



IndoLib: A Natural Language Processing Toolkit for Low-Resource South Asian Languages

Citation

Timalsina, Nitya. 2022. IndoLib: A Natural Language Processing Toolkit for Low-Resource South Asian Languages. Master's thesis, Harvard University Division of Continuing Education.

Permanent link

<https://nrs.harvard.edu/URN-3:HUL.INSTREPOS:37373336>

Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

Share Your Story

The Harvard community has made this article openly available.
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

IndoLib: A Natural Language Processing Toolkit for Low-Resource South Asian
Languages

Nitya Timalsina

A Thesis in the Field of Software Engineering
for the Degree of Master of Liberal Arts in Extension Studies

Harvard University

November 2022

Abstract

Out of 7,151 living languages, 665 languages (9.299%) are spoken by nearly 2 billion people across Southern Asia. Of these, 37.74% (251 languages) are endangered, while the vast majority remain underrepresented in language systems. This thesis presents a new NLP toolkit called IndoLib designed to support natural language processing (NLP) research in South Asian languages, consisting of the Indo-Aryan, Dravidian, and Sino-Tibetan language families, in this case. IndoLib includes four primary components: (i) monolingual and multilingual datasets to expand language modeling and language detection for thirty-one Indic languages, (ii) fine-tuned multilingual models for named entity recognition (NER) and summarization, (iii) a bilingual dataset with Sanskrit-English and English-Sanskrit parallel sentences, and (iv) a fine-tuned machine translation model for two-way translations between Sanskrit and English. The fine-tuned multilingual NER and bilingual translation models outperform current benchmark models upon evaluation. This thesis is intended to aid researchers interested in applying transfer learning to develop or optimize transformer-based models for South Asian languages.

Dedication

After a lifelong search for a cohesive cultural identity, I found myself referred to by similar inquirers as a “third culture kid,” or TCK. TCK is a bittersweet term referring to individuals raised in a culture and environment different from their parents’ native origins. Exposed to a wide variety of cultural experiences, TCKs often struggle to balance the development of unique cultural identities with maintaining the values and traditions of their home cultures (Pollock et al., 2010).

When my family moved to the U.S. in 2002, I quickly adapted to the cultures and customs of the St. Louis suburbs where we found ourselves. My summer breaks were often full of comparative linguistic exercises and memorizing complex Sanskrit grammar, but my life was predominantly ruled by English. I retained my enthusiasm for Nepali cuisine, but my command and pronunciation of the Nepali language began to suffer, dismaying my parents, who taught languages like Nepali and Hindi at the high school and college levels and often cited intricate Indian philosophies that would perplex Descartes. They claimed that my Rs no longer rolled off of my tongue, not to mention my newfound inability to distinctly pronounce variations of the letter S. Addressing these issues meant regaining control of my by-then-Americanized tongue; not a simple task for a 12-year-old. It took several months of daily practice until I could clearly pronounce all 47 Devanagari characters of the Nepali language.

This experience of regaining prior knowledge taught me the importance of retaining important information through attention and well-placed intention. Nearly

two decades later, it inspired me to devote my Master's thesis to enhancing the inclusion of languages like Nepali in natural language processing systems. This thesis is a culmination of the seed that my parents once planted, produced with the motive of maximizing the inclusion and global awareness of my mother tongue, Nepali, as well as several other South Asian languages that may not yet have a spokesperson advocating for them. I am tremendously grateful for the continuous support of my family throughout this process, including Sanskrit grammar coaching sessions by my father. This process has reminded me of the best parts of being a TCK: using one's educational privilege and multicultural experiences to leave a positive mark on the world. As I progress in my research endeavors, I will continue to strive toward that objective.

Acknowledgements

I want to thank my research advisor, Professor Hongming Wang, for her invaluable guidance, support, and reassurance throughout this project. Her feedback was instrumental at various stages of the thesis development process, from narrowing my initial focus areas to optimizing an experimentation process to elicit the most meaningful results possible. I am also grateful to Professor Eric Gieseke for being a reliable sounding board and source of encouragement. Finally, I would like to express my sincere gratitude to my friends and family who supported me throughout my time at Harvard. Steve Merrick, in particular, helped me overcome several periods of uncertainty, starting with the initial conversion of a seemingly endless list of research ideas into a cohesive topic with personal significance. I would not have been able to complete this thesis without all of these incredible people on my side.

This work was also supported by Google’s TPU Research Cloud (TRC), which generously provided TPUs for language model development and optimization.

Contents

Table of Contents	vii
List of Figures	x
List of Tables	xi
List of Code	xii
1 Introduction	1
1.1 Background	2
1.2 Resource Disparities in NLP	3
1.3 State of NLP for South Asian Languages	5
1.4 Research Problem	7
2 Experiment and Computation Setup	9
2.1 Technical Toolkit	9
2.2 Computation	11
2.3 Version Control and Code Availability	11
3 Dataset Development and Optimization	13
3.1 Data Collection	13
3.2 Data Normalization and Cleaning	14

3.3	Downsampling and Upweighting	15
3.4	Direct Utilization of Existing Datasets	16
3.4.1	XLSum	16
3.4.2	Naamapadam	16
3.5	Development of Novel Datasets	17
3.5.1	Monolingual Dataset Development	17
3.5.2	Multilingual Dataset Development	19
3.5.3	Bilingual Dataset Development	20
4	Model Training Pipeline and Evaluation Methods	21
4.1	Training Steps	21
4.2	Data Tokenization	22
4.3	Training Optimization	24
4.3.1	Hyperparameter Tuning	24
4.3.2	Data Collation	25
4.3.3	Data loading	26
4.3.4	Optimizer	26
4.3.5	Training Functions	26
4.3.6	Performance Tracking	27
4.3.7	Fine-Tuning	27
4.3.8	Evaluation Metrics	28
4.4	Task-Specific Evaluation Strategies	30
5	Model Architecture, Selection, and Development	32
5.1	Transfer Learning With Transformer-Based Models	32
5.2	Model Architectures	34
5.2.1	Multilingual BERT	34

5.2.2	ELECTRA	34
5.2.3	XLM-RoBERTa	35
5.2.4	mT5	35
5.2.5	mBART-50	35
5.3	Hyperparameter Optimization	35
5.4	Model Selection	41
5.4.1	Named Entity Recognition	41
5.4.2	Summarization	41
5.4.3	Translation	42
5.5	Model Performance Monitoring	43
6	Model Performance	48
6.1	Model Training	48
6.1.1	Named Entity Recognition	48
6.1.2	Summarization	49
6.1.3	Translation	49
6.2	Model Evaluation and Results	50
6.2.1	Named Entity Recognition	50
6.2.2	Summarization	53
6.2.3	Translation	53
7	Conclusion	56
7.1	Summary	56
7.2	Ethical and Environmental Considerations	57
7.3	Future Work	58
	References	59

List of Figures

1.1	<i>Resource Distribution by Language Type (Joshi et al., 2020)</i>	5
5.1	Optuna Optimization History	38
5.2	Optuna Parallel Coordinate Plot	39
5.3	Optuna Contour Plot	39
5.4	Optuna Slice Plot	40
5.5	Optuna Hyperparameter Importances	40
5.6	<i>Weights & Biases Train Loop Monitoring</i>	44
5.7	<i>Weights & Biases Evaluation Loop Monitoring</i>	45
5.8	<i>Weights & Biases Evaluation Performance Monitoring</i>	46
5.9	<i>Weights & Biases Compute Usage Monitoring</i>	47

List of Tables

1.1	NLP Tasks and Areas of Exploration	7
3.1	Summarization Models Based on XLSum Dataset	16
3.2	Named Entity Recognition Models Based on Naamapadam Dataset	16
3.3	Data Sources for Monolingual Datasets	19
3.4	Novel Datasets for South Asian Language Research	20
4.1	Tokenizer Type by Model	23
4.2	Evaluation Metrics for NLP Tasks	31
5.1	Advantages and Disadvantages of Transfer Learning	33
5.2	Models Evaluated by NLP Task	34
5.3	Performance Comparison of Monolingual NER Models	41
5.4	Performance Comparison of Monolingual Summarization Models	42
5.5	Performance Comparison of Baseline Translation Models	42
6.1	Language-Specific Performance of Fine-Tuned Multilingual IndoLib- NER Model	52
6.2	Language-Specific Performance of Benchmark IndicNER Model	52
6.3	Performance of Fine-Tuned Multilingual Summarization Models	53
6.4	FLoRes Dataset Performance of Fine-Tuned Translation Model Versus Benchmarks	55

List of Code

5.1	Optuna Hyperparameter Training Objective Specification	36
6.1	Evaluation of NER Models	50
6.2	Load FloRes Benchmark Dataset for Two Translation Directions . . .	54
6.3	Compute BLEU Metric for Translation Model	54

Chapter 1: Introduction

In recent years, there has been growing interest in developing technology for understanding and producing human language. As computing power increases, so too do the number of applications that rely on computers to perform linguistic analyses. This trend is especially apparent in Natural Language Processing (NLP), which leverages computational methods to analyze and extract relevant information from written or spoken language. NLP involves an interaction of linguistics, computer science, and artificial intelligence that has become increasingly popular due to its applicability across areas like search engine optimization, question-answering systems, speech recognition, and handwriting recognition.

However, despite the increasing popularity of NLP, most existing NLP systems have been designed and tested in “high-resource” languages like English. A significant portion of current research supports English-centric datasets and models even though most of the world’s population speaks languages other than English. Consequently, very few of these languages are adequately supported by NLP systems. Although some progress has been made toward developing NLP tools for non-English languages, these efforts have predominantly focused on specific domains like machine translation or specialized tasks such as financial transaction processing.

This thesis focuses on developing novel datasets and fine-tuned models to enhance the performance of low-resource South Asian languages for several NLP tasks, including language modeling, language identification, named entity recognition,

summarization, and translation. This chapter defines “low-resource” languages within the context of these objectives and contextualizes the implications for South Asian language families.

1.1. Background

The Sapir-Whorf Hypothesis suggests that a language’s structure affects its speakers’ cognition. Spoken language, as a result, influences how people interact with the world, and individuals’ perceptions are relative to their spoken language (Kihlstrom & Park, 2016). At a fundamental level, language is intricately connected with human identity, perception, interaction, and development. This connection highlights the importance of expanding linguistic support and services to individuals beyond those who understand high-resource languages. Researching under-resourced languages has numerous societal, linguistic, technological, cultural, normative, and cognitive implications that may affect millions, if not billions, of people worldwide (Ruder, 2020).

The broad aim of NLP is to enable computers to comprehend human language. Human language understanding, however, is not a simple task, as genuine comprehension of linguistics requires a robust understanding of complex grammatical structures coupled with the nuances of syntax, semantics, pragmatics, phonology, and morphology. These characteristics are difficult to synthesize into machine-understandable representations, and matters are further complicated by the derivation of datasets from various inputs, ranging from raw text to images and speech. This contributes to an inherent inequality for languages without extensive corpora, written or otherwise. Furthermore, the sheer diversity of languages spoken or written by humans expands the challenges associated with developing inclusive NLP systems (Joshi et al., 2020).

Due to the increasing availability of computing resources, an improved understanding of algorithms, continuous advancements in machine learning techniques, and

the broader availability of annotated data and pre-trained models, the field of NLP has progressed immensely over the last several decades. The global NLP community’s growing acknowledgment of resource distribution gaps has also elicited significant efforts to accommodate under-resourced languages. Although resource disparities and opportunities for improvement remain for the vast majority of languages, considerable progress has been achieved in establishing the tools, benchmark models, and datasets necessary to develop and optimize NLP systems for low-resource and zero-resource languages. For instance, Masakhane, a grassroots organization led by NLP researchers, has published journal articles on developing novel datasets, models, and machine translation methods for various low-resource African languages (Orife et al., 2020). Based on the linguistic diversity represented by South Asian languages and the number of resources recently released, similar opportunities to contribute to the progression of NLP research in other low-resource languages can be anticipated.

As a prominent NLP researcher and the author of numerous academic papers on state-of-the-art multilingual and cross-lingual benchmarks, Sebastian Ruder urges young researchers to delve into this research area. “Given the availability of data and models in different languages,” Ruder states, the “stage is set” for “meaningful progress on languages beyond English” (Ruder, 2021). Researchers across academia and industry have expressed increasing interest in NLP systems that accommodate multilingual communication and cross-cultural understanding. The “stage” is indeed “set” to address existing disparities to build more inclusive, globally expansive digital language ecosystems.

