

# **Learning from the Assignment**

This assignment provided a comprehensive learning experience by combining theoretical knowledge with practical skills in model evaluation, dataset handling, and working with low-resource languages. The key learnings can be categorised into several aspects:

## **1. Understanding Machine Translation Models**

Through this assignment, I gained insights into different machine translation models, including:

1. **Distilled NLLB-200 Model:** A large pre-trained model with 600M parameters, this exposed me to the capabilities of large-scale language models for translation tasks.
2. **IndicTrans Model:** IndicTrans is a powerful open-source NMT model for all 22 scheduled Indic languages in India. It excels at high-quality translation, even for languages with limited resources, thanks to its use of transformer architecture. This makes it a versatile tool for anyone working with translation tasks involving these languages.
3. **ChatGPT Model:** A broader language model, evaluating this model for translation tasks provided insights into the strengths and limitations of more general-purpose language models for specific tasks like machine translation.

Evaluating these diverse models allowed me to understand their unique architectures, training methodologies, and capabilities, enabling a comparative analysis of their strengths and weaknesses in different translation scenarios.

## **2. Evaluation Metrics for Translation Quality**

Understanding BLEU (Bilingual Evaluation Understudy) and ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metrics and using them to evaluate translation quality.

Computing and interpreting these scores helped me:

1. Understand how translation quality is objectively measured and quantified.
2. Identify the factors that contribute to good or poor translation performance, such as handling of morphology, syntax, and semantics.
3. Gain insights into the strengths and weaknesses of different models by analysing their performance across various language pairs and domains.

### **3. Working with Low-Resource Language Data**

The assignment required working with Indian languages, which are often considered low-resource in the context of machine translation. This exposure was invaluable in understanding the challenges and considerations involved in building and evaluating models for low-resource languages, such as:

1. Limited availability of parallel corpora and linguistic resources.
2. Handling language-specific syntax, and cultural context during translation.
3. Techniques for data augmentation and transfer learning to improve model performance.

### **4. Data Handling and Preprocessing**

Working with the Samanantar benchmark dataset and the WAT2021 test dataset involved data handling, preprocessing, and managing large parallel corpora. I learned how to:

Work with different data formats and handle encoding issues.

Preprocess and prepare data for machine translation tasks, ensuring data quality and consistency.

Manage and organise large datasets effectively for model training and evaluation.

### **5. Benchmarking and Model Comparison**

By evaluating multiple models on the same datasets, I gained valuable experience in benchmarking and comparing the performance of different machine translation systems. This is a crucial skill in machine learning research and model development, allowing me to:

1. Establish baselines and set performance benchmarks for different translation tasks.
2. Identify the strengths and weaknesses of each model in handling specific language pairs or domains.
3. Draw insights from comparative analyses to inform future model selection and development.

### **6. Exposure to State-of-the-Art Resources**

The assignment introduced me to valuable resources and tools in the field of natural language processing and machine translation, including:

1. Hugging Face Transformers library: A powerful library for working with pre-trained language models and transformer architectures.

2. IndicTrans model: It tackles all 22 official Indian languages, including those with unique writing systems. Developed by AI4Bharat, this free-to-use model prioritises accessibility and empowers projects focused on diverse Indian languages.
3. OpenAI's ChatGPT API: Exposure to a large language model in translating one language to another, enabling exploration of its translation capabilities.

Familiarity with such resources and tools will be beneficial for future work in computational linguistics and natural language processing.

## **7. Critical Thinking and Analysis**

Throughout the assignment, I exercised critical thinking and analytical skills by:

1. Interpreting evaluation results and drawing insights from the observations.
2. Identifying strengths and weaknesses of different models based on their performance across various language pairs and domains.
3. Analysing the impact of different modelling approaches, dataset characteristics, and evaluation metrics on translation quality.

## 9. Model Comparison

### 1. English to Hindi Translation (en\_hi):

NLLB:

BLEU Score: 0.6179  
ROUGE-1 F1-score: 0.5865  
ROUGE-2 F1-score: 0.3496  
ROUGE-L F1-score: 0.5480

IndicTrans:

BLEU Score: 0.6975  
ROUGE-1 F1-score: 0.6245  
ROUGE-2 F1-score: 0.3939  
ROUGE-L F1-score: 0.5888

ChatGPT:

BLEU Score: 0.6804  
ROUGE-1 F1-score: 0.6056  
ROUGE-2 F1-score: 0.3571  
ROUGE-L F1-score: 0.5585

IndicTrans performs the best in terms of BLEU score (0.6975), indicating higher similarity between translations and reference texts compared to NLLB (0.6179) and ChatGPT (0.6804). IndicTrans also achieves the highest ROUGE scores (ROUGE-1, ROUGE-2, ROUGE-L) among the three models, demonstrating better overlap in unigrams, bigrams, and longest common subsequences with reference translations.

