



Alif الف: An Urdu Generative Pre-trained Language Model

Ali Muhammad Asad, Hammad Sajid
Syeda Haya Fatima, Zainab Haider
Syed Muhammad Ali Naqvi
Supervisor: Dr. Abdul Samad



FYP
FINAL YEAR
PROJECT
SHOWCASE
2025

Introduction

Alif is a revolutionary pre-trained Urdu large language model.

Motivation

- Urdu is a *Low-Resource Language*
- Bias in LLMs due to Lack of Urdu Training Data
- Existing Models: High-Resource & Lack Security

Solution

- Develop an SLM which is smaller, efficient and performs well with low-computational resources for Urdu Language.
- Incurs minimum costs to fine-tune and used for various use cases like banks, customer support, etc in Pakistan.
- Sets a foundational benchmark for LLMs in Urdu Language for future Urdu-NLP research.

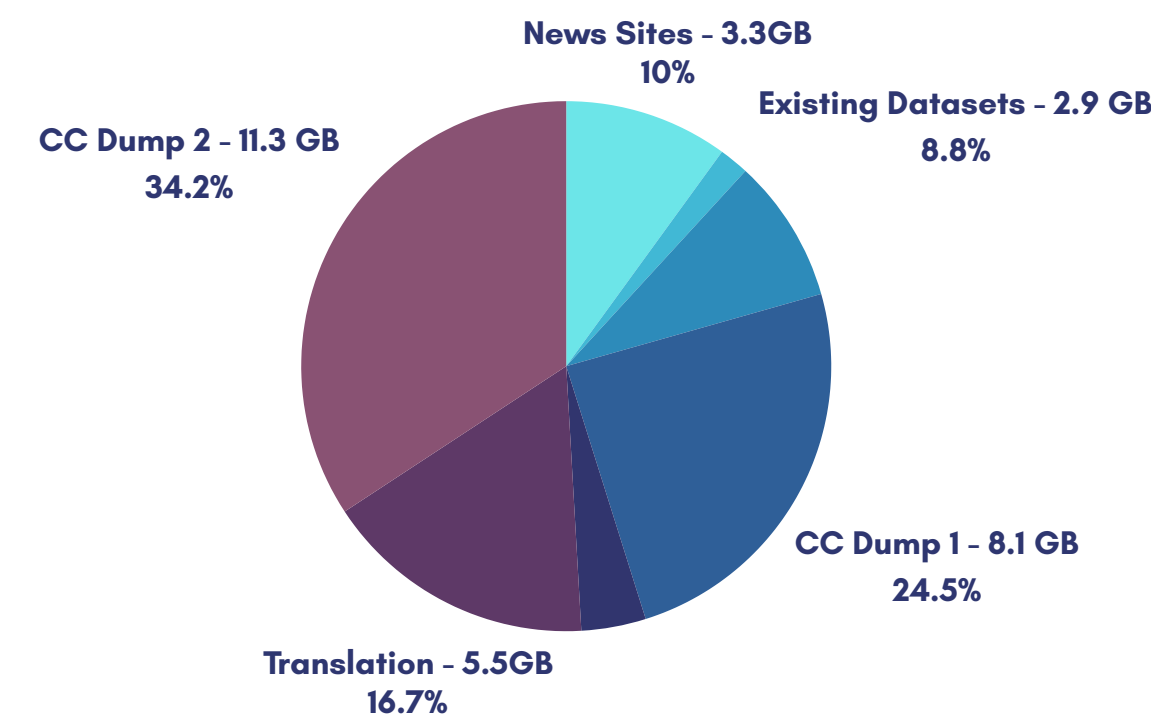
Pretraining Dataset Curation

Collection: Urdu text gathered from web scraping, scanned books, existing Datasets and translation, covering diverse domains like news, blogs, and literature.

Cleaning & Deduplication: Removed non-Urdu content, normalized text (spacing, encoding, Urdu numerals), eliminated noise, and used Jaccard Similarity with Min-Hashing for deduplication.

