# Translation NLP

A comprehensive machine translation system for Arabic-English language pairs, implementing state-of-the-art transformer models with evaluation metrics.

## Table of Contents

## Overview

This implementation provides a complete pipeline for bilingual translation between Arabic and English, using the MarianMT transformer model. The system includes data preprocessing, model training, evaluation, and interactive translation capabilities.

## Dependencies

The system requires the following Python packages:

```
datasets
scikit-learn
pandas
numpy
matplotlib
seaborn
transformers
nltk
emoji
torch
wordcloud
arabic-reshaper
python-bidi
rouge-score
```

## Implementation Details

### 1. Core Components

**Model Configuration**

```python
MODEL_CONFIGS = {
    'ar-en': {
        'model_name': 'Helsinki-NLP/opus-mt-ar-en',
        'src_lang': 'ar',
        'tgt_lang': 'en',
        'src_stopwords': set(stopwords.words('arabic')),
        'tgt_stopwords': set(stopwords.words('english')),
        'preserve_words': {'الى', 'على', 'في', 'عن', 'من', 'يوم', 'سنة', 'شهر'}
    },
    'en-ar': {
        'model_name': 'Helsinki-NLP/opus-mt-en-ar',
        'src_lang': 'en',
        'tgt_lang': 'ar',
        'src_stopwords': set(stopwords.words('english')),
        'tgt_stopwords': set(stopwords.words('arabic')),
        'preserve_words': {'to', 'on', 'in', 'at', 'from', 'day', 'year', 'month'}
    }
}
```

### 2. Main Pipeline Class

The `BilingualTranslationPipeline` class implements the following key components:

### Data Processing

- `load_data()` : Loads and preprocesses the dataset
- `clean_data()` : Cleans text by removing stopwords and normalizing characters
- `filter_data()` : Filters sentences by length
- `split_data()` : Splits data into training and validation sets

### Model Training

- `load_model()` : Loads the pre-trained MarianMT model
- `tokenize_data()` : Tokenizes the dataset for training
- `setup_training()` : Configures training parameters
- `train_model()` : Trains the model

### Translation

- `translate_text()` : Translates individual texts
- `interactive_translation_test()` : Provides an interactive interface for testing translations

# Usage

## Basic Usage

```
# Initialize pipeline
pipeline = BilingualTranslationPipeline('ar-en')

# Run full pipeline
pipeline.run_pipeline()

# Interactive translation
pipeline.interactive_translation_test()
```

## Evaluation

```
# Run evaluation on 100 samples
pipeline.evaluate_translations(num_samples=100)
```

# Evaluation System

The translation system implements a comprehensive evaluation framework using BLEU and ROUGE scores to assess translation quality. This document provides detailed information about the evaluation metrics and their implementation.

## 1. BLEU Score

The BLEU (Bilingual Evaluation Understudy) score is a precision-based metric that measures the similarity between machine translation and human reference translations.

### Implementation Details

```
from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction

def calculate_bleu(reference_text, translated_text):
    # Tokenize the texts
    reference_tokens = [reference_text.split()]
    candidate_tokens = translated_text.split()

    # Calculate BLEU score with smoothing
    smoothie = SmoothingFunction().method1
    bleu_score = sentence_bleu(reference_tokens, candidate_tokens,
                               smoothing_function=smoothie)
    return bleu_score
```

## Key Components

1. **N-gram Matching**
   - Compares n-grams (1-4) between candidate and reference translations
   - Weights: 0.4 for unigrams, 0.3 for bigrams, 0.2 for trigrams, 0.1 for 4-grams

2. **Smoothing Function**
   - Handles cases where there are no n-gram matches
   - Prevents zero scores for partial matches
   - Uses NLTK's smoothing method

3. **Score Range**
   - 0.0 to 1.0 (higher is better)
   - 1.0 indicates perfect translation
   - 0.0 indicates completely different translation

## 2. ROUGE Score

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a recall-based metric that measures the overlap of n-grams between the system and reference translations.