### 2. Hindi to English Translation (hi\_en):

NLLB:

BLEU Score: 0.6677  
ROUGE-1 F1-score: 0.6072  
ROUGE-2 F1-score: 0.3834  
ROUGE-L F1-score: 0.5755

IndicTrans:

BLEU Score: 0.7530  
ROUGE-1 F1-score: 0.6675  
ROUGE-2 F1-score: 0.4529  
ROUGE-L F1-score: 0.6332

ChatGPT:

BLEU Score: 0.7376  
ROUGE-1 F1-score: 0.6545  
ROUGE-2 F1-score: 0.4140  
ROUGE-L F1-score: 0.6165

IndicTrans outperforms both NLLB and ChatGPT in BLEU score (0.7530) and ROUGE scores (ROUGE-1, ROUGE-2, ROUGE-L) for Hindi to English translation.

ChatGPT follows closely behind IndicTrans, showing competitive performance in this translation direction.

### **3. Hindi to Gujarati Translation (hi\_gu):**

NLLB:

BLEU Score: 0.5630  
ROUGE-1 F1-score: 0.4947  
ROUGE-2 F1-score: 0.2592  
ROUGE-L F1-score: 0.4748

IndicTrans:

BLEU Score: 0.6370  
ROUGE-1 F1-score: 0.5026  
ROUGE-2 F1-score: 0.2520  
ROUGE-L F1-score: 0.4791

ChatGPT:

BLEU Score: 0.6064  
ROUGE-1 F1-score: 0.4574  
ROUGE-2 F1-score: 0.1984  
ROUGE-L F1-score: 0.4402

IndicTrans achieves the highest BLEU score (0.6370) and competitive ROUGE scores for Hindi to Gujarati translation.

ChatGPT and NLLB show similar performance but slightly lower compared to IndicTrans in this translation direction.

### **4. Gujarati to Hindi Translation (gu\_hi):**

NLLB:

BLEU Score: 0.6224  
ROUGE-1 F1-score: 0.5848  
ROUGE-2 F1-score: 0.3556  
ROUGE-L F1-score: 0.5524

IndicTrans:

BLEU Score: 0.6726  
ROUGE-1 F1-score: 0.5947  
ROUGE-2 F1-score: 0.3605  
ROUGE-L F1-score: 0.5580

ChatGPT:

BLEU Score: 0.7000  
ROUGE-1 F1-score: 0.6319  
ROUGE-2 F1-score: 0.4303  
ROUGE-L F1-score: 0.6178

ChatGPT achieves the highest BLEU score (0.7000) and competitive ROUGE scores for Gujarati to Hindi translation.

IndicTrans follows closely with similar performance across BLEU and ROUGE metrics.

NLLB shows slightly lower scores compared to ChatGPT and IndicTrans in this translation direction.

IndicTrans consistently demonstrates strong performance across all translation directions, particularly excelling in Hindi to English translation.

ChatGPT performs competitively, especially in Gujarati to Hindi translation.

NLLB shows decent performance but tends to be slightly behind IndicTrans and ChatGPT in most translation tasks.

## **10. Conclusion**

This comparative analysis of machine translation models—NLLB, IndicTrans, and ChatGPT—across various language pairs highlights important findings for translation quality evaluation. IndicTrans consistently demonstrates superior performance, particularly excelling in Hindi to English translation with the highest BLEU and ROUGE scores. ChatGPT shows competitive performance, especially in Gujarati to Hindi translation, while NLLB performs decently but lags behind in overall translation quality metrics.

Key takeaways include the importance of selecting specialised models like IndicTrans for low-resource languages and leveraging broader language models like ChatGPT for specific translation tasks. The study underscores the significance of objective evaluation metrics such as BLEU and ROUGE in guiding model selection and development. Future efforts should focus on further model fine-tuning, dataset expansion, and integration of advanced techniques to enhance translation quality and address linguistic diversity challenges effectively. This report contributes valuable insights to the field of machine translation, facilitating progress towards more accurate and inclusive multilingual communication technologies.

