

Improving Yorùbá Diacritic Restoration

Iroko Fred Ọ̀nòmẹ̀ Orife¹ David I. Adélaní^{1,4} Timi Fasubaa¹
Victor Williamson⁵ Wuraola Fisayo Oyewusi³ Ọ́lámílékan Wahab¹ Kọ́lá Túbòsún²

¹Niger-Volta Language Technologies Institute,
²Yorùbá Name, ³Data Science Nigeria,
⁴Saarland University, ⁵University of Wisconsin-Milwaukee

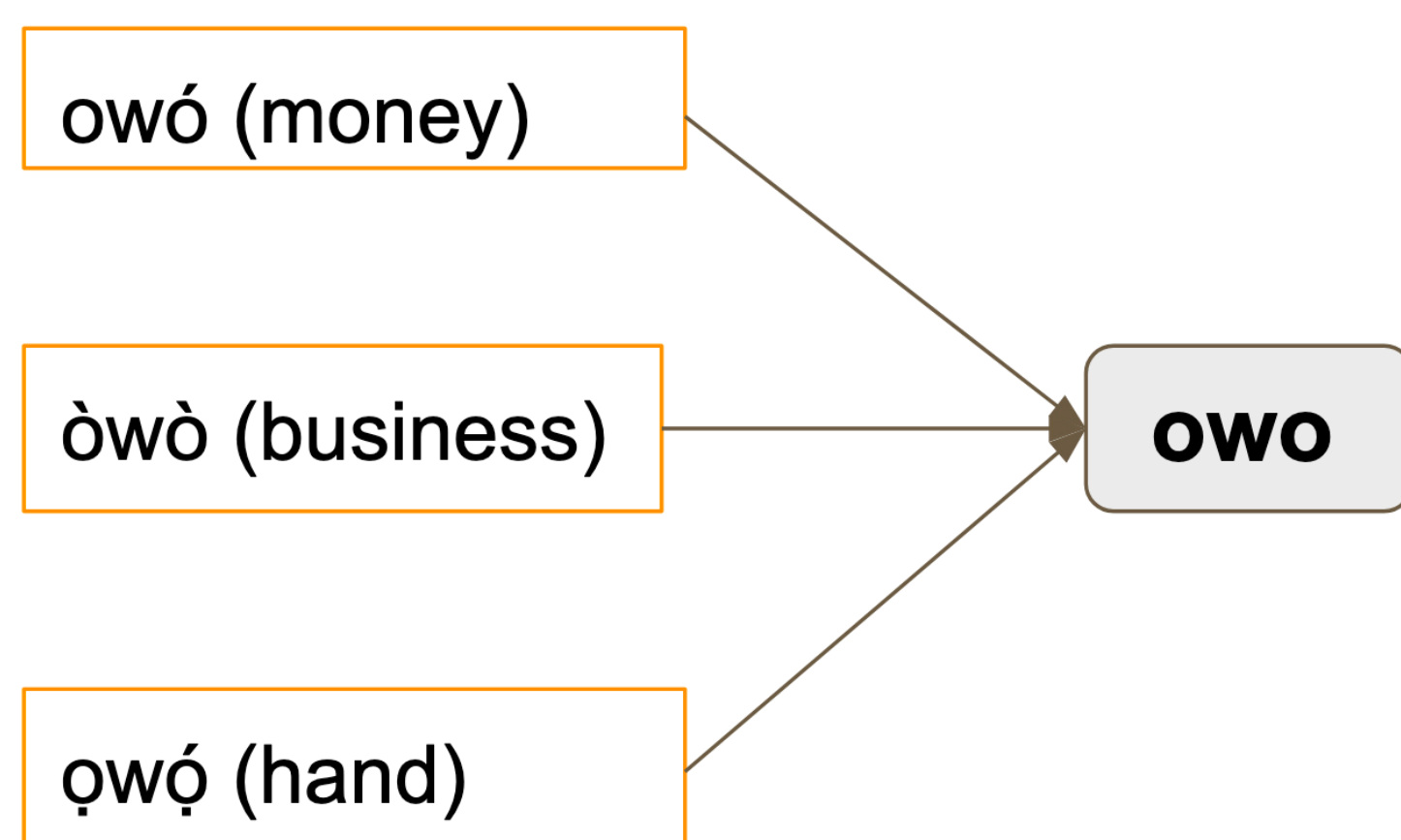
The Yorùbá Language

- 40M speakers in Nigeria & West Africa
- A **tonal** language
- Uses the Latin alphabet with diacritics
- Non-diacritized text is **useless** for computing

Yorùbá text must be correctly represented in computing environments, to build a robust ecosystem of Yorùbá-first language technologies including NLP, ASR & TTS applications.