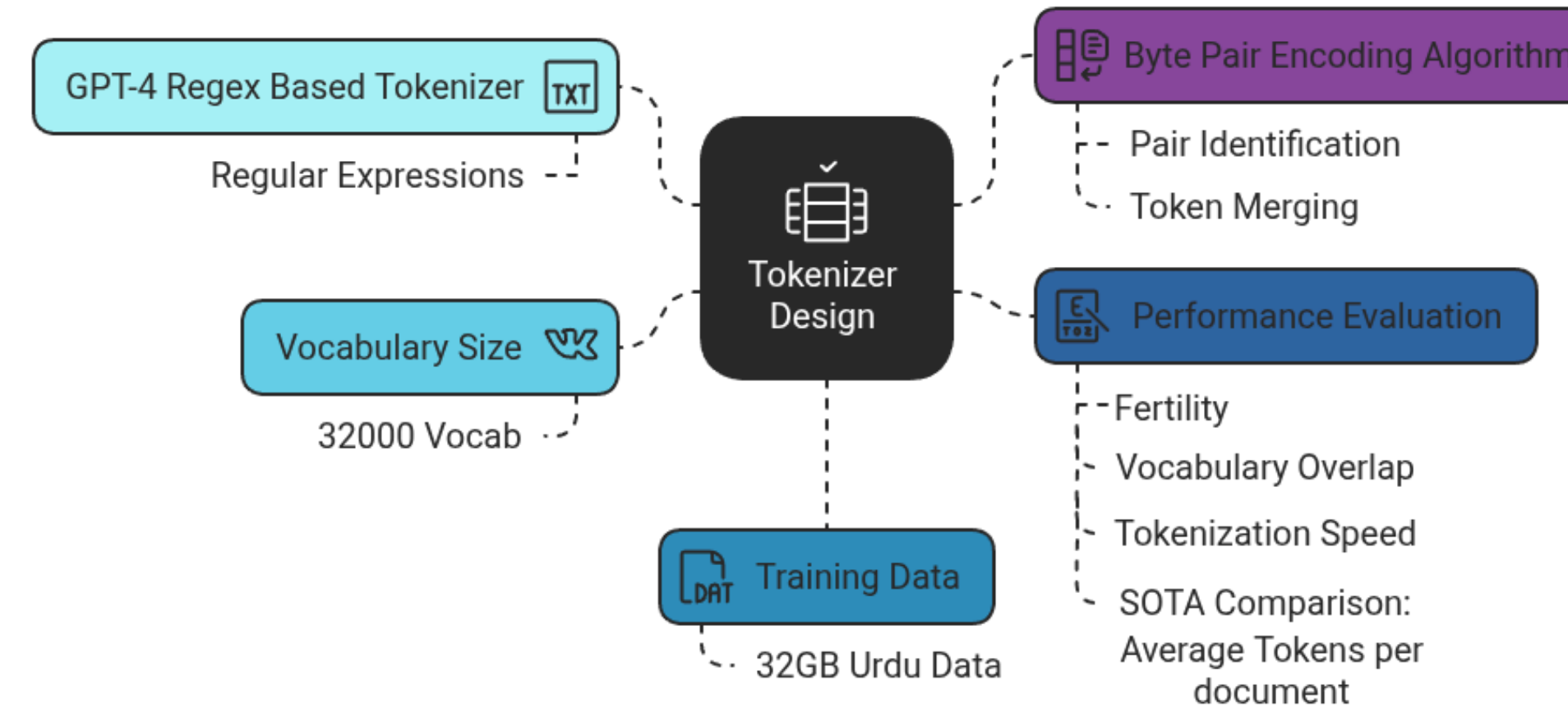
ALIF Pre-training Dataset:

- 33GB
- ~13 Million Rows
- ~5-6 Billion Tokens

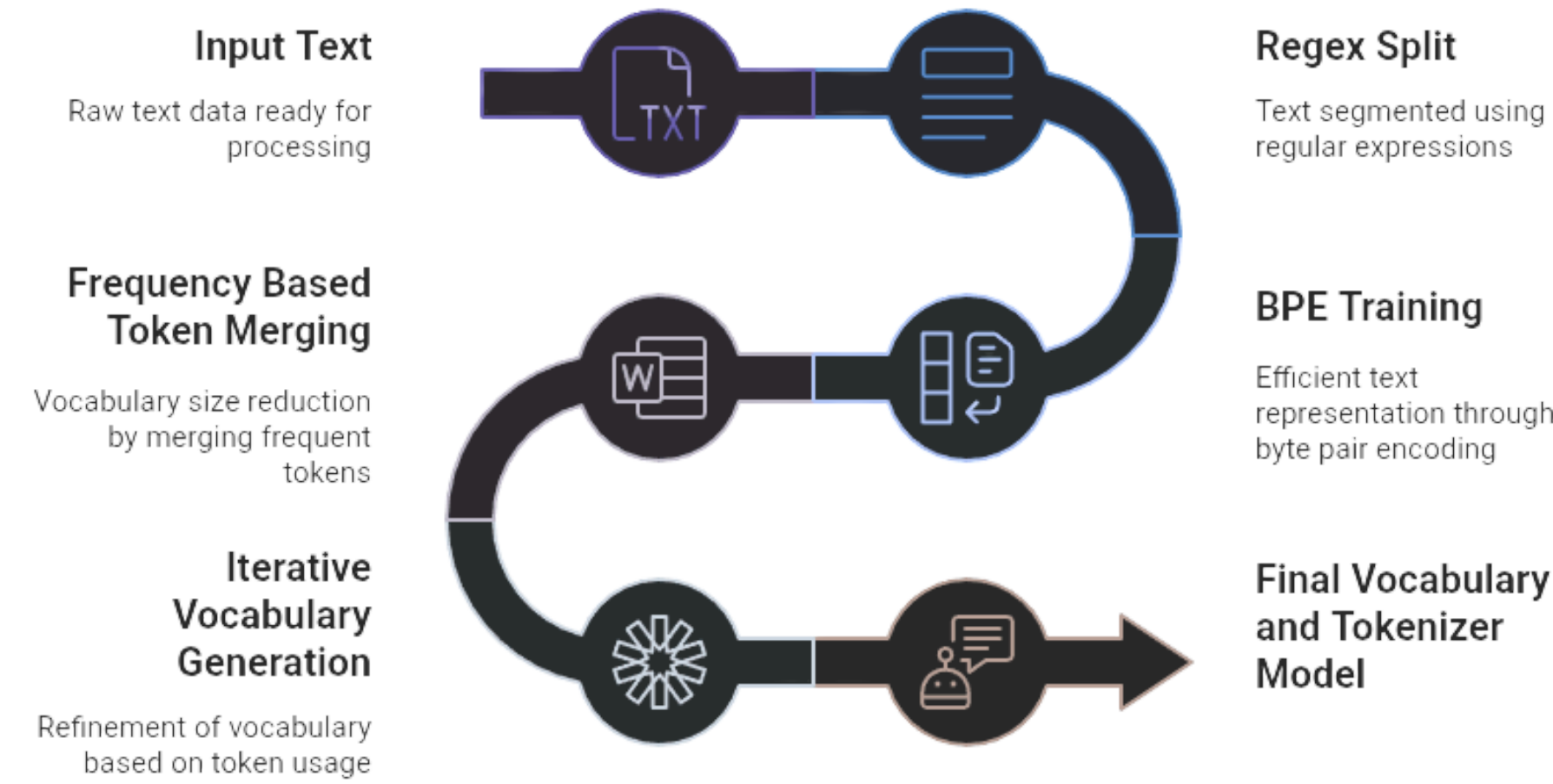


Instruct Dataset Curation

