

HOW EFFECTIVE ARE AI MODELS IN TRANSLATING ENGLISH SCIENTIFIC TEXTS TO NIGERIAN PIDGIN: A LOW-RESOURCE LANGUAGE?

Flora Oladipupo, Anthony Soronnadi, Ife Adebara, & Olubayo Adekanmbi

Research & Innovation Department

Data Science Nigeria

Lagos, Nigeria

{flora, anthony, ife, olubayo}@datasciencenigeria.ai

ABSTRACT

This research explores the challenges and limitations of applying deep learning models to the translation of scientific texts from English to Nigerian Pidgin, a widely spoken but low-resource language in West Africa. Despite advancements in machine translation, translating domain-specific content such as biological research papers presents unique obstacles, including data scarcity, linguistic complexity, and model generalization issues. We investigate the performance of AI models, including Pidgin-UNMT, mt5-base model, AfriTeVa base, Afri-mt5 base model and GPT 4.0 model through a comparative analysis using BLEU scores, CHRF, TER, Africomet metrics on a newly created Eng-PidginBioData dataset of biological texts. Our findings reveal significant gaps in model performance, emphasizing the need for more domain-specific fine-tuning, improved dataset creation, and collaboration with native speakers to enhance translation accuracy. By presenting real-world challenges encountered in applying deep learning to low-resource languages this research suggests strategies to overcome these barriers. Our study provides valuable insights into the persistent challenges faced by AI-driven translation systems, from limited data to domain mismatches, and highlights ways to enhance their effectiveness for underrepresented languages. By addressing these constraints, we offer actionable strategies for more inclusive and impactful scientific knowledge dissemination.

1 INTRODUCTION

Nigeria Pidgin, a Creole language spoken widely by 75 million people in West Africa Kasraee (2017), often serves as a bridge between different language speakers. It is also the 4th most spoken language in Nigeria according to Statista (2021). Despite its widespread use, scientific literature in Nigerian Pidgin remains underrepresented, limiting access to critical knowledge by speakers of this language. The translation of scientific literature is crucial for the dissemination of knowledge across linguistic and cultural boundaries. However, scientific literature in Nigerian Pidgin remains an underexplored area in AI research and academics. This study focuses on development of a translation system for Nigerian Pidgin language based on biological research papers, shedding light on how deep learning struggles to adapt to domain-specific and low-resource contexts.

Despite state-of-the-art performance in high-resource settings, AI-driven machine translation systems frequently fall short in low-resource scenarios. Data scarcity, contextual nuances, and scientific terminology complexity combine to limit the performance of existing Neural Machine Translation (NMT) models. In this paper, we explore machine translation in the domain of biological research texts. Our contributions are as follows:

- **Dataset Creation:** We present **Eng-PidginBioData**, a domain-specific high quality English-to-Pidgin parallel corpus focusing on biological research texts.

- **Comparative Analysis:** We evaluate **mt5 base**, **AfriTeVa base**, **Afri-mt5 base**, **Pidgin UNMT** and **GPT 4.0** models using BLEU, TER, CHRF, Africomet scores metrics, highlighting performance gaps and unexpected failures.
- **Insights & Strategies:** We identify real-world challenges (limited data, cultural nuances, scientific jargon) and propose practical interventions (domain adaptation, expanded corpora, native speaker collaboration) to narrow the performance gap.

2 RELATED WORKS

Incorporating local languages into science education enhances comprehension, engagement, and equity Babaci-Wilite (2016). Despite progress in Neural Machine Translation (NMT), translation models continue to face challenges in low-resource settings such as Nigerian Pidgin (pcm) due to the scarcity of large-scale parallel corpora Nwafor (2022). A major breakthrough came with Pidgin-UNMT Ogueji & Ahia (2019), which leveraged unsupervised learning techniques to train an NMT system without parallel data, laying the foundation for Creole language translation research.

The Afri-mT5 model Adelani et al. (2022), an adaptation of mT5 for African languages, demonstrated limited performance on Nigerian Pidgin, primarily due to the dominance of French in its training data. While studies suggest that fine-tuning with in-domain data enhances translation quality, Afri-mT5 remains constrained in technical domains. AfriTeVa Oladipo et al. (2023) advanced African NLP but underperformed in Nigerian Pidgin MT, largely due to the absence of a dedicated Pidgin corpus.

Cheetah Adebara et al. (2024), which supports 517 African languages, achieved a BLEU score of 32.64 for English-to-Pidgin MT but lacked domain-specific optimization. Its successor, Toucan Elmadany et al. (2024), was fine-tuned on AfroLingu-MT, Africa’s largest MT benchmark. While effective for general Nigerian Pidgin translation, it still requires fine-tuning for specialized domains such as biomedical, technical, and legal texts.