1.2. Resource Disparities in NLP

There are 7,151 active languages distributed across 142 different language families around the world (Rowe & Levine, 2018). However, only around 100 languages

have moderate to high levels of natural language processing (NLP) resources available, and only 20 languages are considered high-resource (Maxwell & Hughes, 2006). The vast majority of the world’s languages lack the resources required to train usable NLP systems, thereby widening the vast technological divide affecting most of the world’s citizens (Maryatt, 2018). Furthermore, most under-resourced languages have historically lacked large labeled corpora needed to develop functional NLP systems, despite a growing global interest in expanding linguistic inclusion (Le & Besacier, 2009). The problem is exacerbated by difficulties in collecting necessary datasets and engaging researchers in developing NLP models for low-resource, zero-resource, and unseen languages (Baumann & Pierrehumbert, 2014).

As a result, the majority of progress made in NLP has been devoted to high-resource languages such as English, thereby widening the “digital language divide” between dominant languages (primarily Western) and others (Ruder, 2020). Despite the number of people who rely on them, relatively minimal research efforts have been undertaken for low-resource and zero-resource NLP development. Only a few hundred languages are even represented on the Internet, and the speakers of under-resourced and underrepresented languages have limited digital access to information and services in their native languages (Pimienta et al., 2009).

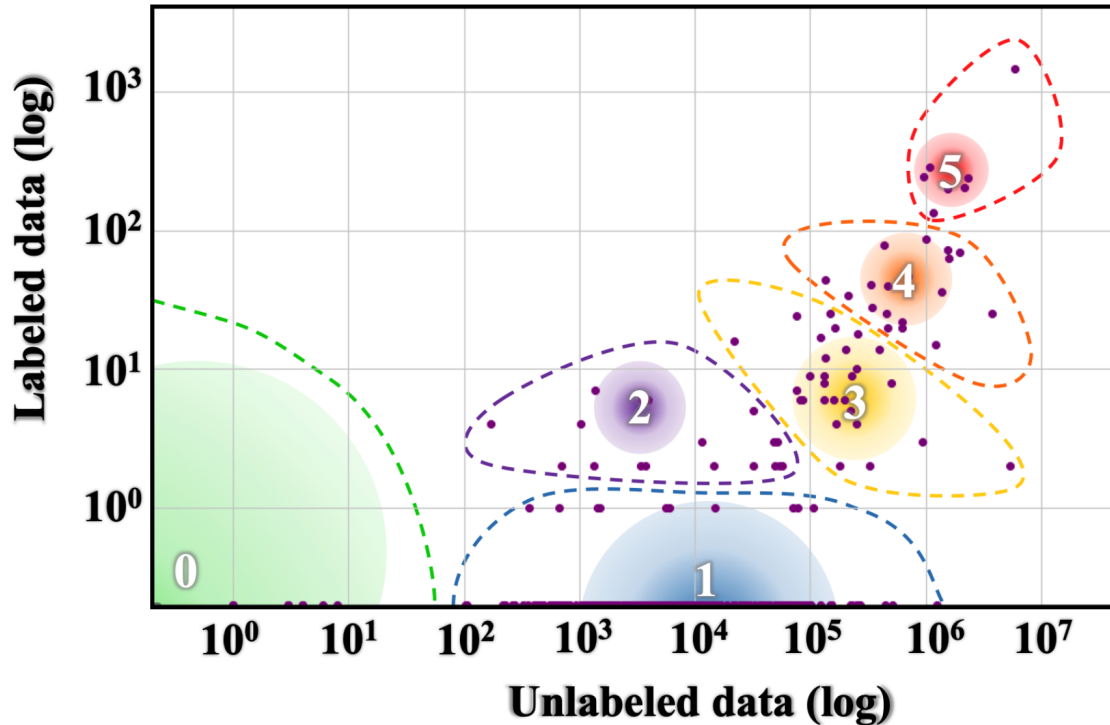


Figure 1.1: *Resource Distribution by Language Type* (Joshi et al., 2020)

In Figure 1.1 above, the size and color of each circle represent the “number of languages” and “speakers” by category. Colors represent the total speaker population size, ranging from low (i.e., violet) to high (i.e., red), using the VIBGYOR (Violet-Indigo-Blue-Green-Yellow-Orange-Red) spectrum. The availability of labeled data is disproportionate to the number of speakers, such that a minority of languages receive the vast majority of resources available (Joshi et al., 2020).

1.3. State of NLP for South Asian Languages

Of 7,151 living languages, 665 (9.299%) are spoken by 1,913,790,000 billion people across Southern Asia. While 68 of these languages are considered institutional, 170 are labeled as “developing,” 176 are “vigorous,” 210 are “in trouble,” and 41 are

“dying” (Eberhard et al., 2022).

Many of these languages share similar historical roots, cultural underpinnings, and structural relationships, and a minimum of 20 actively spoken South Asian languages share a significant portion of their vocabulary with their primary root language, Sanskrit. Additional similarities include but are not limited to their syntactical structure, morphology, and phonology (Sengupta & Saha, 2015).

Considering the diversity and widespread usage of South Asian languages, there is a pressing need to develop NLP applications and resources to overcome low-resource barriers. Novel opportunities to expand machine translation to under-resourced communities exist, and numerous corporations, governmental agencies, and research organizations have made significant efforts to release datasets and state-of-the-art compatible resources to further South Asian language research. This includes the development of IndicNLPSuite, as well as a variety of cutting-edge resources.

Furthermore, there has been renewed interest in exploring inherent linguistic relationships among South Asian languages, including links to root languages like Sanskrit, thereby paving the way toward multilingual, cross-lingual, and ultimately language-agnostic NLP systems (Salian, 2019). There are unprecedented opportunities to contribute to developing NLP models and resources for under-resourced South Asian languages, making this a promising area of research.

1.4. Research Problem

In an effort to enhance the representation of South Asian languages, this thesis incorporates several NLP tasks: language modeling, language identification, named entity recognition (NER), summarization, and machine translation. Chapter 2 includes a complete description of these tasks, and Table 1.1 below lists specific research questions associated with each task:

Natural Language Processing Task	Areas of Exploration
Language Modeling	<ol style="list-style-type: none">1. Do existing datasets used to train benchmark models like Multilingual BERT accurately represent a diverse set of South Asian languages?2. Are specific languages overrepresented or underrepresented compared to others?
Language Identification	<ol style="list-style-type: none">3. Do language models correctly identify and differentiate between different South Asian languages?4. Is a specific model architecture better suited for named entity recognition tasks for South Asian languages?
Named Entity Recognition (NER)	<ol style="list-style-type: none">5. How does a monolingual model trained on the full dataset compare to the multilingual model trained on evenly distributed and shuffled splits of various South Asian languages?6. Does fine-tuning a pre-trained model on monolingual corpora enhance summarization quality?
Summarization	<ol style="list-style-type: none">7. Is a specific model architecture better suited for summarization tasks for South Asian languages?8. Is a specific model architecture better suited for translating a South Asian language to English?
Translation	<ol style="list-style-type: none">9. How does a bilingual model trained on custom corpora compare to benchmark translation models?

Table 1.1: NLP Tasks and Areas of Exploration

These tasks are addressed across seven chapters in this thesis:

- (i) *Introduction* gives an overview of the problem of resource scarcity in NLP and describes the importance of addressing NLP for low-resource languages spoken in South Asia.
- (ii) *Experiment and Computation Setup* delves into the technical requirements of developing datasets and training fine-tuned models for several NLP tasks.
- (iii) *Dataset Development and Optimization* outlines the process of creating datasets for language modeling and language identification tasks.
- (iv) *Model Training Pipeline and Evaluation Methods* discusses how to train models for named entity recognition, summarization, and machine translation. It also presents evaluation metrics that can be used to evaluate model performance and compare fine-tuned models to existing benchmarks.
- (v) *Model Architecture, Selection, and Development* explores the problem of choosing a suitable architecture and developing hyperparameter-tuned models for named entity recognition, summarization, and translation models.
- (vi) *Model Performance* compares the performance of the models developed in the thesis with state-of-the-art benchmarks.
- (vii) *Conclusion* summarizes the results from all experiments conducted throughout the thesis and describes avenues for future work.
- (viii) *References* includes all sources utilized in the thesis.

Chapter 2: Experiment and Computation Setup

The following sections describe the technical and data-specific requirements for developing the datasets covered in Chapter 3 and building the models described in Chapter 5.

2.1. Technical Toolkit

Building a functional, reproducible, and expandable NLP system requires utilizing a robust set of software tools and maximizing the simplicity of the training and evaluation pipelines. Models and datasets should be easy to use, quick to deploy and test, and have minimal dependencies. Furthermore, training methods and pipelines should be replicable, with clear hardware and software requirements.

In this thesis project, two machine learning frameworks were utilized for training and experimentation: PyTorch and JAX. PyTorch was leveraged for CPU-based dataset development and GPU-based training of customized training pipelines (Imambi et al., 2021). JAX, on the other hand, was specifically utilized for TPU-based language modeling use-cases (Bradbury et al., 2018). These frameworks were selected based on user familiarity, ease of use, and hardware-specific performance. In addition, the Python language was selected due to the overwhelming number of Python-based frameworks, libraries, and methods available for NLP.

Along with these frameworks, various libraries and tools were also utilized:

- (1) *Hugging Face* provides easily integrable datasets, pre-trained models, and APIs for various NLP use cases (Yue, 2022). It offers simple access to several libraries, including:
 - (a) *Transformers*, which provides APIs to download and train pre-trained models seamlessly (Wolf et al., 2020).
 - (b) *Datasets*, which allows users to build and manipulate custom datasets (Lhoest et al., 2021).
 - (c) *Accelerate*, which simplifies the process of running models on single and distributed hardware configurations, including multiple GPUs or TPUs, without requiring significant code modification (Gugger et al., 2022).
- (2) *NLTK* offers utilities for tasks like text processing, tokenization, tagging, and parsing (Bird, 2006). Similarly, *iNLTK* provides NLTK functionalities specifically customized for various South Asian languages (Arora, 2020).
- (3) *Pandas*, *Numpy*, *Matplotlib*, *Seaborn*, and *Scikit-learn* are standard Python libraries utilized for computation, data manipulation, visualization, and analysis (McKinney et al., 2011; Harris et al., 2020; Bisong, 2019; Pedregosa et al., 2011).
- (4) *Weights & Biases* is a platform that logs, monitors, and compares experiment model performance (Morris, 2022).
- (5) *Beautiful Soup* parses HTML/XML documents into usable objects, such as lists, dictionaries, and strings (Patel, 2020).
- (6) *Indic NLP Library* provides various tools developed explicitly for South Asian languages (Kunchukuttan, 2020).

- (7) *Flax* is a neural network library designed to optimize JAX models with CPU and GPU acceleration support (Heek et al., 2020).
- (8) *Optuna* is an open-source framework that enables researchers to find optimal hyperparameter settings for their machine learning models (Akiba et al., 2019).
- (9) *Jupyter Notebook* is a web application used to create and share interactive notebooks in Python (Kluyver et al., 2016).

2.2. Computation

NLP research is generally computation-intensive, often involving millions to billions of data points. The problem domain is frequently so large that it requires specialized hardware or software for efficient processing. Even when fine-tuning pre-trained models, computation costs and time requirements can be substantial.

Experimentation was primarily conducted on cloud-based Jupyter Notebook instances. Google’s TPU Research Cloud (TRC) program generously provided TPU credits to minimize the overhead costs of model training and optimization. To optimize resource utilization, code runs were also tested on various CPU, GPU, and TPU instances, with a focus on minimizing overall training time, costs, and memory allocation.

2.3. Version Control and Code Availability

Version control is an essential aspect of software trackability and transparent code management. This is especially important in low-resource NLP, which hinges on collaboration among researchers and developers. To that end, all datasets and finalized models developed in this thesis work are available on Hugging Face (upon

request). In addition, Python scripts and Jupyter Notebooks for each covered task and dataset are available in a *GitHub repository*.

The repository also contains all code used to develop, fine-tune, train, and evaluate the model variants described in this thesis. The Hugging Face and GitHub repositories will be updated with subsequent versions as they are released.

Chapter 3: Dataset Development and Optimization

A corpus is a collection of texts annotated with metadata or other information such as part-of-speech tags, named entities, and lemmas. Corpora can be divided into two broad categories: labeled and unlabeled. Unlabeled corpora contain only text without any annotation. On the other hand, labeled corpora have some additional information associated with each sentence. While part-of-speech tagging is the most commonly-used annotation technique, annotation schemes may also involve objectives like word sense disambiguation, morphosyntactic analysis, and semantic role labeling.