Ambiguous non-diacritized text

Problem:



Automatic diacritic restoration (ADR) systems facilitate text entry and correction that promotes the correct orthography and quotidian usage of the language in electronic media.

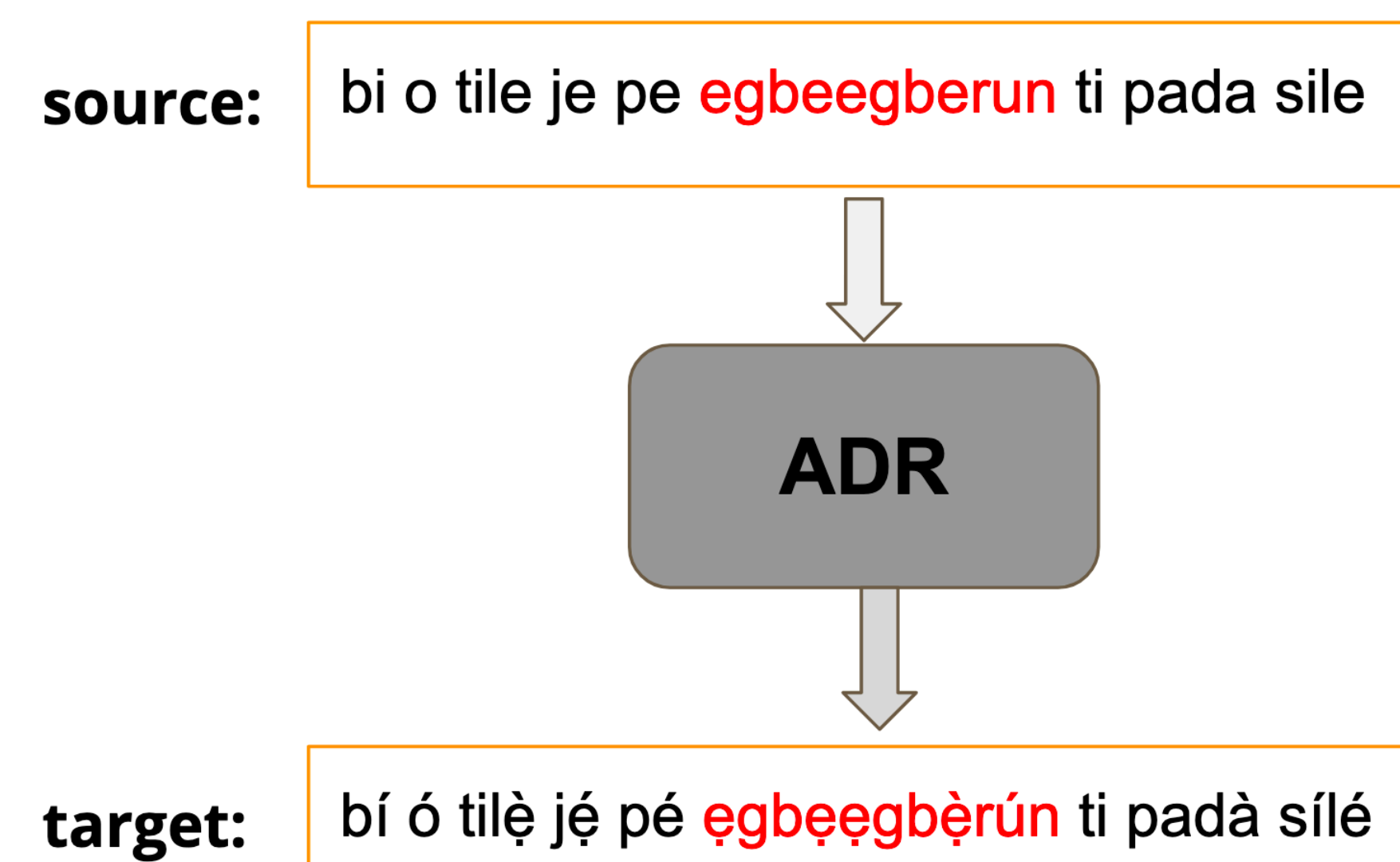
Datasets

Data sources, prevalence and category of text

# words	Source	Description
24,868	rma.nwu.ac.za	Lagos-NWU corpus
50,202	theyorubablog.com	language blog
910,401	bible.com/versions/911	Biblica (NIV)
11,488,825	opus.nlpl.eu	JW300