Tokenizer Architecture & Training

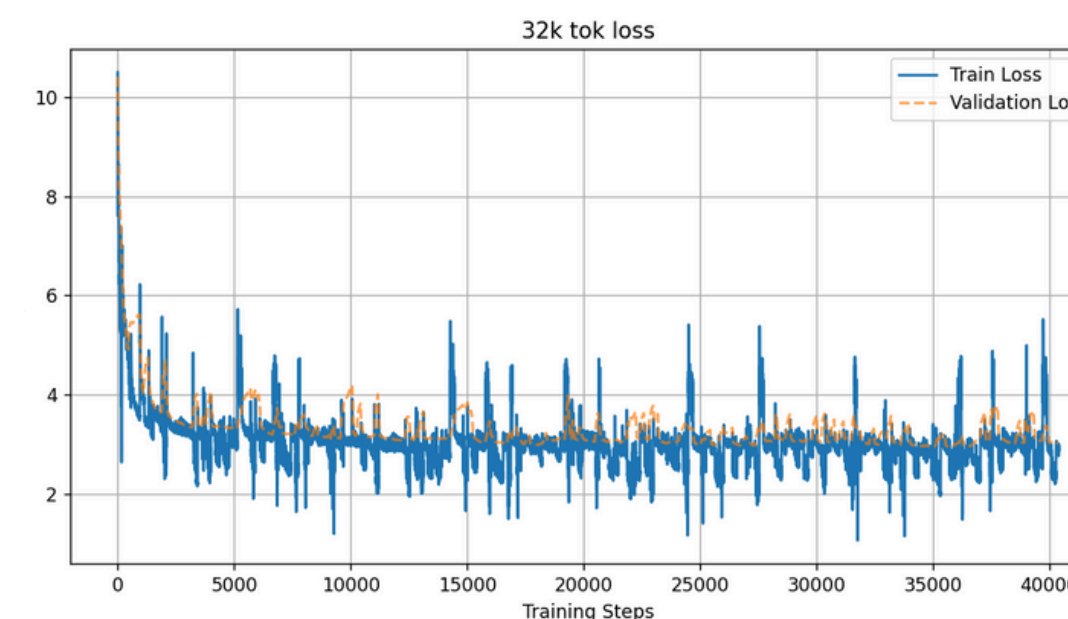
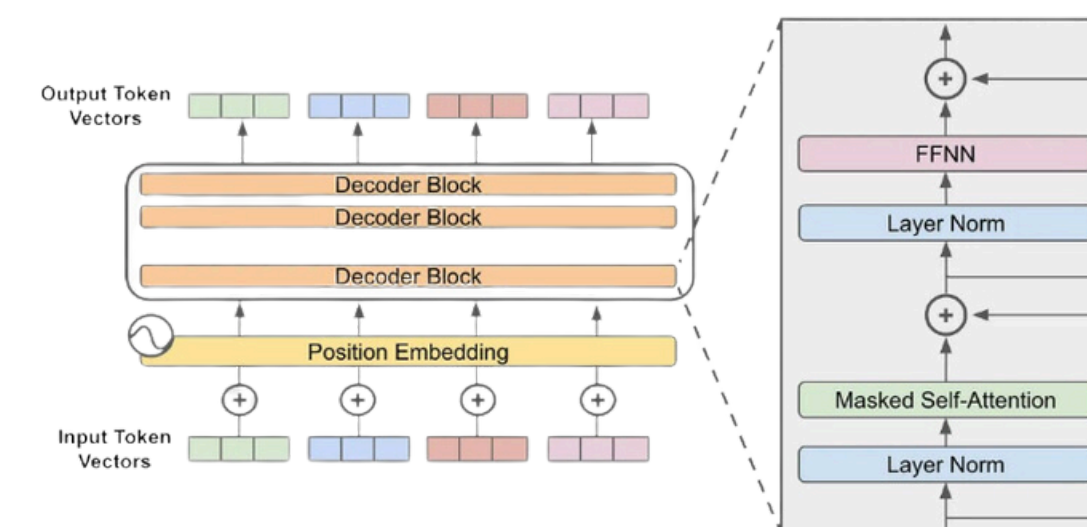