The creation of corpora is a time-consuming process involving significant human effort. There has been extensive research on developing high-quality labeled corpora. However, creating new corpora is still labor-intensive and requires substantial manual intervention. This thesis aims to simplify the development of novel datasets for various South Asian languages through a customized data normalization and cleaning pipeline.

3.1. Data Collection

The Hugging Face Datasets library was instrumental in downloading, processing, and uploading datasets through the Hugging Face Hub. A list of all datasets used in this thesis can be found in Chapter 2.

The first step in data collection is downloading the dataset. The Hugging Face

Datasets Library allows users to download datasets through a simple ‘load_dataset’ function. First, a user specifies the location of the dataset, specific splits (typically train, test, and validation), and an authentication token if needed. Then, the downloaded dataset can be assigned to a variable, downloaded to a local disk, and so on.

In the case of low-resource languages with sparse data and the unavailability of monolingual datasets, multilingual datasets were downloaded for the specific language (or language pair) and filtered to extract monolingual data. For instance, the XLSum corpus includes article-summary pairs for 45 languages. The Marathi subset of the XLSum dataset was downloaded to extract the Marathi content. Then, articles and summaries for Marathi were individually extracted and combined with other corpora obtained from sources like BLOOM, OSCAR, and Wikipedia.

3.2. Data Normalization and Cleaning

Data pre-processing is an essential step in any NLP task, involving steps like cleaning raw text, filtering out noise, normalizing values, and converting strings to numbers. This section describes how each of the datasets described in Chapter 2 was processed.

After collecting the dataset, the first step is to clean it up. In many cases, there may be some issues with the quality or format of the collected data. For example, if the dataset consists of multiple files, it may contain different formats. If the file is corrupted, the entire dataset must be re-downloaded. Similarly, if the data does not match the expected structure, e.g., if it contains URLs instead of plain text, it should be cleaned before proceeding further.

Datasets like XL-Sum and Bloom are pre-processed but may require language-specific division and sentence splitting. Datasets like OSCAR and CC-100 have un-

dergone some preprocessing but require additional steps to remove unnecessary punctuation marks, special characters, URLs, HTML characters, and empty lines. Finally, the text is combined and divided into (roughly) even sentences to be processed.

This thesis used Beautiful Soup, Python ReGex, NLTK, iNLTK, JSON, and Pandas to complete various data-preprocessing tasks. The entire data pre-processing pipeline is available on GitHub, and the processed datasets are available on Hugging Face.

3.3. Downsampling and Upweighting

Downsampling refers to reducing a training set's size by removing samples. Downsampling can help prevent overfitting by limiting the number of parameters in the model. However, downsampled data will not always generalize well across domains because it lacks diversity.

Upweighting refers to increasing the importance of specific samples during training. For example, upweighting may increase the weight given to rarer examples that occur less often in the training set. As a result, the model learns to pay more attention to these examples and becomes better able to predict them.

This thesis incorporates downsampling and upweighting techniques to create diverse datasets with higher coverage of vocabulary and syntactic structures. Downsampling techniques utilized include random sampling, sampling weighted by frequency, and sampling weighted by length. Upweighting techniques included are uniform sampling and sample selection.

In effect, the vocabulary size of the datasets is reduced, and sentences are randomly shuffled to minimize bias or overfitting, particularly for languages with noisy or minuscule datasets. These methods facilitate the development of more balanced datasets that provide competitive results on the trained models.

3.4. Direct Utilization of Existing Datasets

Existing datasets have been used extensively throughout this thesis. The following table provides a brief overview of each dataset, including its origin, size, characteristics, and how they were utilized.

3.4.1 XLSum

XL-Sum is a dataset created for text summarization by the BUET CSE NLP Group. It was utilized in this thesis to create the following monolingual and multilingual summarization models, following extraction of South Asian language components from the full XL-Sum dataset (Hasan et al., 2021).

Model Type	Languages Covered
Monolingual Summarization	Sinhala
Multilingual Summarization	Bengali, Gujarati, Hindi, Marathi, Nepali, Punjabi, Sinhala, Tamil, Telugu, Urdu

Table 3.1: Summarization Models Based on XLSum Dataset

3.4.2 Naamapadam

Naamapadam is a dataset developed by AI4Bharat for named entity recognition (NER) tasks. It was utilized in this thesis to create the following monolingual and multilingual NER models (Mhaske et al., 2022).

Model Type	Languages Covered
Monolingual NER	Oriya
Multilingual NER	Assamese, Bengali, Gujarati, Kannada, Hindi, Malayalam, Marathi, Oriya, Punjabi, Tamil, Telugu

Table 3.2: Named Entity Recognition Models Based on Naamapadam Dataset

3.5. Development of Novel Datasets

Labeled corpora is an essential aspect of training and evaluating language models. This thesis work incorporated existent corpora for named-entity recognition (Naamapadam) and summarization (XLSum) and developed novel corpora for translation (Sanskrit to English). In addition, three types of novel datasets were created for future language modeling and language identification objectives: monolingual, multilingual, and bilingual. This section outlines the datasets developed, along with their respective sources and pre-processing methods.

3.5.1 Monolingual Dataset Development

Monolingual datasets represent a single language and consist of unlabeled text samples. Although monolingual datasets do not require labels, they provide a foundation for building multilingual datasets. In this thesis, monolingual datasets were established for 31 languages through a combination of data sources:

- (1) *Wikipedia*, which offers corpora based on Wikipedia pages in 329 languages (Tunstall et al., 2022).
- (2) *Oscar Corpus*, which includes filtered data for 166 languages from the Common Crawl corpus (Suárez et al., 2020).
- (3) *CC-100*, which includes monolingual data for 100+ languages (Conneau et al., 2019).
- (4) *Itihasa*, which includes 93,000 Sanskrit verses with corresponding English translations (Aralikatte et al., 2021).
- (5) *Large-Scale Nepali Corpus*, which includes over 6.5 million sentences (Lamsal, 2020).

- (6) *The Bloom Library Dataset*, which includes stories from 364 languages across 31 language families (Whitenack, 2022).
- (7) *XL-Sum*, which includes 1.35 million article-summary pairs across 45 languages (Hasan et al., 2021).

Data collected from these sources were divided by language, filtered, cleaned, and normalized using a customized pipeline. A thorough data-processing pipeline resulted in monolingual datasets across 31 South Asian languages, including several previously unrepresented languages listed in the table below.

Language	Language Family	Dataset Sources
Sanskrit	Indo-Aryan	Oscar, Wikipedia, CC-100, Itihasa, Bloom
Nepali	Indo-Aryan	Oscar, Wikipedia, CC-100, Large Scale Nepali Corpus, Bloom, XLSum
Newari	Sino-Tibetan	Oscar, Wikipedia, Bloom
Marathi	Indo-Aryan	Oscar, Wikipedia, Bloom, XLSum
Western Punjabi	Indo-Aryan	Oscar, Wikipedia
Punjabi	Indo-Aryan	Oscar, Wikipedia, XLSum
Gujarati	Indo-Aryan	Oscar, Wikipedia, XLSum
Bishnupriya Manipuri	Indo-Aryan	Oscar, Wikipedia
Sinhala	Indo-Aryan	Oscar, Wikipedia, XLSum
Oriya	Indo-Aryan	Oscar, Wikipedia
Sindhi	Indo-Aryan	Oscar, Wikipedia
Maithili	Indo-Aryan	Oscar, Wikipedia
Assamese	Indo-Aryan	Oscar, Wikipedia
Fiji Hindi	Indo-Aryan	Wikipedia
Bhojpuri	Indo-Aryan	Oscar, Wikipedia, Bloom
Goan Konkani	Indo-Aryan	Oscar, Wikipedia
Doteli	Indo-Aryan	Wikipedia, Bloom
Pali	Indo-Aryan	Wikipedia

Awadhi	Indo-Aryan	Wikipedia, Bloom
Tharu (Chitwania and Dangaura variants)	Indo-Aryan	Bloom
Kashmiri	Indo-Aryan	Wikipedia
Romani	Indo-Aryan	Wikipedia
Hindi	Indo-Aryan	Oscar, Bloom, XLSum
Kannada	Dravidian	Oscar, Bloom
Tamil	Dravidian	Oscar, Bloom
Western Tamang	Sino-Tibetan	Oscar, Bloom
Malayalam	Dravidian	Oscar, Bloom
Bengali	Indo-Aryan	XLSum, Bloom
Athpariya	Sino-Tibetan	Bloom
Lohorung	Sino-Tibetan	Bloom
Telugu	Sino-Tibetan	Oscar, XLSum

Table 3.3: Data Sources for Monolingual Datasets

3.5.2 Multilingual Dataset Development

Multilingual datasets contain labeled texts from different languages. They can be either parallel or non-parallel. Parallel datasets contain equivalent translations of the same sentences in two or more languages, whereas non-parallel datasets contain texts from different languages without translation.

This thesis utilizes existing multilingual datasets for named-entity recognition (Naamapadam) and summarization (XL-Sum), as described in Section 3.4 above. In an effort to maximize the inclusion of low-resource South Asian languages, two novel multilingual datasets were created:

Dataset Type	Relevant NLP Tasks	Number of Languages
Multilingual Unlabeled	Language Modeling, Refinement for downstream tasks	31
Multilingual Labeled	Language Identification, Domain adaptation	31

Table 3.4: Novel Datasets for South Asian Language Research

3.5.3 Bilingual Dataset Development

Bilingual datasets combine two languages that share some common vocabulary. The use of bilingual datasets has become increasingly popular due to their ability to improve the quality of machine learning models. However, parallel corpora are limited in number, especially for low-resource languages, and available datasets have not yet been extensively tested. This thesis involved creating a bilingual dataset for Sanskrit to English translations, derived from:

- (1) *Itihasa*, which contains Sanskrit-English parallel translation pairs (Aralikatte et al., 2021).
- (2) The English-Sanskrit and Sanskrit-English subsets of the *No Language Left Behind (NLLB-200)* dataset includes 148 English-centric and 1465 non-English-centric language pairs (Heffernan et al., 2022; Costa-jussà et al., 2022).

All novel datasets developed in this thesis are intended to be used by other researchers seeking to expand NLP systems for various South Asian languages.

Following the development of novel datasets and customization of pre-existing datasets for the languages covered in this thesis project, a customized training pipeline was developed for three NLP tasks: named entity recognition (NER), summarization, and translation. The following chapter covers the development of a training pipeline, along with appropriate evaluation metrics for each task.

Chapter 4: Model Training Pipeline and Evaluation

Methods

Each NLP task requires specific training, evaluation, and testing pipelines. In this thesis, three different pipelines are developed for named entity recognition, summarization, and translation. Each pipeline consists of four primary components: preprocessing, feature extraction, model training, and post-processing. These components and additional steps are described in this chapter.

4.1. Training Steps

Preprocessing is the first step in any machine learning project and involves converting raw textual data into a format suitable for use within the neural network architecture. Preprocessing may involve adding necessary labels, filtering irrelevant information and extraneous characters, tokenizing text into individual words, removing stop words, converting all characters to lowercase, or standardizing the dataset format. For example, in the case of language identification, a preprocessing step might involve removing HTML characters from all text samples and attaching language-specific labels to a multilingual dataset.

Feature extraction may also be utilized to transform raw textual data into a form that can be fed directly into the neural network. Some features can be extracted using simple statistical calculations, while others may require more complex trans-

formations. Commonly used statistical features include n-grams, unigrams, bigrams, and trigrams. More advanced features often combine multiple statistics from different parts of the document. For example, one common technique for extracting named entities involves counting the number of times each word appears in the context of other words, followed by calculating the ratio between these counts and the total number of instances where the word appeared.

Training natural language models involves creating an algorithm capable of learning patterns in the data and synthesizing them to perform the desired task. A summarization model, for instance, accepts a document or set of documents as input and outputs a summary. Models may be trained using supervised or unsupervised methods. Supervised training pairs known examples of correct output (e.g., summaries) with their corresponding inputs (e.g., documents). Unsupervised training is used when no such pairs exist. Semi-supervised training is a hybrid approach where some labeled data is available but not enough to train a complete model.

Post-processing is the final stage of the pipeline and consists of methods for improving the quality of the output produced by the model. For example, in named entity recognition, post-processing might involve re-ranking the results according to confidence scores or removing any false positives. Post-processing also includes steps like ranking predictions, generating reports, and visualizing results.

