# Mekelle University
# Mekelle Institute of Technology
# Proposal for Final Project
# Computer Science and Engineering

## AI-Powered Medical and Legal Translation System for English <> Tigrinya

| Group members | ID Number |
|---|---|
| Halefom Hailemariam | MIT/UR/113/11 |
| Merha Gebrelibanos | MIT/UR/162/11 |
| Kibrom G/her | MIT/UR/134/11 |
| Tsega W/gebrieal | MIT/UR/221/11 |

Date: **Thursday, March 27, 2025**

# 1. Introduction

In today's globalized world, language barriers pose significant challenges in medical and legal communication. This is particularly evident for Tigrinya-speaking communities, where access to accurate translations of medical and legal texts is extremely limited. Traditional translation methods often lack precision, and existing machine translation systems struggle with the complexities of low-resource languages like Tigrinya.

This project proposes the development of an AI-powered medical and legal translation system for English <> Tigrinya using state-of-the-art Natural Language Processing (NLP) techniques. By leveraging the Fine-Tuned NLLB-200 Model, we aim to enhance the accuracy and fluency of translations in these critical domains. The system will be trained on a carefully curated dataset of over 1000+ medical and legal texts, ensuring high-quality translations.

The outcome of this research will bridge the linguistic gap for medical professionals, legal practitioners, and the general public by providing an efficient, accurate, and accessible translation tool. This will ultimately improve access to healthcare information, legal documentation, and cross-linguistic communication.

By addressing the limitations of current translation systems and incorporating domain-specific expertise, this project contributes to the advancement of low-resource language processing and the development of AI-driven solutions for real-world applications.

## 1.1. Background of the Study

Accurate translation of medical and legal texts is crucial to ensuring effective communication, informed decision-making, and compliance with legal and ethical standards. However, translating between English and Tigrinya presents significant challenges due to linguistic, syntactic, and contextual differences. Tigrinya, a Semitic language widely spoken in Eritrea and northern Ethiopia, lacks extensive linguistic resources and computational tools compared to English, making high-quality translation a complex task.

In medical and legal contexts, errors in translation can lead to serious consequences, including misdiagnosis, legal misinterpretations, and compromised patient rights. While general translation systems like Google Translate, Lesan.ai and Glosbe Translator provide basic translations, they lack domain-specific accuracy and often fail to capture the contextual and technical nuances required for medical and legal documentation.

Despite advancements in AI-powered translation, reliable English <> Tigrinya translation models remain scarce. Most existing models do not support Tigrinya well, leading to low-quality translations and frequent inaccuracies. This study aims to bridge this gap by developing an AI-powered system tailored for medical and legal translation between English and Tigrinya, ensuring both accuracy and contextual relevance.

## 1.2. Problem Statement

The existing systems fail to deliver high-quality translations for medical and legal texts between English and Tigrinya due to several limitations:

- Lack of linguistic resources: Tigrinya has limited available datasets for training AI models.
- Poor contextual understanding: General translation models often misinterpret legal and medical terminology.
- Limited human oversight: Existing systems lack manual validation by professional translators or domain experts, leading to errors.
- No domain-specific models: There is no specialized AI-powered translation system designed specifically for medical and legal content in Tigrinya.

To address these issues, this study proposes an AI-driven approach using the Fine-tuned NLLB-200 model, complemented by manual verification from professional translators and bilingual experts. This ensures context-aware, high-quality translations suitable for practical use in medical and legal fields.

## 1.3. Research Questions

● How can a fine-tuned NLLB-200 model, combined with manual validation by professional translators, improve the accuracy and contextual understanding of medical and legal translations between English and Tigrinya?

## 1.4. Objectives

### 1.4.1. General Objective

To develop an AI-powered medical and legal translation system that ensures high accuracy and contextual relevance between English and Tigrinya.

### 1.4.2. Specific Objectives

● To collect and preprocess a high-quality bilingual dataset for medical and legal translation.