831,820	bible.com/versions/207	KJV
177,675	GitHub	Embeddings dataset
142,991	GitHub	Language ID corpus
47,195		Yorùbá lexicon
29,338	yoruba.unl.edu	Proverbs
2,887	unicode.org/udhr	Human rights edict
150,360	Private sources	Conversations
15,281	Private sources	Short stories
20,038	OCR	Háà Ènìyàn (Fiction)
28,308	yo.globalvoices.org	Global Voices news

Approach

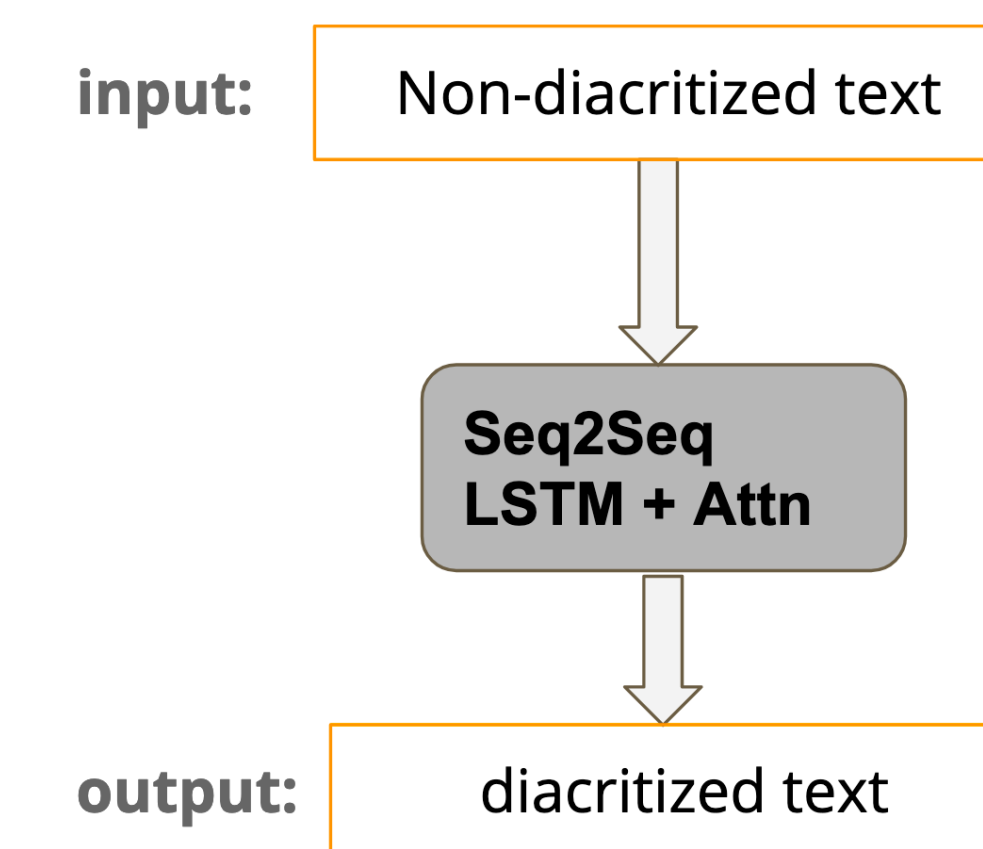
Training idea: Predict diacritics based on context



This Work

Baseline [Orife, 2018]