4.2. Data Tokenization

Tokenization involves breaking down the text into smaller units called tokens, which typically consist of words, phrases, or text chunks representing meaningful concepts. Tokens are typically identified by whitespace boundaries, such as spaces, tabs, line breaks, and newlines. Tokenization is performed to convert raw text data into a format that enables training machine learning algorithms on large amounts

of unstructured text. Among many types of tokenizers, three options are primarily used: Byte-Pair Encoding (BPE), WordPiece, and SentencePiece.

Byte-Pair Encoding (BPE) is a method that breaks the input text into substrings based on byte sequences. Each substring represents one word or character in the original sentence. The number of unique words/characters can be reduced by merging similar sub-sequences (Sennrich et al., 2015).

WordPiece is an alternative tokenizer that uses a vocabulary of subword units instead of single characters. It first splits the input text into subwords using a predefined dictionary. Then it replaces each subword with its most frequent unit from the dictionary. Finally, it merges all these pieces back into the original string. This approach has demonstrated improved performance for some NLP tasks, such as named entity recognition (Wu et al., 2016).

SentencePiece is another type of tokenizer that applies a hierarchical segmentation algorithm to break the input text into sentences. It then segments each sentence into tokens using a predefined dictionary of subwords. Finally, it merges all these pieces back into a single string (Kudo & Richardson, 2018).

The Transformer models covered in this thesis utilize two types of tokenization: WordPiece and SentencePiece. Table 4.1 below lists the models used in this thesis and their corresponding tokenization types.

Pre-Trained Model	Tokenizer Type
Multilingual BERT (mBERT)	WordPiece
ELECTRA	WordPiece
XLM-RoBERTa	SentencePiece
mBART-50	SentencePiece

Table 4.1: Tokenizer Type by Model

4.3. Training Optimization

The effective training of a model requires optimization of hyperparameters and data processing methods. Furthermore, tracking variations of all models, datasets, and code helps maximize reproducibility and rapid prototyping. In addition to the NLP pipeline listed above, several minor steps are involved in building a complete NLP application, including: (i) hyperparameter tuning, (ii) data collation, (iii) data loading, (iv) optimization, (v) training function, and (vi) performance tracking.

4.3.1 Hyperparameter Tuning

The first step in building a successful NLP application is selecting the right set of parameters. Hyperparameters control various aspects of the model, including the size and type of hidden layers, the amount of dropout applied to each layer, the number of epochs to train over, and many others. Optuna is a hyperparameter optimization framework that vastly simplifies the process of hyperparameter tuning. It allows users to easily create models, evaluate them against a validation dataset, and tune their hyperparameters based on an objective function that evaluates the model's performance across various combinations of hyperparameters.

Specific parameters tuned in this thesis include learning rate, weight decay, and epoch:

- (a) *Learning rates* determine how quickly or slowly the weights of neural networks change during training. A high learning rate enables a network to learn quickly, while a lower learning rate extends convergence time. Higher learning rates generally allow for faster convergence but can lead to unstable behavior. Conversely, lower learning rates require more iterations to converge but may be less prone to oscillation.

- (b) *Weight decay* is a regularization technique that prevents the weights from growing too large and inhibiting the network from becoming overly sensitive to small changes in input data. Weight decay works by penalizing large values of the weights with a penalty term. This means that the larger the weight value, the greater the penalty.
- (c) *Epochs* refer to the number of times through the entire dataset that the algorithm should run. A typical approach would be to start with a low number of epochs (e.g., 10) and increase it until the model achieves its best accuracy.

Before training each model, Optuna was used to build, train, and evaluate deep learning models using gradient-based optimizers. It provides a simple interface for creating models, evaluating them against a validation dataset, and tuning hyperparameters based on an objective function that evaluates model performance across different combinations of hyperparameters (i.e., the hyperparameter search space).

4.3.2 Data Collation

Data collators are objects that form batches of dataset elements and apply pre-processing steps like padding where necessary. Model-specific data collators are available through the Hugging Face Transformers library. For instance, `DataCollatorForSeq2Seq` dynamically pads inputs and labels as required by an encoder-decoder model like `mT5`, which requires shifting labels to the right by one during decoding. This is required to limit the decoder to seeing past ground truth labels rather than current or future labels that it could memorize. Using the correct data collator suited for a given model and task is critical for effective training.

4.3.3 Data loading

Training pipelines often require access to large amounts of raw text data. In order to make efficient use of memory, it is essential to load only the portions of the dataset that are needed at any given time. `DataLoader` classes provide a convenient way to do so. They enable the specification of the portion of a dataset that should be loaded into memory, along with other information like the total number of samples in the dataset, the batch size, and the maximum sample length.

4.3.4 Optimizer

An optimizer class is responsible for adjusting the weights and biases of a neural network model. Common optimizers used in deep learning include gradient descent, Adam, RMSProp, Adadelta, and RMSprop. Optimizers are typically supplied as part of a larger optimization package but can also be independently implemented. In this thesis project, the AdamW optimizer is utilized, as it implements the Adam algorithm with weight decay. An optimizer makes two primary contributions to overall model performance: reducing the variance of gradients and making sure the loss decreases monotonically. By default, the optimizer will automatically choose the best learning rate for each parameter group. However, for more fine-grained control over learning rates, the learning rate can be manually configured using the `learning_rate` argument.

4.3.5 Training Functions

Once all necessary components have been assembled, we must write a training function to iterate over the entire dataset and update the model's parameters accordingly. Training functions are written in Python using Keras, Tensorflow, PyTorch, or

another appropriate language. They take care of all the boilerplate work that would otherwise slow down development.

There are two main types of training functions: evaluation functions and prediction functions. Evaluation functions are used to measure the performance of a model during training. The most commonly used evaluation metric is accuracy, which measures the percentage of correctly classified instances. Prediction functions, on the other hand, use a trained model to predict outcomes for previously unseen examples. These predictions can be binary (true/false) or categorical (one of several possible values).

4.3.6 Performance Tracking

Tracking is a powerful tool for debugging neural networks because it provides clear visualizations of the contribution of each component to the overall loss value. Tracking also makes it easy to compare different versions of the same network. With automatic checkpointing and performance logging, Weights & Biases makes it easy to track model training sessions and GPU utilization.

4.3.7 Fine-Tuning

After training a model, there may still be room for further improvement. Fine-tuning is a technique that trains a new model using a pre-trained model as a starting point. It can improve the performance of the original model by learning from the data it has already seen and making better use of its existing knowledge.

In this thesis, pre-trained models are loaded from the Hugging Face Hub, fine-tuned on custom data and optimized hyperparameters, and evaluated. Additional fine-tuning techniques may further improve the performance of these models.

4.3.8 Evaluation Metrics

A wide variety of metrics are available for evaluating NLP models. This thesis incorporates the loss, accuracy, precision, recall, F1, BLEU, and ROUGE metrics:

A. Loss: The loss metric measures the difference between the predicted and actual outputs. This thesis uses cross entropy as the loss function for all models trained. In binary classification, binary cross-entropy is calculated as:

$$-(y \log(p) + (1 - y) \log(1 - p))$$

In multi-class classification (i.e. $M > 2$), a separate loss is calculated for each class label per observation and the result is summed (Zhang & Sabuncu, 2018). The formula for calculating loss is:

$$-\sum_{c=1}^M y_{o,c} \log(p_{o,c})$$

B. Accuracy: The accuracy metric measures the percentage of correctly classified examples in the test set, using true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). The higher the score, the better the model is at classifying data (Pedregosa et al., 2011). The formula for calculating accuracy is:

$$\frac{TP + TN}{TP + TN + FP + FN}$$

C. Precision: The precision metric measures the proportion of correct predictions among all predictions made by the model. Precision is calculated by dividing the number of correct predictions by the total number of predictions made by the model on the test set (Powers, 2020). The formula for calculating precision is:

$$\frac{TP}{TP + FP}$$

D. Recall: The recall metric measures the proportion of positive examples that are correctly identified. Recall is calculated by dividing the true positives by the total number of positive examples in the test set (true positives plus false negatives) (Davis & Goadrich, 2006). The formula for calculating recall is:

$$\frac{TP}{TP + FN}$$

E. F1 Score: The F1 score is a metric used to evaluate the performance of machine learning models on text classification tasks. It combines Precision and Recall metrics into one single number that takes values between 0 and 1. A value of 0 indicates no predictive power, while a value of 1 means perfect prediction (Goutte & Gaussier, 2005). The formula for calculating the F1 score is:

$$\frac{2 * Precision * Recall}{Precision + Recall} = \frac{2 * TP}{2 * TP + FP + FN}$$

F. BLEU: The BLEU (BiLingual Evaluation Understudy) metric calculates the similarity between two sequences based on their word overlap and penalizes short machine translations using the brevity penalty (BP). The higher the score, the more similar the sentences are (Papineni et al., 2002). The formula for calculating BLEU is:

$$BP \cdot \exp \left(\sum_{n=1}^N w_n \log p_n \right)$$

G. ROUGE: The ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metric evaluates sentence summarization systems. It computes the overlap

between the summary and the original text. The higher the score, the better the summary (Lin, 2004). Several variations of ROUGE are used in this work: ROUGE-1, ROUGE-2, ROUGE-L, and ROUGE-LSum. ROUGE-1 evaluates the overlap of unigrams (individual words) between the generated and reference summaries, ROUGE-2 considers bigrams (two consecutive words), and ROUGE-L considers the longest common sequences of words. ROUGE-LSum is similar to ROUGE-L but splits the text based on “ η ” characters instead of spaces. The formula for calculating ROUGE-N, which evaluates the overlap of n-grams, is:

$$\frac{\sum_r \sum_s \text{match}(\text{gram}_{s,c})}{\sum_r \sum_s \text{count}(\text{gram}_s)}$$

These metrics are used to evaluate model performance in generating summaries based on text input. Task-specific evaluation strategies are described in Section 4.5 below.

4.4. Task-Specific Evaluation Strategies

There is a limited number of benchmark evaluation datasets available for low-resource languages. Therefore, specific strategies are followed for each NLP task to ensure a fair comparison with existing benchmarks and compare performance between novel, fine-tuned models and state-of-the-art models. Table 4.2 below lists evaluation metrics and evaluation strategies for each scenario:

Task	Evaluation Metric	Evaluation Strategy
Summarization	Loss Rouge1 Rouge2 RougeL RougeLSum	<ol style="list-style-type: none"> 1. Test each language’s performance on an unseen “test” set, following the example of the current benchmark models. 2. Compare the performance of the fine-tuned model to those of existing benchmark models on the test dataset.
Named Entity Recognition	Loss Accuracy Precision Recall F1 Score	<ol style="list-style-type: none"> 3. Test each language’s performance on an unseen “test” set, following the example of current benchmark models. 4. Compare the performance of the fine-tuned model to those of existing benchmark models on the test dataset.
Translation	BLEU	<ol style="list-style-type: none"> 5. Test translation performance for Sanskrit-to-English and English-to-Sanskrit on a previously unseen test dataset. 6. Evaluate the performance of the fine-tuned model on the Facebook Low Resource (FLoRes) machine translation benchmark evaluation dataset, which contains professionally refined translations across a wide variety of low-resource languages, including Sanskrit. 7. Compare the performance of the fine-tuned model to those of existing benchmark models on the test dataset and the FLoRes evaluation benchmark dataset.

Table 4.2: Evaluation Metrics for NLP Tasks

After establishing datasets and evaluation metrics for each NLP task addressed in the thesis, a model selection and training pipeline was developed. This process is described in detail in Chapter 5.

Chapter 5: Model Architecture, Selection, and Development

Training language models from scratch can be resource-intensive, costly, and impractical from a temporal efficiency standpoint. When developing a language model for a specific language or objective, a choice must be made between pre-training the model from scratch and fine-tuning a pre-trained model.

Due to the resource-intensiveness of training a model from scratch, fine-tuning a pre-trained model is a suitable option for a majority of research use-cases. With the development and release of groundbreaking transformer-based models like BERT and its derivatives, such as ALBERT and RoBERTa, there is a wealth of options to select from. The next section delves into the specific advantages and capabilities of transformer-based models and describes the models explored in this thesis at length.

5.1. Transfer Learning With Transformer-Based Models

The process of fine-tuning a pre-trained language model involves a technique referred to as "transfer learning." Transfer learning essentially involves "dropping" the head of a pre-trained model while training its body, thereby facilitating the "transfer" of the model's knowledge. The fine-tuning process can be customized for a specific task, such as masked language modeling or named entity recognition, based on the research objective.

Transfer learning vastly accelerates training time by "transferring" the knowledge acquired by the pre-trained model. It is important, however, to ensure that the pre-trained model is as similar to the new task it is fine-tuned on as possible.