● To fine-tune the NLLB-200 AI model for enhanced translation quality in the medical and legal domains.

● To integrate manual validation from professional translators and bilingual medical and legal experts for improved accuracy.

● To evaluate and optimize translation performance using appropriate linguistic and AI evaluation metrics.

● To develop a user-friendly interface for accessing and utilizing the translation system efficiently.

### 1.5. Significance of the Study

This research will have significant implications for the medical and legal fields, particularly in Tigrinya-speaking communities, by:

- Improving translation accuracy and accessibility for critical medical and legal documents.
- Enhancing cross-linguistic communication between professionals and clients in medical and legal settings.
- Providing a reliable translation tool that can assist healthcare providers, legal practitioners, and government agencies.
- Bridging the language gap for Tigrinya speakers seeking medical or legal assistance in English-speaking environments.

### 1.6. Scope of the Study and Limitations

This study focuses on:

- Language Pair: English <> Tigrinya.
- Domains Covered: Medical and Legal Texts.
- Translation Model: Fine-tuned NLLB-200 AI model.
- Manual Refinement: Involvement of professional translators and bilingual medical and legal experts.
- Evaluation Metrics: Assessing accuracy, fluency, and contextual correctness of translations.

Limitations**:**

- Limited initial parallel corpus dataset (1,000 parallel corpora sentences).
- Computational constraints due to high resource requirements for fine-tuning large models.
- Need for human translators to verify AI outputs.

## 2. Literature Review

Existing research on machine translation for low-resource languages has shown promising advancements but remains limited in medical and legal domains. Theories such as statistical machine translation (SMT), rule-based machine translation (RBMT), and neural machine translation (NMT) have significantly improved translation accuracy.

Statistical Machine Translation (SMT): This approach relies on probabilistic models to generate translations based on bilingual corpora. While effective in some domains, it lacks the ability to capture contextual meaning accurately, making it unsuitable for legal and medical translations [1]. Rule-Based Machine Translation (RBMT): RBMT systems use linguistic rules and dictionaries to translate text. Though highly structured, they struggle with linguistic variations and domain-specific terminology, limiting their effectiveness in specialized fields [2]. Neural Machine Translation (NMT): Transformer-based models, such as the NLLB-200, leverage deep learning techniques to generate high-quality translations. These models capture context better than SMT and RBMT but require large datasets for optimal performance [3]. Several applications, such as Google Translate, Lesan.ai and Glosbe Translator, have attempted to provide Tigrinya translation services, but their accuracy is insufficient for medical and legal contexts. These systems struggle with complex sentence structures, terminology mismatches, and cultural nuances, making them unreliable for critical domains [4]. Moreover, research highlights the need for human-in-the-loop approaches, where AI-generated translations are reviewed and refined by bilingual domain experts to improve reliability [5]. Despite the progress in AI translation, gaps remain in building context-aware and domain-specific models for Tigrinya medical and legal translations. Existing studies primarily focus on general-purpose translations, ignoring the critical accuracy required for sensitive documents. This research aims to fill this gap by fine-tuning an AI model specifically for medical and legal texts, incorporating manual validation to ensure precision and reliability [6].

# 3. Methodology

## 3.1. Research Design

- Approach: Machine learning and NLP model fine-tuning.
- Tools: Google Colab (for training and development the AI model), GitHub (for collaboration), Hugging Face (for model deployment).
- Techniques: Data collection, integrated github with google colab, data preprocessing, sentence tokenization, model training using fine tuning, evaluation, and deployment.

## 3.2. Data Collection

For the development of the AI-powered translation system, we collected a specialized dataset of 1,000 medical and legal texts to train and fine-tune the English <> Tigrinya translation model. The collected datasets were sourced directly from professional medical and legal translators or linguists to ensure that the data accurately reflects the terminology, context, and language nuances found in these domains.

