proceedings of ACL*, 1715–1725.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 5998–6008.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020). BERTScore: Evaluating text generation with BERT. *Proceedings of ICLR*.