Tokenizer Architecture and Training



Model Architecture & Training

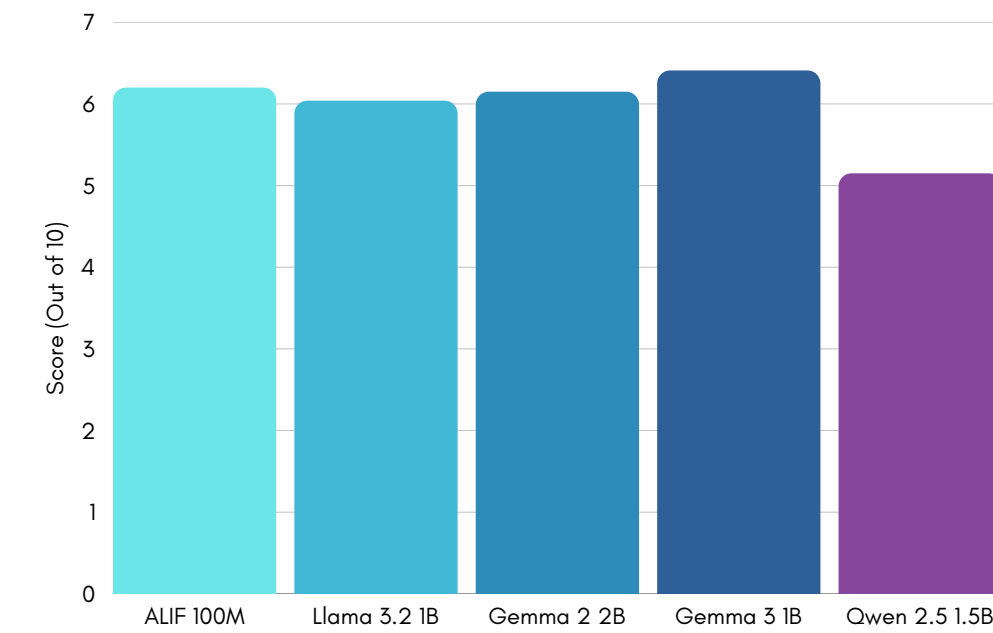
	Alif-Small	Alif-Large
Model Size	116M	1B
Block Size	1024	4096
Epochs	5	5
N. Heads	12	24
N.Layers	12	24



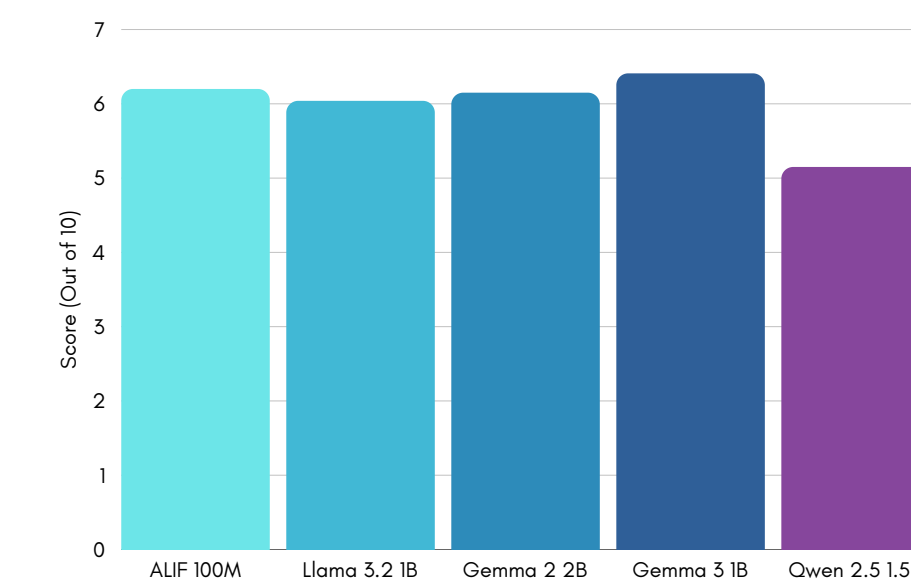
Evaluation

LLM-as-a-judge was used to score text generation on coherence, fluency, and relevance against SOTA models