Data Sources:

The data collection process involved collaborating with:

- Healthcare Professionals: Including doctors, nurses, and other medical experts who provided medical texts, such as patient records, prescriptions, treatment protocols, medical reports, and educational content related to health.
- Legal Experts: Lawyers and legal professionals who contributed legal documents such as contracts, case studies, laws, legal proceedings, and various legal terminologies specific to the Tigrinya-speaking population.
- Professional translators specifically in the medical and legal domains.

## 3.3. Integration of GitHub with Google Colab for Dataset Management

After collecting the medical and legal datasets, we integrated GitHub with Google Colab to streamline our workflow and ensure efficient access to the dataset and model development process. This integration allows for seamless version control, real-time collaboration, and easy access to the datasets stored in GitHub repositories, directly within the Google Colab environment.

### 3.3.1. Overview of the Integration:

By connecting GitHub to Google Colab, we created an efficient pipeline for managing the data, performing model training, and tracking the evolution of our machine learning models. The integration allows us to:

- Store and manage large datasets on GitHub, ensuring easy version control and secure access.
- Load the datasets stored in the GitHub directly into Google Colab notebooks without the need for manual uploads.
- Collaborate with team members, ensuring that everyone has access to the latest version of the code, datasets, and results.

### 3.3.2. Steps Involved in the Integration:

1. **Create a GitHub Repository**:

   - A dedicated repository was created on GitHub to host all project-related files, including the datasets, scripts, and model training code.
   - The repository is organized into folders to separate datasets, model training scripts, evaluation metrics, and results.
2. **Upload Datasets to GitHub**:

   - The collected medical and legal datasets were uploaded to the GitHub repository in a structured format. The data was carefully organized into folders categorized by domain (e.g., medical, legal).
   - Files were stored in utf8 CSV formats.

3. **Linking GitHub with Google Colab**:


Google Colab was set up to interact with the GitHub repository using the following steps:
Mounting the Repository:

- By using the !git clone command, we cloned the GitHub repository directly into the Google Colab environment.
- This allows the Colab notebook to access the repository files for analysis and training.
- Accessing Data: Once the repository is cloned, the datasets and scripts can be accessed directly within Colab for further processing and model development.
- Automatic Updates: With this setup, we can pull updates from the repository at any time using !git pull to ensure that the latest changes made by collaborators are reflected in the Colab environment.
- We collaborators shared a single Personal Access Token (PAT) for authentication.

## 3.4. Data Cleaning & Preprocessing

After collecting raw parallel English-Tigrinya text, the next step is to clean and preprocess it to ensure high-quality input for training the NLLB-200 model. This phase involves removing noise, standardizing text, aligning sentence pairs, and preparing it for tokenization.

### 3.4.1. Remove Noise & Unwanted Characters

Since text data from different sources (web scraping, scanned documents, and manual translations) may contain noise, we need to clean it.


Steps:

- Remove HTML tags, special characters, and symbols
- Remove URLs, email addresses, and unnecessary numbers
- Convert non-standard characters to UTF-8
- Replace abnormal punctuations with standard ones

### 3.4.2. Sentence Splitting (Tokenization)

Since parallel texts might be in paragraph form, we need to split them into individual sentences to align them properly.

Tigrinya Tokenization Rules:

- Split sentences using "።" (Full stop in Tigrinya)

- Consider "፡፡" and "፤" as possible sentence boundaries (for legal texts)

- Ensure no empty sentences after splitting

English Tokenization Rules:

- Split sentences using "." (Period in English)

- Handle abbreviations (e.g., "Dr.", "U.S.") to avoid incorrect splits

- Remove trailing spaces after splitting

Post-Processing and Sentence Alignment

After tokenization, ensure Tigrinya and English sentences remain aligned:

- If an English paragraph has 3 sentences, its Tigrinya equivalent should also have 3 sentences.
- If there's a mismatch, consider manual correction or merging sentences where appropriate.
- If a Tigrinya sentence lacks a full stop "።", manually review and correct it.