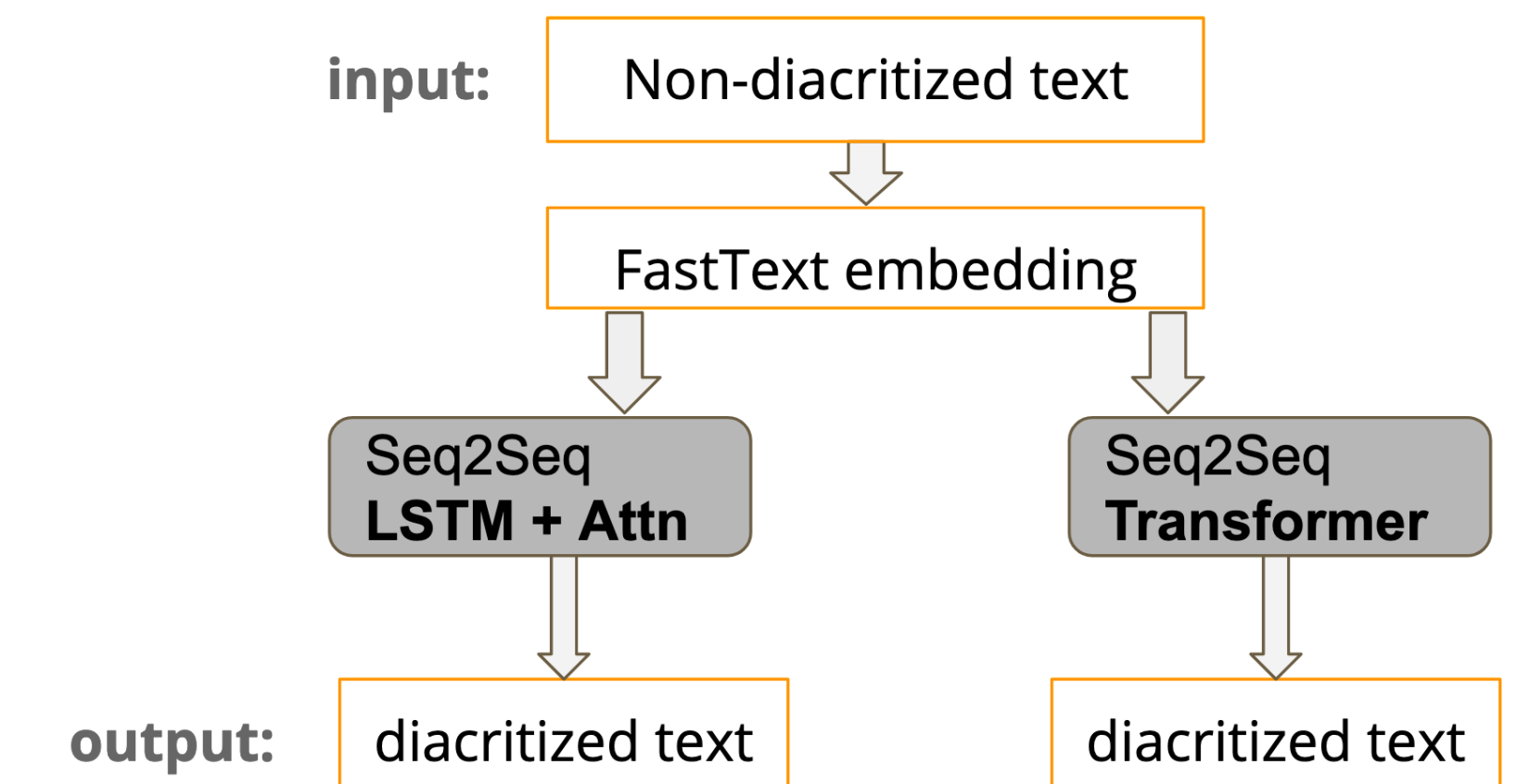
corpus: Bible + LagosNWU+ blog



Improved ADR:

corpus: Bible + LagosNWU+ blog + JW300+interviews + ...

Embeddings: FastText Embeddings [Alabi et al 2020]



Results

BLEU scores, predicted perplexity & WER on the Global Voices testset

Model	BLEU	Perplexity	WER%
Soft-attention model from (Orife, 2018)	26.53	1.34	58.17
+ Language ID corpus	42.52	1.69	33.03
++ Interview text	42.23	1.59	32.58
+ All new text <i>minus</i> JW300	43.39	1.60	31.87
+ All new text	59.55	1.44	20.40
+ All new text with FastText embedding	58.87	1.39	21.33
Transformer model			
+ All new text <i>minus</i> JW300	45.68	1.95	34.40
+ All new text	59.05	1.40	23.10
+ All new text + FastText embedding	59.80	1.43	22.42



Paper (arXiv)
2003.10564



Code (Github)
Niger-Volta-LTI/yoruba-adr