3 METHODOLOGY

3.1 DATASET

The Nigerian Pidgin dataset used in this study was meticulously curated by scraping open-sourced English biological research papers. To maintain privacy, the authors’ names were removed from the scraped data through an anonymization process. This corpus was then segmented into 2,300 sentences, forming the foundational dataset for translation. These sentences were then translated manually into Nigerian Pidgin. To ensure linguistic accuracy and cultural relevance, these translations underwent a manual correction process by human annotators proficient in Nigerian Pidgin. This dual-step translation process aimed to refine the quality of the dataset, making it a reliable benchmark for evaluating machine translation models. This dataset is called **Eng-PidginBioData**.

Split	Size	TTR (English)	TTR (Pidgin)
Train	1380	16.6378	11.9739
Dev	460	27.5226	19.1845
Test	460	28.4705	20.0848

Table 1: Type-Token Ratio (TTR) for English and Pidgin on Eng-PidginBioData

To assess the linguistic richness and diversity of the curated dataset Eng-PidginBioData, we computed the Type-Token Ratio (TTR) for both English and Nigerian Pidgin text across the training, development, and test splits. The results show that TTR values are consistently lower in Nigerian Pidgin compared to English, reflecting Pidgin’s characteristic reliance on a smaller, more flexible vocabulary with frequent code-switching and word reuse. Specifically, in the training set, TTR for English is 16.64, whereas Nigerian Pidgin has a lower TTR of 11.97, suggesting that Pidgin exhibits more lexical repetition. Similarly, in the development and test sets, TTR values for English (27.52 and 28.47) remain higher than for Pidgin (19.18 and 20.08), reinforcing this pattern.

3.2 MODEL FINE-TUNING

The Eng-PidginBioData dataset was used to fine-tune existing machine translation models. The models includes: **Pidgin-UNMT 220M parameters**, Ogueji & Ahia (2019), **mt5 base 580M parameters** Raffel et al. (2023), **AfriTeVa base 229M parameters** Oladipo et al. (2023), **Afri-mt5 base 580M parameters** Adelani et al. (2022)

A uniform set of hyperparameters was used to ensure consistency across all models. Specifically, a learning rate of $2e-5$ was employed with AdamW optimization, and a batch size of 4 was used. Additionally, a maximum sequence length of 256 tokens was maintained throughout training. These hyperparameter choices align with findings from Nag et al. (2024), who demonstrated the effectiveness of low learning rates in fine-tuning transformer models for low-resource languages.

3.3 ANALYSIS

Model	BLEU		TER		CHRF		AfriComet	
	(pcm → en)	(en → pcm)	(pcm → en)	(en → pcm)	(pcm → en)	(en → pcm)	(pcm → en)	(en → pcm)
Afri-mt5-base	34.43	23.37	43.30	66.09	63.78	43.46	0.56	0.42
AfriTeVa-base	18.47	25.55	63.40	53.03	50.84	59.05	0.41	0.51
Pidgin-UNMT-base	52.10	30.34	22.02	41.36	82.04	67.29	0.73	0.70
mt5-base	38.90	35.02	43.28	37.88	63.95	66.72	0.56	0.62

Table 2: Translation performance across different models. Nigerian Pidgin (pcm), English (eng).

Model	BLEU	TER	CHRF	AfriComet
GPT-4.0	37.64	51.36	61.84	0.72

Table 3: Evaluation of GPT-4.0 on the Eng-PidginBioData dataset.

The table presents the evaluation of various machine translation models on the Eng-PidginBioData dataset, incorporating BLEU, Translation Error Rate (TER), CHRF, and AfriComet metrics to assess their translation accuracy, fluency, and semantic preservation. The results highlight significant disparities in performance across the evaluated models, particularly in their ability to translate between English and Nigerian Pidgin. Pidgin-UNMT emerged as the most effective model, achieving the highest BLEU score (52.10) for Pidgin-to-English (pcm-eng) translation, the lowest TER (22.02), and the highest CHRF (82.04). These results confirm that Pidgin-UNMT is best suited for translating Nigerian Pidgin to English, producing fluent, accurate, and semantically aligned translations with minimal errors. However, for English-to-Pidgin translation, Pidgin-UNMT’s performance dropped significantly (BLEU: 36.98, TER: 41.36), indicating difficulties in generating high-quality Nigerian Pidgin text from English scientific content. mT5-base demonstrated a relatively balanced performance across both translation directions, scoring BLEU 38.90 and 35.02. While it performed better than Afri-mT5 and AfriTeVa, its high TER values (43.28 pcm → en, 37.88 en → pcm) suggest that its translations require substantial post-editing. The model’s CHRF and AfriComet scores indicate that it captures character-level fluency but struggles with precise meaning preservation, particularly for domain-specific terms. Afri-mT5-base exhibited the weakest performance overall, particularly in English-to-Pidgin translation, with a BLEU score of just 23.37 and the highest TER (66.09), indicating a severe misalignment with reference translations. The CHRF score (43.46) and AfriComet score (0.42) further confirm its inability to generate meaningful Nigerian Pidgin translations. Despite being an adaptation of mT5 for African languages, Afri-mT5 still underperformed significantly, suggesting that it lacks sufficient exposure to Nigerian Pidgin training data and does not fully capture the linguistic structure of Nigerian Pidgin. AfriTeVa, while slightly better than Afri-mT5, still struggled with fluency and accuracy in Nigerian Pidgin translation. It achieved a BLEU score of 25.55 for en → pcm, significantly lower than Pidgin-UNMT (36.98) and mT5-base (35.02). However, it