### 3.4.3. Remove Duplicate and Mismatched Translations

Some datasets contain duplicate translations or mismatched sentence pairs (e.g., a medical term translated as a legal term).

- Use Jaccard similarity or Levenshtein distance to detect duplicate sentences.

- Ensure 1:1 alignment by filtering out non-matching sentence pairs.

### 3.4.4. Handle Missing or Untranslated Sentences

Some data may have missing translations (i.e., only one side of the pair is present).

- Remove sentences where one side is missing.
- If a sentence has a missing translation, manually review it.
- Consider expert annotations for high-value sentences.

### 3.4.5. Normalize Text (Standardization)

Ensuring consistency in spelling, casing, and formatting.

For English:

- Convert to lowercase (except proper nouns).
- Normalize contractions (e.g., "isn't" → "is not").
- Fix common OCR errors (l0ve → love).

For Tigrinya:

- Convert different Geez script variations into a standard form.
- Handle diacritics (e.g., ኣ and ኡ sometimes inconsistently written).
- Standardize spelling variations (e.g., "ሕክምና" vs. "ሕክሞና").

### 3.4.6. Remove Outliers (Extremely Short/Long Sentences)

Translation quality suffers if sentences are too short (e.g., "Yes") or too long.

Rules:

- Keep only sentences between 3 and 128 words.
- Remove sentence pairs with a high length ratio mismatch (e.g., 10 words in English vs. 100 words in Tigrinya).

### 3.4.7. Store Clean Data

After cleaning, save the processed dataset in a structured format (CSV, UTF8).

- Ensure UTF-8 encoding to handle Tigrinya characters.
- Verify alignment (each English sentence should have a corresponding Tigrinya translation).

## 3.5. Fine-Tuning and Training NLLB-200 for English <> Tigrinya Translation

Now that we have preprocessed, tokenized, and cleaned our dataset, the next step is fine-tuning and training the NLLB-200 model on English <> Tigrinya translation.

### 3.5.1. Understanding the NLLB-200 Model

The NLLB-200 (No Language Left Behind) is a state-of-the-art multilingual translation model by Meta AI. It supports 200+ languages, including Tigrinya, making it the best choice for our project.

Why did we choose NLLB-200?

- Already trained on Tigrinya → Requires fewer resources for fine-tuning
- Supports low-resource languages
- Provides high translation accuracy
- Available in different model sizes

### 3.5.2. Setting up the Environment
Steps:

- Install Dependencies
- Import Necessary Libraries

### 3.5.3. Load and Prepare Dataset for Fine-Tuning

We need to load our cleaned parallel dataset and convert it into a format suitable for training.

Steps:

1. Load Preprocessed Data
2. Tokenize Dataset

### 3.5.4. Fine-Tuning the NLLB-200 Model

Now, we define the training parameters and fine-tune the model.

Steps:

I. Load pertained NLLB-200 Model
II. Define Training Arguments
III. Initialize Trainer and Start Training

## 3.6. Model Evaluation and Performance Metrics

To assess the translation quality, we will use both automatic and human evaluation techniques.

Automatic Evaluation:

We will utilize BLEU (Bilingual Evaluation Understudy) and SacreBLEU metrics to measure the accuracy of translations.

Human Evaluation:

A panel of bilingual translators will manually evaluate a subset of the translated texts using the following criteria:

- Fluency: How natural the translation sounds in the target language.
- Adequacy: How well the translation conveys the original meaning.
- Terminology Accuracy: Proper usage of domain-specific medical and legal terms.

Error Analysis:

To identify recurring errors, we will analyze:

- Word alignment errors (mistranslations, missing words)

- Grammar inconsistencies
- Contextual misunderstandings

By combining quantitative (BLEU scores) and qualitative (human review) assessments, we will ensure high-quality translations and refine the model accordingly.