The benefits and drawbacks of fine-tuning a pre-trained model (i.e. transfer learning) rather than pre-training a model from scratch are summarized in the following table:

Advantages of Transfer Learning	Disadvantages of Transfer Learning
Efficiency due to faster processing speed	Transfer of any bias existing in the original model
Lower environmental impact and computing costs	Skewing of results and/or performance if pre-trained model doesn't accurately match the downstream task

Table 5.1: Advantages and Disadvantages of Transfer Learning

This thesis focuses on transfer learning using various transformer-based models, described in Section 5.2. Transformers are based on an "attention" mechanism and consist of an encoder, decoder, or both. An encoder serves to convert input text into numerical representations (i.e. features or embeddings), while a decoder "decodes" representations from the encoder and outputs probabilities. Specifically, encoders are bi-directional and rely on a "self-attention" mechanism, while decoders are uni-directional, auto-regressive, and utilize masked attention and cross-attention. A model may be based on an encoder, decoder, or both in the case of encoder-decoder, or sequence-to-sequence, models (Vaswani et al., 2017).

In order to determine the model architecture best suited for each NLP task, various types of models are evaluated and compared. The following sections describe the model types and model selection process for named entity recognition, summarization, and translation. Section 5.2 includes a description of the model types utilized,

Section 5.3 details the model selection process for each type of task, and Section 5.4 covers the hyperparameter-optimized training process for the selected model.

5.2. Model Architectures

The following pre-trained models are loaded as baselines, fine-tuned, and compared for various natural language processing scenarios:

Task	Models Evaluated
Named entity recognition	Multilingual BERT (mBERT) ELECTRA XLM-RoBERTa
Summarization	mT5 mT5-XLSum IndicBART IndicBART-XLSum
Translation	mT5 mBART-50

Table 5.2: Models Evaluated by NLP Task

The individual models are briefly described below.

5.2.1 Multilingual BERT

Multilingual BERT (mBERT) is the multilingual variation of BERT that was simultaneously trained on 104 languages. This model has been shown to achieve state-of-the-art results across many NLP tasks (Devlin et al., 2018).

5.2.2 ELECTRA

ELECTRA is a simple but effective neural network architecture explicitly designed for sequence classification problems. Its main characteristics include self-attention and position embeddings (Clark et al., 2020). In terms of multilingual performance, it outperforms previous approaches such as RoBERTa and XLNet ().

5.2.3 XLM-RoBERTa

XLM-RoBERTa combines the multilingual pre-training method from XLM with the RoBERTa model. As a result, it achieves better performance than multilingual BERT (mBERT) on various cross-lingual benchmarks (Conneau et al., 2019).

5.2.4 mT5

Multilingual T5 (mT5) is a multilingual variation of T5 that was trained on 101 languages. It outperforms other multilingual variants such as XLM and XLM-RoBERTa (Xue et al., 2020).

5.2.5 mBART-50

Multi-Task Bidirectional Attention Flow Transformer (mBART-50) is a multi-task variant of BART that was trained on 104 different languages. It is a bidirectional transformer whose input consists of WordPiece representations (Tang et al., 2020).

The models described above were chosen because of their superior performance on various NLP tasks, as referenced by numerous benchmark papers. However, each model has advantages and disadvantages compared to others, particularly on specific NLP tasks. Therefore, several model candidates are trained on a subset or sample of the dataset for each NLP task covered, in order to optimize the hyperparameters before further training.

5.3. Hyperparameter Optimization

Before the full-scale model training process is implemented, hyperparameters are tuned for each model candidate, resulting in the optimal learning rate, weight decay, and epoch parameters. Training arguments incorporate these hyperparameters,

and the performance of all model candidates is compared following training. Finally, the best-performing model is selected for further training on the larger-scale dataset.

First, the datasets to utilize for optimization must be specified. In this case, a shuffled subset of the full dataset is used for faster computation:

```
1 tuning_subset_train = tokenized_datasets['train'].shuffle(seed=73).select(range
    (10000))
2 tuning_subset_validation = tokenized_datasets['test'].shuffle(seed=73).select(
    range(1000))
3
4 # Project constants
5 LR_MIN = 1.9e-5
6 LR_CEIL = 3.9e-5
7 WD_MIN = 0.008
8 WD_CEIL = 0.015
9 MIN_EPOCHS = 3
10 MAX_EPOCHS = 7
11 PER_DEVICE_EVAL_BATCH = 4
12 PER_DEVICE_TRAIN_BATCH = 4
13 NUM_TRIALS = 3
14 SAVE_DIR = f"optuna-test-translation"
15 NAME_OF_MODEL = f"sa-en-translation"
16 MAX_LENGTH = 512
17 study_name = "hp-search-multilingual-xlm-roberta"
18
19 def objective(trial: optuna.Trial):
20     training_args = TrainingArguments(
21         output_dir=SAVE_DIR,
22         learning_rate=trial.suggest_loguniform('learning_rate', low=LR_MIN, high=
LR_CEIL),
23         weight_decay=trial.suggest_loguniform('weight_decay', WD_MIN, WD_CEIL),
```

```

24     num_train_epochs=trial.suggest_int('num_train_epochs', low = MIN_EPOCHS,
high = MAX_EPOCHS),
25     per_device_train_batch_size=PER_DEVICE_TRAIN_BATCH,
26     per_device_eval_batch_size=PER_DEVICE_EVAL_BATCH,
27     gradient_accumulation_steps = 2,
28     save_steps=100000,
29     fp16=True,
30     disable_tqdm=False,)
31
32     trainer = Trainer(
33         model=model,
34         args=training_args,
35         train_dataset=tuning_subset_train,
36         eval_dataset=tuning_subset_validation,
37         data_collator=data_collator,
38         tokenizer=tokenizer,
39         compute_metrics=compute_metrics,)
40
41     result = trainer.train()
42     return result.training_loss

```

Listing 5.1: Optuna Hyperparameter Training Objective Specification

In this case, a translation model is tuned using a specified hyperparameter search space. The Optuna objective specifies the search to evaluate learning rates between $1.9e-5$ and $3.9e-5$, weight decay between 0.008 and 0.015, and epochs between 3 to 7. A total of three trials are conducted, with Trial 2 resulting in the best study parameters.

Trial 2 suggests using a learning rate of $3.053506639389683e-05$ and a weight decay of 0.011824988783648205 over 4 epochs. These parameters can now be used for full training.

Optuna also offers insights regarding optimization history for the best values, as well as the importance of each hyperparameter in terms of the objective value. The following figures represent this information visually:

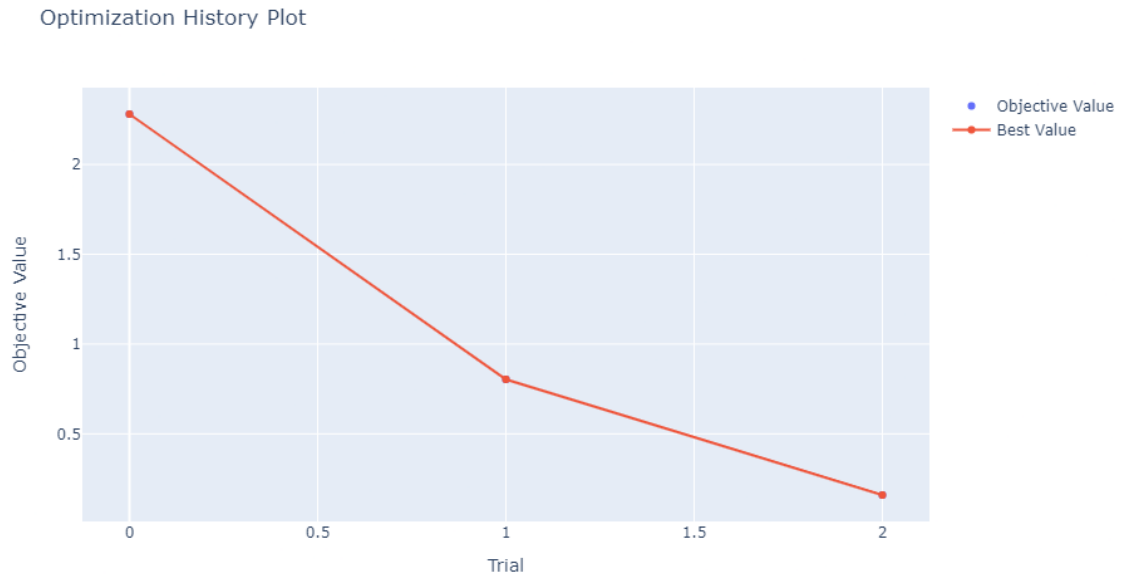


Figure 5.1: Optuna Optimization History

Parallel Coordinate Plot

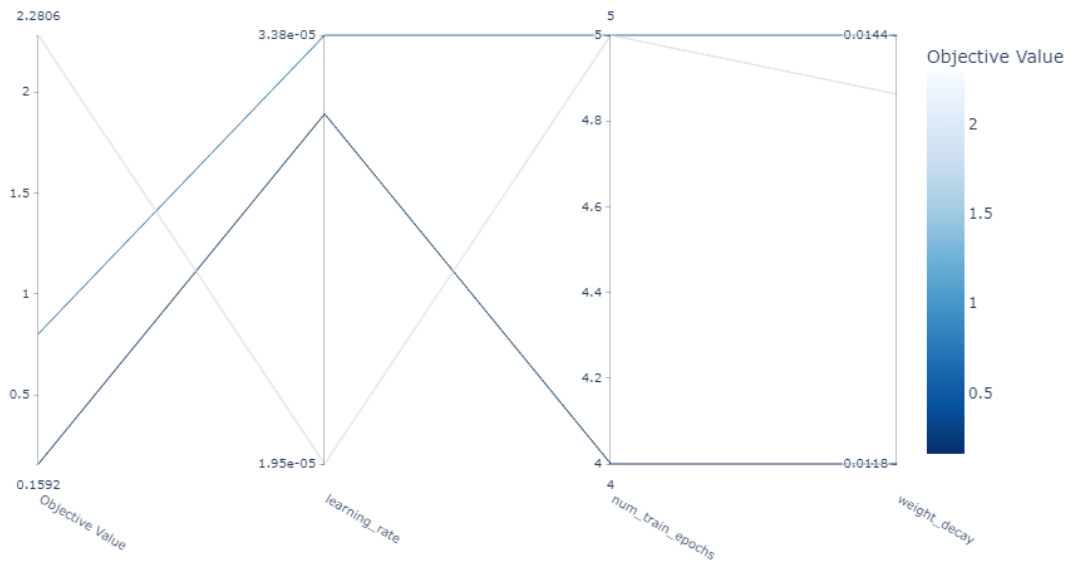


Figure 5.2: Optuna Parallel Coordinate Plot

Contour Plot

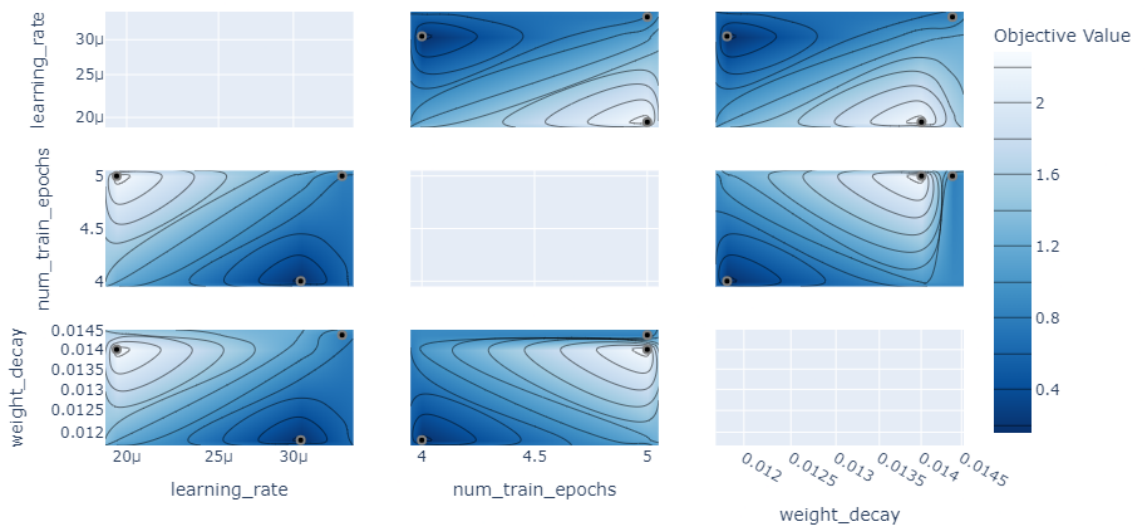


Figure 5.3: Optuna Contour Plot

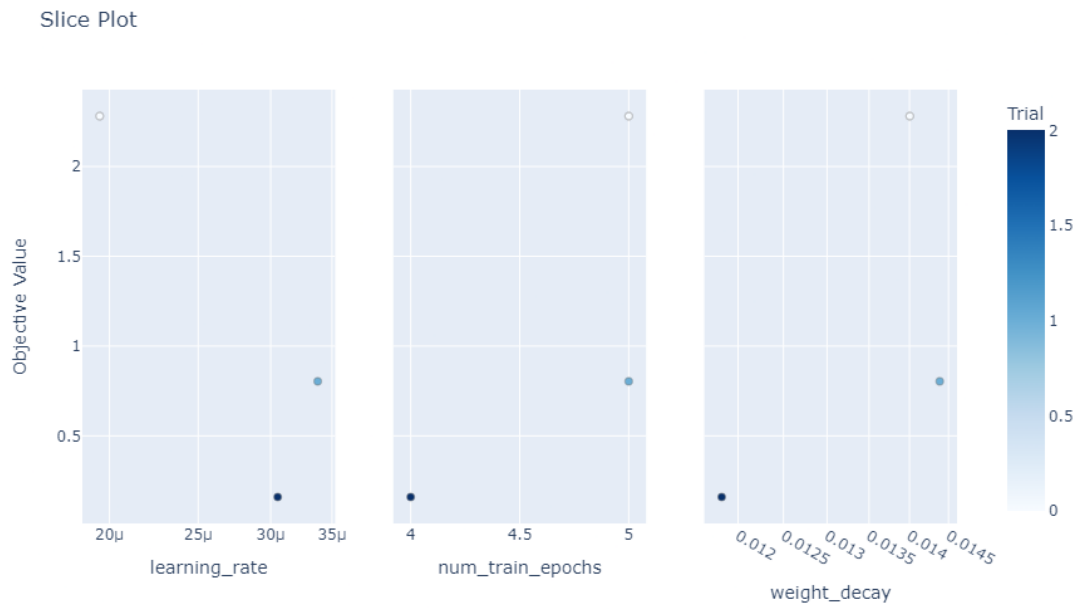


Figure 5.4: Optuna Slice Plot

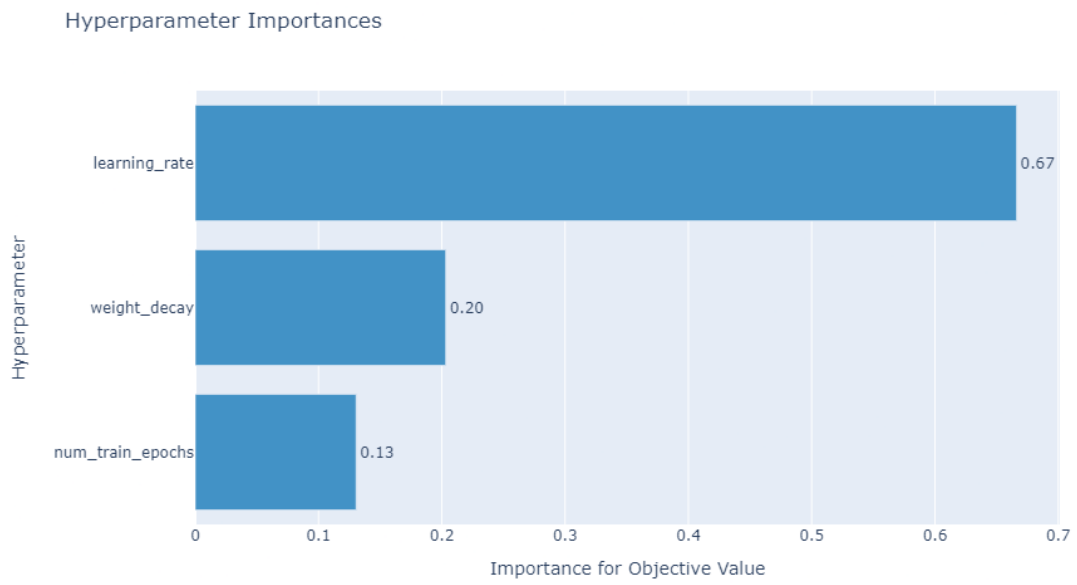


Figure 5.5: Optuna Hyperparameter Importances

5.4. Model Selection

The hyperparameter-based model selection process was implemented through the evaluation of several sets of models per task. Specifics and results are included in the following sections.

5.4.1 Named Entity Recognition

Specifically, the named entity recognition (NER) task involved training mBERT, ELECTRA, and XLM-RoBERTa models on the Oriya language subset of the Naama-padam dataset:

Model	Loss	Accuracy	Precision	Recall	F1
mBERT	0.5971	0.8820	0.5433	0.5397	0.6378
ELECTRA	0.5886	0.8835	0.5480	0.5377	0.6754
XLM-RoBERTa	0.2299	0.928	0.8191	0.8286	0.8239

Table 5.3: Performance Comparison of Monolingual NER Models

XLM-RoBERTa outperformed mBERT and ELECTRA and was thus selected for the training of the multilingual NER model.

5.4.2 Summarization

The summarization task involved training mT5, mT5-XLSum, IndicBART, and IndicBART-XLSum on the Sinhala language subset of the XLSum dataset. This facilitated comparisons between i. the mT5 and mBART models for NER and ii. models pre-trained on the XLSum dataset versus models without prior exposure to XLSum for Summarization.

The mT5-XLSum model outperformed other model variations and was thus selected for the training of the multilingual summarization model:

Model	Loss	Rouge1	Rouge2	RougeL	RougeLSum
mT5	1.905	9.987	2.304	8.907	8.756
mT5-XLSum	1.607	13.7865	5.0629	11.3648	11.3862
IndicBART	2.405	4.056	1.305	3.402	3.477
IndicBART-XLSum	3.489	5.839	2.487	3.897	3.897

Table 5.4: Performance Comparison of Monolingual Summarization Models

5.4.3 Translation

Finally, the translation task involved evaluating 15 pre-trained models on the full FLoRes (dev+devtest) dataset to compare translation performance on a previously unseen language pair. Both Sanskrit-to-English and English-to-Sanskrit translations are evaluated:

Model	BLEU (Sanskrit to English)	BLEU (English to Sanskrit)	Average BLEU Score
XLM-ProphetNet	0.011879377928836	0.005861833242	0.008870605585496
t5 Large	0.000000000624575	0.133689755	0.066844877793178
mT5 (base)	0.022259849073041	0.01555443717	0.018907143119520
m2m-100	0.074590953161494	0.472830417	0.273710685083422
NLLB-200 (Distilled-600M)	0.398453950111655	0.4310856611	0.414769805622258
mBART-50	0.649924249166765	0.6633522396	0.656638244372801
IndicBART	0.178649619263800	0.1375794077	0.158114513457517
Helsinki NLP - Opus EN-MUL	0.110251327811496	0.4235989039	0.266925115879636
Helsinki NLP - OPUS EN-HI	0.092890830261022	0.3812469298	0.237068880009920
banglat _{5_nmt_en_bn}	0.046909449315937	0.07511954087	0.061014495094270
banglat _{5_nmt_bn_en}	0.616880848055648	0.6145097945	0.615695321286366

Table 5.5: Performance Comparison of Baseline Translation Models

Among all the models compared, mBART-50 emerged as the best-performing

option and was selected for the training of the full-scale translation model. The decision to use the Facebook Low Resource (FLoRes) benchmark dataset for evaluation was based on the fact that it contains high-quality translations and is a widely accepted standard for evaluating machine translation systems. It also provides a uniform evaluation metric to utilize across various models.

5.5. Model Performance Monitoring

Model performance, training and evaluation time, and computation intensity are monitored through Weights & Biases. Model variations that required significantly longer to train while producing results that are equivalent, or even marginally better, than more efficient models were discarded. Thorough examination of these results aided the determination of ideal computation and hyperparameter variations, thereby expediting training and evaluation times as well as facilitating the development of more robust models.

Following are examples of outputs from training a multilingual named entity recognition model, with performance recorded across different runs.

For instance, various metrics are measured through the train loop, indicating which model variations are likely to achieve higher metrics. Named entity recognition, in particular, lends itself well to monitoring several tokens like person, place, organization, and others. Using Weights & Biases, it is possible to understand these metrics more closely in order to make a more informed model selection decision:

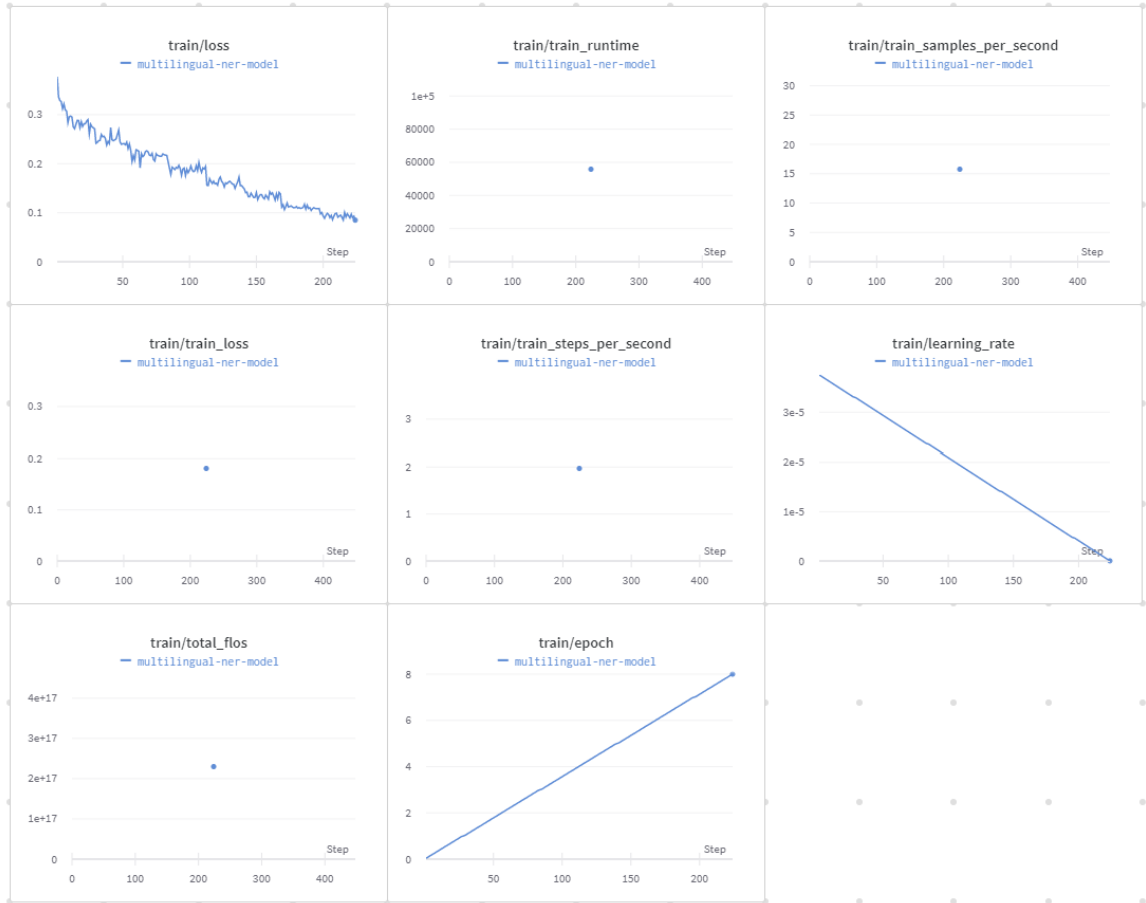


Figure 5.6: *Weights & Biases Train Loop Monitoring*

The evaluation process is also monitored, indicating changes in metrics like loss and overall accuracy as the trained model is evaluated at each epoch or specified step:



Figure 5.7: *Weights & Biases Evaluation Loop Monitoring*

Finally, the computation overhead, including the GPU, CPU, and disk utilization, is also monitored:

The models trained using the process described in Chapter 5 deliver strong performance on the test datasets compared to baseline models. Furthermore, a comparison with current benchmark models demonstrates that these fine-tuned models significantly improve over previous state-of-the-art methods.