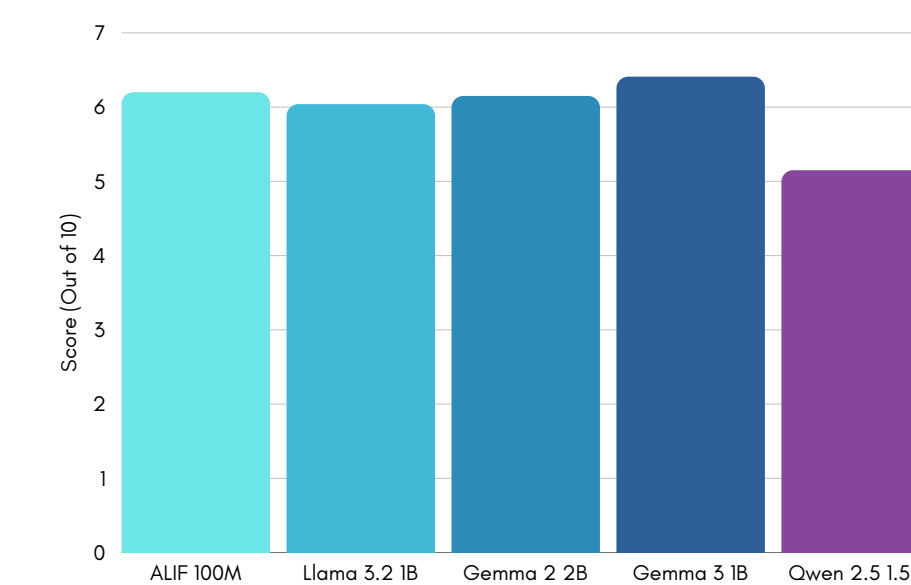
- Our model demonstrated comparable performance to SOTA models 10× its size.



Few Shot Testing

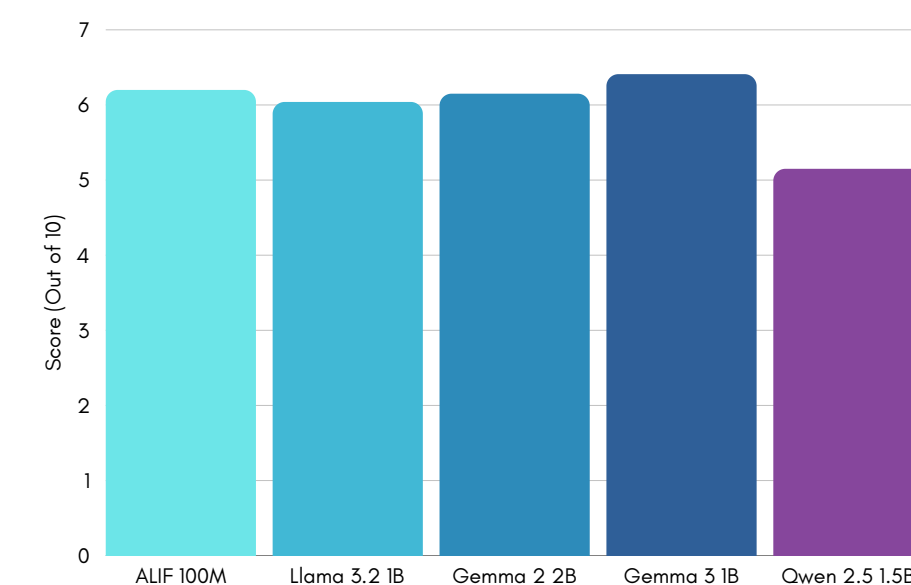


Sentiment Classification



Question Answering

Grammar Correction



Summarisation

Results

- The model performance indicates that monolingual Urdu model surpasses Multilingual Models
- The model has inherent understanding of the language and can be finetuned on urdu downstream tasks with minimal fine-tuning
- In Future, ALIF can be scaled to (8B, 16B) for commercial use cases.