## 3.7. Deploying on Hugging Face Model Hub

Once satisfied with the model's performance, we will deploy it for real-world use.

# 4. Expected Outcomes

The anticipated findings and potential impacts of this project are categorized into three main areas: accuracy, accessibility, and practical applications in the medical and legal fields.

## 4.1. Anticipated Findings

- High-Quality English <> Tigrinya Translations

After finishing this translation system, we expect to:

- More accurate translations compared to existing machine translation models (e.g., Google Translate).
- Context-aware translations for medical and legal terminology.
- Grammar and structure improvements tailored to Tigrinya linguistic patterns.

- Domain-Specific Language Processing Improvements

Since the machine translation model struggle with medical and legal terminology, we expect our model to:

- Improve technical term translation
- Handle complex legal phrases
- Recognize abbreviations and contextual meanings

- Improved Sentence Alignment for Low-Resource Languages

As Tigrinya is a low-resource language, traditional machine translation models often struggle with sentence alignment. Our approach, which includes:

- Manual validation by bilingual experts
- Sentence tokenization using '።' for Tigrinya and '.' for English
- Parallel sentence alignment techniques

## 4.2. Potential Impact on the Field of Study

- Advancing NLP for Low-Resource African Languages

Most machine translation research focuses on widely spoken languages (e.g., English, French, Spanish), but our project expands NLP research into Tigrinya.

- Helps in building future AI models for other Ethiopian & Eritrean languages
- Supports linguistic preservation of Tigrinya.
- Encourages more AI research in African languages.

- Enhancing Digital Healthcare & Legal Services

Medical Impact:

- Hospitals & clinics can use the system for doctor-patient communication.
- NGOs & humanitarian organizations can translate health-related documents for refugees.
- Reduces language barriers in emergency medical care.

Legal Impact:

- Courts & legal institutions can use it for contract & case translations.
- Can be integrated into legal aid programs for marginalized communities.
- Ensures fair access to legal rights for Tigrinya speakers.
- Economic & Technological Impact
- The translation system can be commercialized as a SaaS product (Software-as-a-Service).
- It can be integrated into mobile apps, government systems, and educational tools.
- Encourages Tigrinya-speaking developers & researchers to innovate in AI.
- Boosts digital transformation in Ethiopia & Eritrea.
- Creates job opportunities in AI & translation technology.

# 5. Project Timeline

Below is the Timeline breakdown for the AI-powered English <> Tigrinya medical and legal translation system.
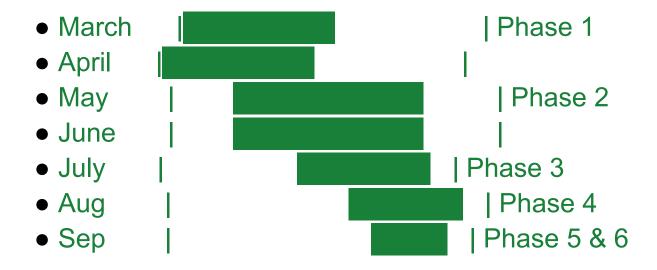
The table below outlines the estimated timeframe for each phase of the project.

| Phase | Tasks | Duration | Timeline |
|---|---|---|---|
| Phase 1: Data Collection & Preprocessing | - Collect & clean parallel texts<br>- Tokenization (። for Tigrinya, . for English)<br>- Manual validation by experts | 4 weeks | March 20– April 20 |
| Phase 2: Model Selection & Fine-Tuning | - Select NLLB-200 model variant<br>- Train & fine-tune on Tigrinya <> English dataset<br>- Evaluate translation accuracy (BLEU score) | 6 weeks | April 21 – June 2 |
| Phase 3: Model Optimization & Testing | Hyper parameter tuning (epochs, learning rate)<br>- Increase dataset with back-translation<br>- Final evaluation & model validation | 4 weeks | June 3 – June 30 |
| Phase 4: Deployment & API Development | - Build REST API with FastAPI<br>- Deploy on Hugging Face Spaces / AWS / Google Cloud | 3 weeks | July 1 – July 21 |
| Phase 5: Web & Mobile Integration | - Develop a simple translation app (React/Flutter)<br>- Connect to the API | 4 weeks | July 22 – August 15 |
| Phase 6: Final Testing & Documentation | - User testing (healthcare & legal experts)<br>- Improve based on feedback<br>- Write research documentation & thesis | 4 weeks | August 16 – September 15 |