## Implementation Details

```python
from rouge_score import rouge_scorer


def calculate_rouge(reference_text, translated_text):
    # Initialize ROUGE scorer
    rouge_scorer_obj = rouge_scorer.RougeScorer(
        ['rouge1', 'rouge2', 'rougeL'],
        use_stemmer=True
    )

    # Calculate ROUGE scores
    scores = rouge_scorer_obj.score(reference_text, translated_text)
    return scores
```

## ROUGE Variants

1. **ROUGE-1 (Unigram Overlap)**
   - Measures overlap of unigrams (single words)
   - Formula: (Number of overlapping unigrams) / (Total unigrams in reference)
   - Best for: General translation quality assessment

2. **ROUGE-2 (Bigram Overlap)**
   - Measures overlap of bigrams (word pairs)
   - Formula: (Number of overlapping bigrams) / (Total bigrams in reference)
   - Best for: Phrase-level translation accuracy

3. **ROUGE-L (Longest Common Subsequence)**
   - Measures longest common subsequence of words
   - Considers word order and sentence structure
   - Best for: Overall translation fluency

## 3. Evaluation Process

The evaluation system follows these steps:

1. **Data Preparation**

   ```python
   # Get validation samples
   eval_samples = dataset['validation'].select(range(num_samples))
   ```

2. **Translation Generation**

   ```python
   # Generate translations
   for sample in eval_samples:
       source_text = sample[src_lang]
       reference_text = sample[tgt_lang]
       translated_text = pipeline.translate_text(source_text)['translated']
   ```

3. **Score Calculation**

```
 # Calculate metrics
bleu_scores = []
rouge_scores = defaultdict(list)

for translation in translations:
    bleu_score = calculate_bleu(reference_text, translated_text)
    rouge_scores_dict = calculate_rouge(reference_text, translated_text)
```

4. **Results Aggregation**

```
 # Calculate averages
avg_bleu = np.mean(bleu_scores)
avg_rouge = {metric: np.mean(scores)
             for metric, scores in rouge_scores.items()}
```

## 4. Visualization

The evaluation results are visualized in two ways:

1. **BLEU Score Distribution**
   - Histogram showing distribution of BLEU scores
   - Helps identify translation quality patterns
   - Includes kernel density estimation

2. **ROUGE Scores Comparison**
   - Bar chart comparing ROUGE-1, ROUGE-2, and ROUGE-L scores
   - Visual representation of different aspects of translation quality
   - Normalized scale (0-1) for easy comparison

## 5. Output Files

The evaluation system generates several output files:

1. **Evaluation Results**

```
evaluation_results_{direction}.txt
```

   Contains:
   - Average BLEU score
   - Average ROUGE scores
   - Number of samples evaluated
   - Timestamp of evaluation

2. **Visualizations**

```
visualizations/evaluation_scores_{direction}.png
```

   Contains:
   - BLEU score distribution plot
   - ROUGE scores comparison plot

## 6. Best Practices

1. **Sample Size**
   - Use at least 100 samples for reliable evaluation
   - Consider using more samples for statistical significance
   - Balance between evaluation time and accuracy

2. **Reference Quality**
   - Ensure high-quality reference translations
   - Use multiple reference translations when possible
   - Validate reference translations for accuracy

3. **Error Handling**
   - Handle empty translations gracefully
   - Log evaluation errors for debugging
   - Provide meaningful error messages

4. **Performance Optimization**
   - Use caching for repeated evaluations
   - Implement parallel processing for large datasets
   - Optimize memory usage for large-scale evaluation

## 7. Interpreting Results

1. **BLEU Score Interpretation**

   - 0.4-0.5: Good translation
   - 0.3-0.4: Acceptable translation
   - Below 0.3: Needs improvement

2. **ROUGE Score Interpretation**

   - ROUGE-1 > 0.4: Good word overlap
   - ROUGE-2 > 0.2: Good phrase matching
   - ROUGE-L > 0.3: Good sentence structure

3. **Combined Analysis**

   - Compare BLEU and ROUGE scores
   - Look for patterns in score distributions
   - Consider language-specific characteristics

# Visualization

The system generates several visualizations:

1. **Word Clouds**

   - Shows most frequent words in both languages
   - Helps in understanding vocabulary distribution

2. **Length Distributions**

   - Displays sentence length patterns
   - Helps in identifying optimal sequence lengths

3. **Training Metrics**

   - Shows training loss over time
   - Helps in monitoring model convergence

4. **Evaluation Scores**

   - BLEU score distribution
   - ROUGE scores comparison
   - Helps in understanding translation quality

# Output Files

The pipeline generates several output files:

- `translation_log_{direction}.txt` : Detailed training logs
- `evaluation_results_{direction}.txt` : Evaluation metrics
- `final_report_{direction}.txt` : Comprehensive pipeline report
- Visualizations in the `visualizations/` directory

# Notes

- The system uses GPU acceleration when available
- Implements caching for improved performance
- Includes comprehensive error handling
- Provides detailed logging throughout the pipeline