Figure 5.8: Weights & Biases Evaluation Performance Monitoring

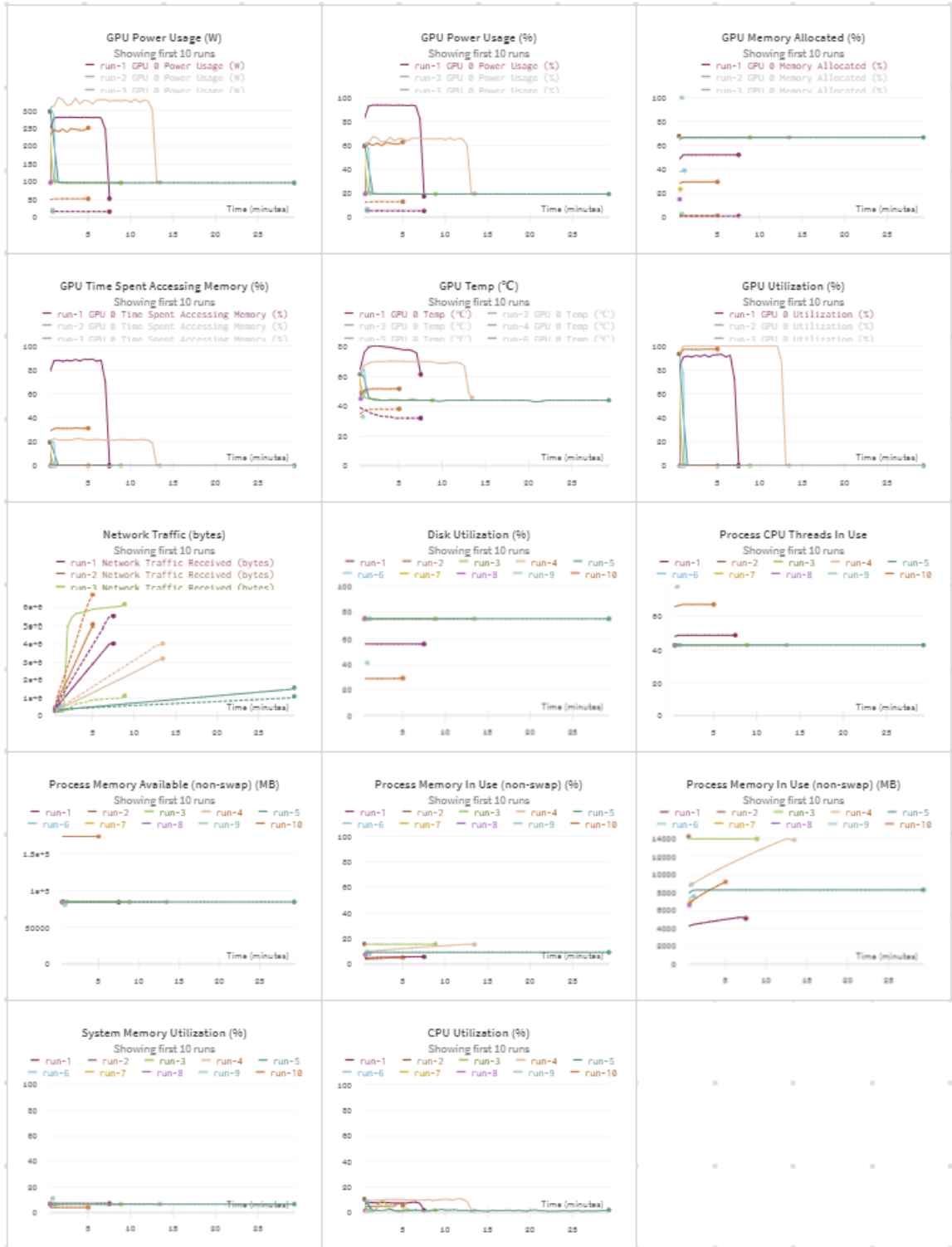


Figure 5.9: Weights & Biases Compute Usage Monitoring

Chapter 6: Model Performance

6.1. Model Training

Once an optimal model architecture has been selected for an NLP task, it is trained on a large shuffled and downsampled dataset. The validation split of the dataset is utilized for evaluation metrics as a part of the training loop, and the test split is used for predictions and further testing.

6.1.1 Named Entity Recognition

For the named entity recognition task, an XLM-RoBERTa model is trained on the full Naamapadam dataset, which includes 11 South Asian languages: Assamese, Bengali, Gujarati, Kannada, Hindi, Malayalam, Marathi, Oriya, Punjabi, Tamil, and Telugu.

Each language’s train and validation splits are concatenated to form comprehensive train and validation sets. Optuna is utilized for hyperparameter tuning, and the optimal hyperparameters obtained are used to train the full-scale XLM-RoBERTa model.

Finally, the test splits within the Naamapadam dataset dictionary are used to evaluate the model’s performance on NER for each language. The results are available in Chapter 6, along with performance comparisons on benchmark models.

6.1.2 Summarization

For the summarization task, an mT5 model previously trained on the full XLSum dataset is fine-tuned on the South Asian language splits of the XLSum dataset, including Bengali, Gujarati, Hindi, Marathi, Nepali, Punjabi, Sinhala, Tamil, Telugu, and Urdu.

Once again, each language’s train and validation splits are concatenated to form comprehensive train and validation sets.

Optuna is utilized for hyperparameter tuning, and the optimal hyperparameters obtained are used to train the full-scale mT5 model.

Finally, the test splits for each South Asian language within the XLSum dataset dictionary are used to individually evaluate the model’s performance on summarization for each language. The results are available in Chapter 6, along with performance comparisons on benchmark models.

6.1.3 Translation

An mBART-50 model is trained on the full dataset containing Sanskrit-English and English-Sanskrit translation pairs for the translation task. Train and validation splits are used during the training and evaluation process.

Optuna is utilized for hyperparameter tuning, and the optimal hyperparameters obtained are used to train the full-scale mBART-50 model.

Finally, the test split of the dataset is used to assess the model’s performance on unseen data. The resulting model is also evaluated on the FLoRes dataset and compared to benchmark models. All metrics and comparisons are available in Chapter 6.

This section covered the selection, optimization, and subsequent training of

NLP models for named entity recognition, summarization, and translation. The following section includes the evaluation results for each of these models, as well as corresponding comparisons with current benchmark models.

6.2. Model Evaluation and Results

6.2.1 Named Entity Recognition

For named entity recognition (NER), loss, accuracy, precision, recall, and F1 metrics are evaluated.

Following the evaluation process of IndicNER, the test dataset is tokenized and evaluated by language:

```
1 tokenized_test_set = {}
2
3 for lang in raw_datasets:
4     tokenized_test_set[lang] = raw_datasets[lang]['test'].map(
5         tokenize_and_align_labels,
6         batched=True,
7         num_proc=32,
8         load_from_cache_file=True,
9         desc="Running tokenizer on test dataset of language {0}".format(lang),
10    )
11
12 final_metrics = {}
13
14 for lang in tokenized_test_set:
15     predictions, labels, metrics = trainer.predict(tokenized_test_set[lang],
16         metric_key_prefix=lang)
17
18 lang_specific_results = {}
```

```

18 for key in metrics:
19     if 'overall_precision' in key:
20         lang_specific_results['Precision'] = metrics[key]
21     elif 'overall_recall' in key:
22         lang_specific_results['Recall'] = metrics[key]
23     elif 'overall_f1' in key:
24         lang_specific_results['F1'] = metrics[key]
25     elif 'overall_accuracy' in key:
26         lang_specific_results['Accuracy'] = metrics[key]
27     elif 'loss' in key:
28         lang_specific_results['Loss'] = metrics[key]
29 final_metrics[lang] = lang_specific_results

```

Listing 6.1: Evaluation of NER Models

Performance on the test dataset for each language trained is listed in the table below:

Language	Loss	Precision	Recall	F1	Accuracy
Assamese	0.169148	0.363636	0.500000	0.421053	0.940299
Bengali	0.165400	0.755203	0.764293	0.759721	0.950631
Gujarati	0.135021	0.746835	0.819444	0.781457	0.957278
Hindi	0.129340	0.810011	0.864078	0.836171	0.960976
Kannada	0.172282	0.789668	0.819296	0.804209	0.943417
Malayalam	0.197746	0.833001	0.831894	0.832447	0.934380
Marathi	0.163358	0.809976	0.813842	0.811905	0.948280
Oriya	0.523524	0.233766	0.155508	0.186770	0.853328
Punjabi	0.263783	0.646846	0.684638	0.665205	0.921533
Tamil	0.203594	0.760870	0.853659	0.804598	0.929577
Telugu	0.115842	0.859155	0.859155	0.859155	0.961616

Table 6.1: Language-Specific Performance of Fine-Tuned Multilingual IndoLibNER Model

The fine-tuned IndoLibNER model outperforms the benchmark IndicNER model on all of the listed languages, as shown above. Evaluation results using the IndicNER model are included in the table below.

Language	Loss	Precision	Recall	F1	Accuracy
Assamese	8.418570	0.000000	0.000000	0.000000	0.018657
Bengali	8.115846	0.032444	0.045135	0.037752	0.024571
Gujarati	6.699370	0.000000	0.000000	0.000000	0.025316
Hindi	7.283190	0.023979	0.032767	0.027692	0.029268
Kannada	6.025383	0.019126	0.021440	0.020217	0.029988
Malayalam	6.131558	0.027012	0.032558	0.029527	0.036208
Marathi	6.739916	0.035337	0.045943	0.039948	0.031099
Oriya	5.844997	0.023114	0.036717	0.028369	0.026776
Punjabi	6.984297	0.030512	0.034965	0.032587	0.035594
Tamil	5.823369	0.028846	0.036585	0.032258	0.070423
Telugu	5.715331	0.033333	0.028169	0.030534	0.028283

Table 6.2: Language-Specific Performance of Benchmark IndicNER Model

6.2.2 Summarization

ROUGE scores (Rouge1, Rouge2, RougeL, and RougeLSum) are reported for summarization models. Because the fine-tuned IndoLibSUM did not outperform the benchmark csebuetnlp/mT5_multilingual_XLSum model at the multilingual scale, it was not further evaluated on individual languages. The evaluation metrics are listed in the table below.

Loss	Rouge1	Rouge2	Rougel	Rougesum
1.607	13.7865	5.0629	11.3648	11.3862
1.6063	13.7156	5.0303	11.2804	11.3094
1.6507	13.6056	5.0362	11.1609	11.1783
1.6503	13.5593	4.9916	11.1269	11.1483
1.6531	13.5221	4.9599	11.0911	11.1095

Table 6.3: Performance of Fine-Tuned Multilingual Summarization Models

These metrics suggest that fine-tuning the benchmark on the same dataset (XL-Sum) does not elicit superior performance. Performance improvements could potentially be achieved by fine-tuning the csebuetnlp/mT5_multilingual_XLSum on new datasets or additional languages.

It is also possible to conclude from the initial model selection process that the mT5 model architecture lends itself best to the task of summarizing text in South Asian languages. Further evaluation and comparison of IndicBART-based models on additional languages may be necessary to definitively conclude the superiority, or otherwise, of mT5-based models.

6.2.3 Translation

Finally, for translation, BLEU scores are reported for both English-Sanskrit and Sanskrit-English translation predictions. The FLoRes-200 benchmark evaluation

dataset is utilized for these comparisons (Costa-jussà et al., 2022; Goyal et al., 2021; Guzmán et al., 2019).

```
1 # Load datasets for Sa-En and En-Sa translation
2 from datasets import load_dataset
3
4 san_to_eng_test = load_dataset("facebook/flores", "san_Deva-eng_Latn",split="dev+
   devtest")
5 eng_to_san_test = load_dataset("facebook/flores", "eng_Latn-san_Deva",split="dev+
   devtest")
```

Listing 6.2: Load FloRes Benchmark Dataset for Two Translation Directions

Both splits of the evaluation dataset are used in order to maximize the accuracy of the comparison.