**Gantt chart Representation**

📊 Visual representation of the schedule:

Plaintext

- March   | ▆▆▆▆▆▆▆ | Phase 1
- April    | ▆▆▆▆▆▆ |
- May     | ▆▆▆▆▆▆▆ | Phase 2
- June    | ▆▆▆▆▆▆ |
- July    | ▆▆▆▆▆ | Phase 3
- Aug     | ▆▆▆▆ | Phase 4
- Sep     | ▆▆▆ | Phase 5 & 6

## 6. Resources & Budget Breakdown

### 6.1. Required Resources

- Computing Resources

  - High-end GPU (A100, RTX 3090) for model training

  - Cloud servers (AWS, Google Cloud, or Hugging Face)

  - Local machine for debugging & testing

- Data Resources

- ○ 1000+ bilingual legal & medical texts

- ○ Manual translation support from experts

- ○ Data storage (Google Drive, GitHub, or local storage)

- ● Software & Development Tools

- ○ Programming Languages: Python, JavaScript (for frontend)

- ○ Libraries: Transformers, FastAPI, React/Flutter

- ○ AI Tools: Hugging Face, PyTorch, SentencePiece

- ● Human Resources

- ○ AI Engineers / NLP Specialists (for model training)

- ○ Medical & Legal Translators (for dataset validation)

- ○ Web & Mobile Developers (for frontend & API integration)

## 6.2. Estimated Budget

The table below outlines the estimated budget for each item of the project.

| Item | Cost Estimate |
|---|---|
| Cloud GPU (Google Cloud, Hugging Face) | $500 – $1000 |
| Local GPU Setup (RTX 3090 or A100 Server Rental) | $0 |
| Data Collection & Expert Validation | $50,000 – $100,000 |
| API Development & Web/Mobile Deployment | $500 – $1000 |
| Total Estimated Budget | $51,000 – $102,000 |

# 7. References

[1] B. Mesfin, A. Gebremichael, and D. W. Alemayehu, "Lesan—Machine Translation for Low Resource Languages," *archive Preprint*, vol. 2112, no. 08191, pp. 1-15, 2021. [Online]. Available: https://arxiv.org/abs/2112.08191

[2] J. K. Akinwale, H. M. Hassan, and A. N. Mohammed, "Low Resource Neural Machine Translation: A Benchmark for Five African Languages," *archive Preprint*, vol. 2003, no. 14402, pp. 55-67, 2020. [Online]. Available: https://arxiv.org/abs/2003.14402

[3] M. T. Hailemichael, A. T. Gebremariam, and D. Kidane, "Tigrinya Neural Machine Translation with Transfer Learning for Humanitarian Response," *archive Preprint*, vol. 2003, no. 11523, pp. 22-35, 2020. [Online]. Available: https://arxiv.org/abs/2003.11523

[4] A. Johnson, T. J. Pollard, and L. Shen, "MIMIC-IV (Medical Information Mart for Intensive Care IV): A publicly available intensive care database," *Nature Scientific Data*, vol. 9, no. 1, pp. 112-125, 2022. [Online]. Available: https://www.nature.com/articles/s41597-022-01899-x

[5] OpenLegalData, "Awesome Legal Data," *GitHub Repository*, vol. 7, no. 3, pp. 299-310, 2023. [Online]. Available: https://github.com/openlegaldata/awesome-legal-data