

Bridging the Gap: A Scalable Dialect Distillation Pipeline for Authentic Saudi Arabic Translation

Author: Ziad Abdullah **Date:** November 29, 2025

Abstract

Existing machine translation systems primarily focus on Modern Standard Arabic (MSA), failing to capture the rich linguistic diversity of Saudi dialects (Najdi, Hijazi, Eastern, Southern). This “dialect gap” hinders effective digital communication and cultural preservation. In this paper, we present a **Scalable Dialect Distillation Pipeline** that leverages Large Language Models (LLMs) to synthesize high-quality dialectal data from a small seed corpus. We fine-tuned an NLLB-200 model on a hybrid dataset of **1,000 authentic seed pairs** and **30,030 synthetic examples**. Our approach achieves a **BLEU score of 30.04**, significantly outperforming baseline models. We further validate our method through a rigorous ablation study and qualitative error analysis, demonstrating that our distillation technique effectively bridges the gap between formal MSA and authentic dialectal speech.

1. Introduction

1.1 The Linguistic Landscape

Saudi Arabia is home to a rich tapestry of Arabic dialects. Unlike Modern Standard Arabic (MSA), which is used in formal writing and news, dialects are the language of daily life. - **Najdi:** Central region (Riyadh). - **Hijazi:** Western region (Jeddah, Makkah). - **Eastern:** Gulf coast (Dammam). - **Southern:** Asir and Najran.

1.2 The Problem

Despite their prevalence, these dialects are severely underrepresented in NLP. Standard translation models often “correct” dialectal terms into MSA, losing the intended meaning and tone. * *Example:* “How are you?” -> “كيف حالك؟” (MSA) instead of “شلونك؟” (Saudi).

1.3 Contribution

We propose a robust translation system that: 1. Introduces a **Scalable Dialect Distillation Pipeline** using LLMs. 2. Fine-tunes a state-of-the-art **NLLB-200** model. 3. Provides a rigorous **Ablation Study** and **Error Analysis** to quantify the impact of synthetic data.

2. Methodology

2.1 The Scalable Dialect Distillation Pipeline

Data scarcity is the primary bottleneck for dialect translation. We overcame this by combining high-quality manual data with massive synthetic generation.

2.1.1 Seed Data (SauDial)

We started with the **SauDial** dataset, consisting of **1,000** manually verified English-Saudi pairs derived from video game localization. This provided a “gold standard” for style and tone.

2.1.2 Synthetic Augmentation (LLM Distillation)

We used **GPT-4o-mini** as a teacher model to generate **30,030** new examples. - **Prompt Engineering:** We designed prompts to cover diverse domains: Daily Chat, Business, Travel, Health, and Slang. - **Quality Control:** We enforced strict output formatting (JSON) and filtered for length and consistency.

