

# LexiLingo AI

Tài Liệu Tham Khảo

References & Resources



Comprehensive Reference Guide

**Version:** 2.0

**Last Updated:** January 2026

**Author:** Nguyen Huu Thang

## Mục lục

<b>1</b>	<b>Language Models</b>	<b>2</b>
1.1	Large Language Models . . . . .	2
<b>2</b>	<b>Speech Models</b>	<b>2</b>
<b>3</b>	<b>Fine-tuning Techniques</b>	<b>2</b>
<b>4</b>	<b>Datasets</b>	<b>2</b>
4.1	Grammar Error Correction (GEC) . . . . .	2
4.2	English Learner Corpora . . . . .	3
4.3	CEFR & Vocabulary . . . . .	3
4.4	Dialogue & Conversation . . . . .	3
4.5	Pronunciation & Speech . . . . .	3
<b>5</b>	<b>Related Research Papers</b>	<b>4</b>
5.1	Grammar Error Correction . . . . .	4
5.2	Language Assessment . . . . .	4
5.3	Multi-task Learning . . . . .	4
<b>6</b>	<b>Tools &amp; Libraries</b>	<b>4</b>
<b>7</b>	<b>Vietnamese NLP Resources</b>	<b>5</b>
<b>8</b>	<b>Mobile Deployment</b>	<b>5</b>

## 1 Language Models

### 1.1 Large Language Models

Model	Paper/Link	Description
<b>Qwen2.5</b>	<a href="#">GitHub</a>   <a href="#">Paper</a>	Alibaba's multilingual LLM, strong on English & code
<b>LLaMA 3</b>	<a href="#">Meta AI</a>   <a href="#">Paper</a>	Meta's open LLM, excellent for fine-tuning

Bảng 1: Language Models

## 2 Speech Models

Model	Paper/Link	Description
<b>Whisper v3</b>	<a href="#">OpenAI</a>   <a href="#">Paper</a>	SOTA speech recognition, multilingual
<b>Faster-Whisper</b>	<a href="#">GitHub</a>	CTranslate2 optimized Whisper (4x faster)
<b>HuBERT</b>	<a href="#">Paper</a>   <a href="#">HuggingFace</a>	Self-supervised speech representation
<b>Piper TTS</b>	<a href="#">GitHub</a>	Fast, offline VITS-based TTS
<b>Silero VAD</b>	<a href="#">GitHub</a>	Voice Activity Detection

Bảng 2: Speech Models

## 3 Fine-tuning Techniques

Technique	Paper	Description
<b>LoRA</b>	<a href="#">arXiv:2106.09685</a>	Low-Rank Adaptation for efficient fine-tuning
<b>QLoRA</b>	<a href="#">arXiv:2305.14314</a>	Quantized LoRA, enables 65B models on single GPU
<b>PEFT</b>	<a href="#">GitHub</a>	HuggingFace Parameter-Efficient Fine-Tuning
<b>SFT</b>	<a href="#">TRL Docs</a>	Supervised Fine-Tuning Trainer

Bảng 3: Fine-tuning Techniques

## 4 Datasets

### 4.1 Grammar Error Correction (GEC)

Dataset	Link	Size	Description
<b>BEA-2019</b>	<a href="#">CodaLab</a>	34K	Write & Improve + LOCNESS shared task
<b>FCE Corpus</b>	<a href="#">Cambridge</a>	1,244	First Certificate in English exam

Dataset	Link	Size	Description
<b>CoNLL-2014</b>	<a href="#">NUS</a>	1,312	Grammatical Error Correction shared task
<b>JFLEG</b>	<a href="#">GitHub</a>	1,501	Fluency-focused GEC benchmark
<b>W&amp;I+LOCNESS</b>	<a href="#">HuggingFace</a>	34K	Combined Write&Improve + LOCNESS
<b>ERRANT</b>	<a href="#">GitHub</a>	Tool	Error annotation toolkit

Bảng 4: Grammar Error Correction Datasets

## 4.2 English Learner Corpora

Dataset	Link	Size	Description
<b>EFCAMDAT</b>	<a href="#">EF-Cambridge</a>	83M words	EF-Cambridge Open Language Database, CEFR-labeled
<b>TOEFL11</b>	<a href="#">ETS</a>	12,100	TOEFL essays with scores
<b>ICNALE</b>	<a href="#">Official</a>	10K	International Corpus Network of Asian Learners
<b>PELIC</b>	<a href="#">GitHub</a>	46K texts	Pitt English Language Institute Corpus