For instance, the fine-tuned model developed in this thesis is evaluated in the Sanskrit-to-English translation direction as follows:

```
1 # Evaluate generated predictions against reference predictions in FloRes dataset
2 from datasets import load_metric
3 metric = load_metric('sacrebleu')
4
5 for r in results:
6     prediction = r['prediction']
7     reference = [r['reference']]
8     metric.add(prediction=prediction, reference=reference)
9 metric.compute()
```

Listing 6.3: Compute BLEU Metric for Translation Model

Performance on the FLoRes evaluation dataset in comparison to all candidate models is included in the table below. Both translation directions, along with the average BLEU score for each model, are also included:

Model	BLEU (Sanskrit to English)	BLEU (English to Sanskrit)	Average BLEU
t5 Large	0.0000000006	0.1336897550	0.0668448778
mT5 (base)	0.0222598491	0.0155544372	0.0189071431
m2m-100	0.0745909532	0.4728304170	0.2737106851
NLLB-200 (Distilled-600M)	0.3984539501	0.4310856611	0.4147698056
mBART-50	0.6499242492	0.6633522396	0.6566382444
IndicBART	0.1786496193	0.1375794077	0.1581145135
IndoLIB-SaEn (fine-tuned model)	3.7573104547	0.5356411511	2.1464758029

Table 6.4: FLoRes Dataset Performance of Fine-Tuned Translation Model Versus Benchmarks

In summation, when tested on their respective test sets and compared to several standard benchmarks, the models developed in this thesis are competitive with, if not superior to, current benchmark models for the languages covered. The summarization model does not outperform the benchmark model, but nevertheless demonstrates strong performance over baseline BART-based models. The NER and translation models demonstrate significant performance improvements over baseline and benchmark models.

Chapter 7: Conclusion

This thesis focuses on the development of novel datasets and fine-tuned models to address a variety of NLP tasks to improve the inclusion of South Asian languages in NLP systems. In particular, it investigates how hyperparameter optimization and model selection can be leveraged to improve the performance of existing models when applied to a target language or multiple languages. It further evaluates the effectiveness of various techniques for building NLP models by comparing them against state-of-the-art models. Finally, it introduces new datasets for use in language modeling and language detection for 31 South Asian languages across three language families: Indo-Aryan, Sino-Tibetan, and Dravidian.

7.1. Summary

The outcome of this thesis is IndoLib, a set of tools designed to enhance the inclusion of South Asian languages in NLP systems. IndoLib includes the development of datasets as well as monolingual and multilingual models to address several tasks: language modeling, language identification, named entity recognition, summarization, and machine translation. Specific contributions of this work include:

- (i) Datasets for language modeling and language identification, encompassing 31 languages spoken in South Asia, including languages previously unrepresented or underrepresented in NLP datasets and models.

- (ii) Monolingual models for Sinhala summarization and Oriya named entity recognition.
- (iii) Multilingual models for named entity recognition and summarization.
- (iv) A bidirectional machine translation model to translate Sanskrit to English and vice versa.

7.2. Ethical and Environmental Considerations

NLP research has various environmental and ethical implications, including but not limited to: substantial energy consumption to train massive language models; biased models and datasets; and underrepresentation or misrepresentation of certain groups or cultures. To minimize these issues, NLP researchers must optimize transparency regarding datasets used, processing steps, model training pipelines, and known biases.

Fine-tuning pre-trained models significantly optimizes the environmental footprint of training language models. The fine-tuning process involves running a small number (e.g., 100) of epochs on the original dataset while modifying only a few parameters (e.g., weights). For example, one can use the pre-trained model as a feature extractor by removing all layers except the last layer before passing the data through a dense network that predicts sentiment polarity. This approach reduces the computation required to classify text and improves performance because the final classification layer is trained from scratch rather than using the pre-trained weights.

Biases in language models can be challenging to address, especially as sizeable pre-trained language models function as “black boxes” that are challenging to decipher or accurately predict, particularly when exposed to new scenarios. Ensuring proper representation of diverse populations within the datasets used to fine-tune

models is a sensible starting point in understanding and addressing potential biases.

7.3. Future Work

This thesis covers developing fine-tuned models for named entity recognition, summarization, and translation of various South Asian languages. It also included the development of labeled and unlabeled datasets for 31 languages, including several languages that are either underrepresented or entirely unrepresented by current language models. Future work may expand upon this thesis by:

- (i) Developing language models based on the unlabeled multilingual dataset to include additional South Asian languages in future BERT-like models.
- (ii) Optimizing language identification systems using the labeled multilingual dataset to better differentiate between different South Asian languages, particularly languages with similar writing systems and phonologies.

Opportunities to contribute to the inclusion of South Asian languages in NLP are abundant. The results of this thesis work could facilitate further research and provide a baseline for further analyses, optimizations, and comparisons.

References

- Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-Generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (pp. 2623–2631).
- Aralikatte, R., de Lhoneux, M., Kunchukuttan, A., & Søggaard, A. (2021). Itihasa: A Large-Scale Corpus for Sanskrit to English Translation.
- Arora, G. (2020). iNLTK: Natural Language Toolkit for Indic Languages.
- Baumann, P. & Pierrehumbert, J. (2014). Using Resource-Rich Languages to Improve Morphological Analysis of Under-Resourced Languages. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation* (pp. 3355–3359).
- Bird, S. (2006). NLTK: The Natural Language Toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions* (pp. 69–72).
- Bisong, E. (2019). Matplotlib and Seaborn. In *Building Machine Learning and Deep Learning Models on Google Cloud Platform* (pp. 151–165). Springer.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M. J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., et al. (2018). JAX: Composable Transformations of Python+ NumPy Programs.

- Clark, K., Luong, M.-T., Le, Q. V., & Manning, C. D. (2020). Electra: Pre-Training Text Encoders as Discriminators Rather Than Generators.
- Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2019). Unsupervised Cross-Lingual Representation Learning at Scale.
- Costa-jussà, M. R., Cross, J., Çelebi, O., Elbayad, M., Heafield, K., Heffernan, K., Kalbassi, E., Lam, J., Licht, D., Maillard, J., et al. (2022). No Language Left Behind: Scaling Human-Centered Machine Translation.
- Davis, J. & Goadrich, M. (2006). The Relationship Between Precision-Recall and ROC Curves. In *Proceedings of the 23rd International Conference on Machine Learning* (pp. 233–240).
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-Training of Deep Bidirectional Transformers for Language Understanding.
- Eberhard, D. M., Simons, G. F., & Fennig, C. D. (2022). Ethnologue: Languages of the World. *SIL International, Dallas*, 25.
- Goutte, C. & Gaussier, E. (2005). A Probabilistic Interpretation of Precision, Recall and F-Score, With Implication for Evaluation. In *European Conference on Information Retrieval* (pp. 345–359).: Springer.
- Goyal, N., Gao, C., Chaudhary, V., Chen, P.-J., Wenzek, G., Ju, D., Krishnan, S., Ranzato, M., Guzmán, F., & Fan, A. (2021). The FLORES-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation.
- Gugger, S., Debut, L., Wolf, T., Schmid, P., Mueller, Z., & Mangrulkar, S. (2022). Accelerate: Training and Inference at Scale Made Simple, Efficient and Adaptable.

- Guzmán, F., Chen, P.-J., Ott, M., Pino, J., Lample, G., Koehn, P., Chaudhary, V., & Ranzato, M. (2019). Two New Evaluation Datasets for Low-Resource Machine Translation: Nepali-English and Sinhala-English.
- Harris, C. R., Millman, K. J., Van Der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., et al. (2020). Array Programming With NumPy. *Nature*, 585(7825), 357–362.
- Hasan, T., Bhattacharjee, A., Islam, M. S., Samin, K., Li, Y.-F., Kang, Y.-B., Rahman, M. S., & Shahriyar, R. (2021). XL-Sum: Large-Scale Multilingual Abstractive Summarization for 44 Languages.
- Heek, J., Levskaya, A., Oliver, A., Ritter, M., Rondepierre, B., Steiner, A., & van Zee, M. (2020). Flax: A Neural Network Library and Ecosystem for JAX.
- Heffernan, K., Çelebi, O., & Schwenk, H. (2022). Bitext Mining Using Distilled Sentence Representations for Low-Resource Languages.
- Imambi, S., Prakash, K. B., & Kanagachidambaresan, G. (2021). PyTorch. In *Programming With TensorFlow* (pp. 87–104). Springer.
- Joshi, P., Santy, S., Budhiraja, A., Bali, K., & Choudhury, M. (2020). The State and Fate of Linguistic Diversity and Inclusion in the NLP World.
- Kihlstrom, J. F. & Park, L. (2016). *Cognitive Psychology: Overview*.
- Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B. E., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J. B., Grout, J., Corlay, S., et al. (2016). *Jupyter Notebooks - A Publishing Format for Reproducible Computational Workflows*, volume 2016.

- Kudo, T. & Richardson, J. (2018). Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.
- Kunchukuttan, A. (2020). The IndicNLP Library. *Indian Language NLP Library*.
- Lamsal, R. (2020). A Large Scale Nepali Text Corpus. *IEEE dataport*.
- Le, V.-B. & Besacier, L. (2009). Automatic Speech Recognition for Under-Resourced Languages: Application to Vietnamese Language. *IEEE Transactions on Audio, Speech, and Language Processing*, 17(8), 1471–1482.
- Lhoest, Q., del Moral, A. V., Jernite, Y., Thakur, A., von Platen, P., Patil, S., Chaumond, J., Drame, M., Plu, J., Tunstall, L., et al. (2021). Datasets: A Community Library for Natural Language Processing.
- Lin, C.-Y. (2004). Rouge: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out* (pp. 74–81).
- Maryatt, E. (2018). Bringing Low-resource Languages and Spoken Dialects Into Play With Semi-Supervised Universal Neural Machine Translation.
- Maxwell, M. & Hughes, B. (2006). Frontiers in Linguistic Annotation for Lower-Density Languages. In *Proceedings of the Workshop on Frontiers in Linguistically Annotated Corpora* (pp. 29–37).
- McKinney, W. et al. (2011). Pandas: A Foundational Python Library for Data Analysis and Statistics. *Python for High Performance and Scientific Computing*, 14(9), 1–9.

- Mhaske, A., Khapra, M., Kedia, H., Rudramurthy, V., Kunchukuttan, A., & Kumar, P. (2022). Naamapadam: A Large-Scale Named Entity Annotated Data for Indic Languages.
- Morris, J. (2022). A Step by Step Guide to Tracking Hugging Face Model Performance.
- Orife, I., Kreutzer, J., Sibanda, B., Whitenack, D., Siminyu, K., Martinus, L., Ali, J. T., Abbott, J., Marivate, V., Kabongo, S., et al. (2020). Masakhane–Machine Translation For Africa.
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). Bleu: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics* (pp. 311–318).
- Patel, J. M. (2020). Web Scraping in Python Using Beautiful Soup Library. In *Getting Structured Data from the Internet* (pp. 31–84). Springer.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-Learn: Machine Learning in Python. *The Journal of Machine Learning Research*, 12, 2825–2830.
- Pimienta, D., Prado, D., & Blanco, Á. (2009). Twelve Years of Measuring Linguistic Diversity in the Internet: Balance and Perspectives.
- Pollock, D. C., Van Reken, R. E., & Pollock, M. V. (2010). *Third Culture Kids: The Experience of Growing Up Among Worlds*. Hachette UK.
- Powers, D. M. (2020). Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness and Correlation.

- Rowe, B. M. & Levine, D. P. (2018). *A Concise Introduction to Linguistics*. Routledge.
- Ruder, S. (2020). Why You Should Do NLP Beyond English.
- Ruder, S. (2021). ML and NLP Research Highlights of 2020.
- Salian, I. (2019). GPUs Help Researcher Decipher Ancient Sanskrit.
- Sengupta, D. & Saha, G. (2015). Study on Similarity Among Indian Languages Using Language Verification Framework. *Advances in Artificial Intelligence*.
- Sennrich, R., Haddow, B., & Birch, A. (2015). Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Suárez, P. J. O., Romary, L., & Sagot, B. (2020). A Monolingual Approach to Contextualized Word Embeddings for Mid-Resource Languages.
- Tang, Y., Tran, C., Li, X., Chen, P.-J., Goyal, N., Chaudhary, V., Gu, J., & Fan, A. (2020). Multilingual Translation With Extensible Multilingual Pretraining and Finetuning.
- Tunstall, L., Barham, M., Wolf, T., Lhoest, Q., & Platen, P. (2022). Wikipedia Dataset on Hugging Face.
- Vaswani, A., Shazeer, N. M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In *NIPS*.
- Whitenack, D. (2022). Bloom Library Language Modeling Dataset. *SIL International*.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al. (2020). Transformers: State-of-the-art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations* (pp. 38–45).

- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., Barua, A., & Raffel, C. (2020). mT5: A Massively Multilingual Pre-Trained Text-to-Text Transformer.
- Yue, D. (2022). Hugging Face: Embracing Natural Language Processing.
- Zhang, Z. & Sabuncu, M. (2018). Generalized Cross Entropy Loss for Training Deep Neural Networks With Noisy Labels. *Advances in Neural Information Processing Systems*, 31.