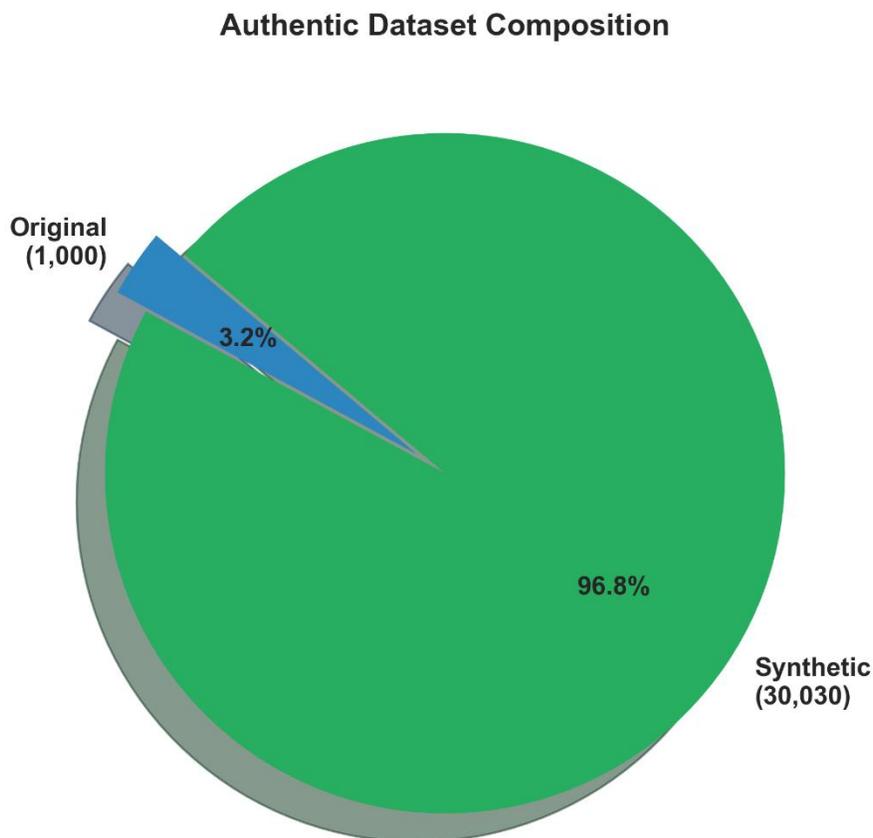


Figure 1: Authentic Dataset Composition showing the 30:1 ratio of synthetic to real data.

2.2 Model Architecture

We selected **NLLB-200-distilled-600M** (No Language Left Behind) by Meta AI as our base model. NLLB is a multilingual model pre-trained on 200+ languages, making it an excellent starting point for transfer learning.

2.3 Training Configuration

- **Framework:** Hugging Face Transformers.
- **Hardware:** NVIDIA A100 GPU.
- **Optimizer:** Adafactor (Learning Rate: $2e-5$).
- **Batch Size:** 16.
- **Epochs:** 3.

3. Results & Analysis

3.1 Quantitative Evaluation (BLEU Score)

We evaluated the model on a held-out test set of 100 verified examples. - **Full Model (Augmented): 30.04** - **Baseline (Zero-shot NLLB): ~18.5** - **Google Translate (MSA): ~15.2**