Bảng 5: English Learner Corpora

## 4.3 CEFR & Vocabulary

Dataset	Link	Description
<b>CEFR-J Wordlist</b>	<a href="#">Official</a>	Japanese CEFR word frequency lists
<b>English Profile</b>	<a href="#">Cambridge</a>	Official CEFR vocabulary lists
<b>Oxford 3000/5000</b>	<a href="#">Oxford</a>	Core vocabulary lists
<b>EVP</b>	<a href="#">Cambridge</a>	CEFR-graded vocabulary

Bảng 6: CEFR &amp; Vocabulary Resources

## 4.4 Dialogue & Conversation

Dataset	Link	Size	Description
<b>Intel/orca_dpo_pairs</b>	<a href="#">HuggingFace</a>	13K	DPO training pairs
<b>Anthropic HH-RLHF</b>	<a href="#">HuggingFace</a>	170K	Human preference data
<b>OpenAssistant</b>	<a href="#">HuggingFace</a>	161K	Multilingual conversation
<b>Tatoeba</b>	<a href="#">Official</a>	10M+	Parallel sentences, good for translation

Bảng 7: Dialogue &amp; Conversation Datasets

## 4.5 Pronunciation & Speech

Dataset	Link	Size	Description
TIMIT	LDC	6,300	Phoneme-aligned speech corpus
LibriSpeech	OpenSLR	1,000 hours	Clean speech for ASR
CommonVoice	Mozilla	17K+ hours	Multilingual speech corpus
L2-ARCTIC	GitHub	26 hours	Non-native English speech

Bảng 8: Pronunciation &amp; Speech Datasets

## 5 Related Research Papers

### 5.1 Grammar Error Correction

Paper	Year	Link	Key Contribution
GECToR	2020	arXiv:2005.12592	Efficient sequence tagging for GEC
T5 for GEC	2021	ACL	Transfer learning approach
GrammarT5	2022	arXiv:2203.07442	Grammar pre-training
LLM-GEC	2023	arXiv:2303.13648	LLMs for GEC

Bảng 9: Grammar Error Correction Papers

### 5.2 Language Assessment

Paper	Year	Link	Key Contribution
CEFR Classification	2018	ACL	Automatic CEFR level prediction
Automated Essay Scoring	2020	arXiv:2012.13958	BERT for essay scoring
Fluency Assessment	2021	Interspeech	Speech fluency metrics

Bảng 10: Language Assessment Papers

### 5.3 Multi-task Learning

Paper	Year	Link	Key Contribution
MT-DNN	2019	arXiv:1901.11504	Multi-Task Deep Neural Networks
UniLM	2020	NeurIPS	Unified Language Model pre-training
T5	2020	arXiv:1910.10683	Text-to-Text Transfer Transformer

Bảng 11: Multi-task Learning Papers

## 6 Tools & Libraries

Tool	Link	Purpose
<b>Transformers</b>	<a href="#">HuggingFace</a>	Model loading & inference
<b>PEFT</b>	<a href="#">HuggingFace</a>	Parameter-efficient fine-tuning
<b>TRL</b>	<a href="#">HuggingFace</a>	Transformer Reinforcement Learning
<b>BitsAndBytes</b>	<a href="#">GitHub</a>	8-bit/4-bit quantization
<b>vLLM</b>	<a href="#">GitHub</a>	High-throughput LLM serving
<b>ERRANT</b>	<a href="#">GitHub</a>	Error annotation toolkit
<b>Language Tool</b>	<a href="#">GitHub</a>	Rule-based grammar checking
<b>Sentence-Transformers</b>	<a href="#">GitHub</a>	Sentence embeddings

Bảng 12: Development Tools &amp; Libraries

## 7 Vietnamese NLP Resources

Resource	Link	Description
<b>VinAI PhoGPT</b>	<a href="#">GitHub</a>	Vietnamese generative model
<b>VietAI ViT5</b>	<a href="#">HuggingFace</a>	Vietnamese T5
<b>UIT-ViNewsQA</b>	<a href="#">GitHub</a>	Vietnamese QA dataset
<b>VLSP</b>	<a href="#">Official</a>	Vietnamese Language and Speech Processing

Bảng 13: Vietnamese NLP Resources

## 8 Mobile Deployment

Resource	Link	Description
<b>ONNX Runtime</b>	<a href="#">GitHub</a>	Cross-platform inference
<b>TensorFlow Lite</b>	<a href="#">TensorFlow</a>	Mobile deployment
<b>llama.cpp</b>	<a href="#">GitHub</a>	CPU inference for LLMs
<b>whisper.cpp</b>	<a href="#">GitHub</a>	CPU inference for Whisper
<b>MLC LLM</b>	<a href="#">GitHub</a>	Universal LLM deployment

Bảng 14: Mobile Deployment Resources