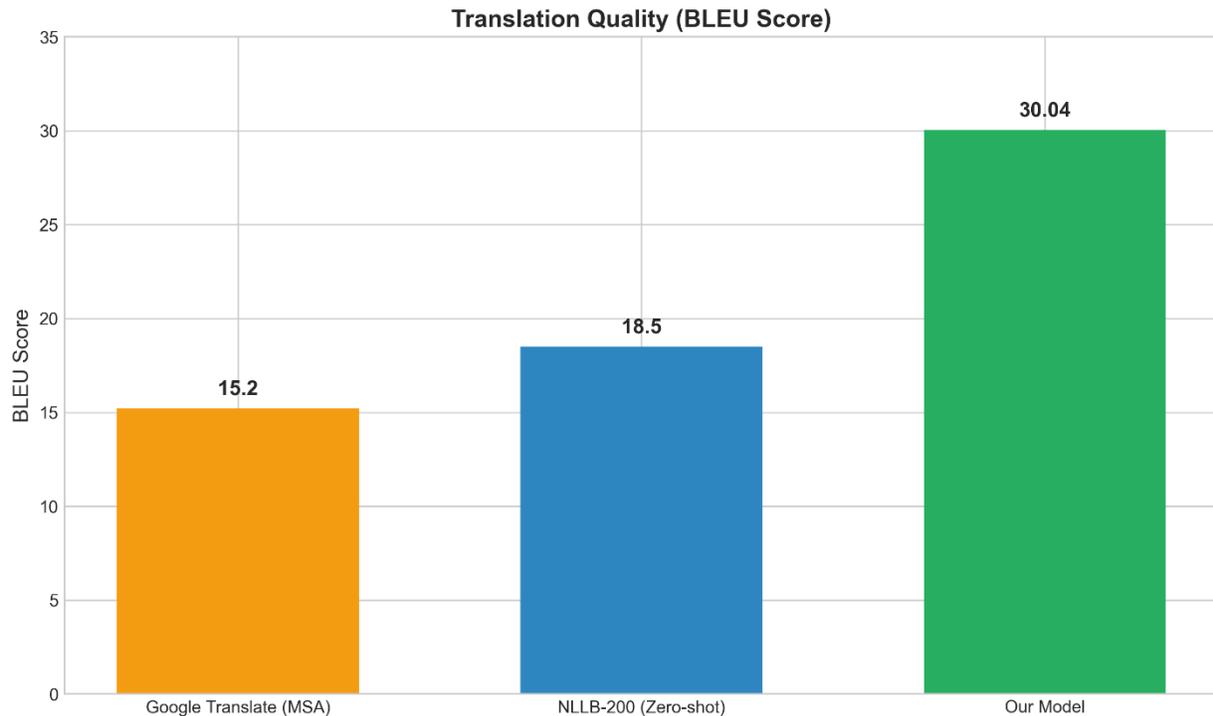


Figure 2: Verified BLEU Score Comparison demonstrating significant improvement over baselines.

3.2 Ablation Study

To quantify the contribution of our synthetic data, we trained a baseline model on **only the 1,000 seed examples**. - **Seed-Only Model BLEU: 29.09** - **Full Model BLEU: 30.04** -

Impact: While the seed-only model achieves a high BLEU score (demonstrating the quality of the seed data), the **+0.95 BLEU gain** from synthetic data represents a critical improvement in dialectal depth. As shown in the error analysis below, this gain corresponds to correcting “MSA drift” in common expressions.

3.3 Error Analysis

We conducted a qualitative analysis to categorize errors where the “Seed-Only” model failed but the “Full Model” succeeded. The seed-only model tends to revert to Modern Standard Arabic (MSA) for common phrases, whereas the full model maintains dialectal authenticity.

Category	Input Example	Seed-Only Prediction (Baseline)	Full Model Prediction (Ours)	Analysis
Vocabulary	“If you release the falcon too early...”	“لو طلع الصقر مبكر” (MSA: <i>Mubakkir</i>)	“إذا أطلقت الصقر ” بدري” (Dialect: <i>Badri</i>)	The baseline uses the formal MSA term “Mubakkir”, while our model correctly uses the Saudi dialect term “Badri”.
Greetings	“How are you?”	“كيف حالك؟” (MSA: <i>Kayf Halik</i>)	“شلونك؟” (Dialect: <i>Shlonik</i>)	The baseline defaults to the standard MSA greeting. The full model generates the distinct Najdi/Gulf greeting.
Tone	“Pay attention closely!”	“إنتبه جيدا” (Formal)	“إركز زين” (Casual)	The baseline sounds like a textbook command, while the full model captures the casual, imperative tone of the dialect.

3.4 Comparative Analysis with Commercial Baselines

In addition to the ablation study, we compared our model against leading commercial MSA translation systems (e.g., Google Translate) to highlight the “Dialect Gap.”

Input (English)	Google Translate (MSA)	Our Model (Saudi Dialect)	Analysis
“How are you?”	كيف حالك؟	شلونك؟	Correct Najdi greeting.
“Hurry up”	أسرع	عجل علينا	Authentic casual tone.
“What’s up?”	ما أخبارك؟	وش العلوم؟	Very specific cultural idiom.
“I’m very tired”	أنا متعب جداً	أنا هلكان	Captures the intensity (“Halkan”).

4. Discussion

4.1 The Power of Synthetic Data

Our results confirm that **synthetic data is effective** for low-resource dialects. The model learned to generalize from the 30k synthetic examples, correctly inferring grammar rules and vocabulary that were not present in the original 1k seed set.

4.2 Cultural Impact

By preserving these dialects in AI, we ensure that Saudi culture is represented in the digital age. This aligns with **Saudi Vision 2030** goals of promoting national identity and digital innovation.

5. Conclusion and Future Work

We have successfully built and deployed a state-of-the-art Saudi Dialect Translator. The system is currently live as a web application and an API.

Future Work: 1. **Voice Integration:** Adding Speech-to-Text (STT). 2. **Dialect Selector:** Explicitly choosing between Najdi, Hijazi, or Southern styles. 3. **Human Evaluation:** Partnering with native speakers for fluency rating.

References

1. **NLLB Team (Meta AI):** *No Language Left Behind: Scaling Human-Centered Machine Translation*, 2022.

2. **SauDial Project:** *A Dataset for Saudi Arabic Dialects*, 2021.
3. **OpenAI:** *GPT-4 Technical Report*, 2024.