outperformed Afri-mT5 in CHRF (59.05 vs. 43.46) and had a lower TER (53.03 vs. 66.09), suggesting that its broader exposure to African languages provided some advantages, though it remains suboptimal for Nigerian Pidgin translation. The TER values across all models for en \rightarrow pcm confirm that translating into Nigerian Pidgin is significantly more challenging than translating into English. The flexible, informal, and code-mixed nature of Nigerian Pidgin likely contributes to this difficulty, as models struggle with generating fluent and contextually appropriate translations. Future improvements should focus on fine-tuning models with larger, domain-specific corpora and exploring hybrid approaches that combine the strengths of Pidgin-UNMT with the adaptability of multilingual models like mT5-base. Based on the results in Table 3, GPT-4.0 demonstrates greater fluency than other models on the Eng-PidginBioData dataset, as reflected in its high BLEU score. However, its elevated TER score highlights challenges in scientific accuracy, necessitating extensive post-editing. While the model effectively preserves general meaning, its translation of technical terms requires improvement. Future research should focus on fine-tuning GPT-4.0 with a larger, domain-specific Nigerian Pidgin corpus to enhance its precision. A hybrid approach that integrates GPT-4.0 with a domain-adapted model could help address accuracy gaps. Additionally, incorporating human-in-the-loop evaluation and post-editing workflows would further enhance translation quality, ensuring that scientific texts remain both accessible and reliable for Nigerian Pidgin speakers.

3.4 CHALLENGES AND INSIGHTS

One of the major challenges encountered in this study was the lack of large-scale, high-quality parallel corpora for Nigerian Pidgin, particularly in the biological domains. The absence of well-structured datasets made it necessary to curate domain-specific data manually, a process that involved scraping, translating, and manually editing scientific texts. This data limitation directly impacted the fine-tuning efficiency of the models, as they lacked extensive Nigerian Pidgin text exposure during pretraining. A key observation from this study is the significant disparity in performance between Pidgin-specific models and adapted multilingual models. AfriMT5, despite being an African language adaptation of mmt5, struggled with fluency and accuracy, achieving a high TER (66.09) and lower CHRF (43.46) than AfriTeVa. These results suggest that African multilingual models require additional domain adaptation to effectively translate scientific texts into Nigerian Pidgin. Another observation during evaluation was that Nigerian Pidgin translations frequently mixed English words where no direct Pidgin equivalent existed. While this is common in Nigerian Pidgin usage, it posed a challenge for evaluation metrics such as BLEU, which struggle with code-mixed sentences. The lack of a standardized Nigerian Pidgin orthography further complicated model evaluation, as multiple valid translations could exist for the same phrase, leading to inconsistencies in scoring. Future models should incorporate context-aware embeddings and specialized tokenization strategies to handle code-switching more effectively. Also, while automatic metrics provided insights into model performance, the study lacked real-time human feedback loops to assess fluency, coherence, and cultural relevance. Incorporating human evaluators for post-editing and ranking translations could further refine the models. To enhance translation quality, future research should focus on expanding the dataset with domain-specific glossaries, fine-tuning models with a larger Nigerian Pidgin scientific corpus, and incorporating transfer learning techniques. Also, introducing human-in-the-loop translation validation for refining Nigerian Pidgin translations in real-world applications could further improve scientific text translation accuracy, ensuring that Pidgin speakers can access critical information in education, and research.

3.5 CONCLUSION

This study assessed fine-tuned machine translation models for Nigerian Pidgin in the scientific domain, highlighting the strengths and limitations of multilingual and African language-adapted models. Pidgin-UNMT emerged as the most effective for pcm-eng (BLEU: 30.34, AfriComet: 0.71), while mmt5 performed poorly (BLEU: 0.03), emphasizing the need for targeted fine-tuning on low-resource languages. While AfriMT5 and AfriTeVa showed improvements through adaptation, they struggled with fluency, semantic adequacy, and technical terms. High TER scores and evaluation challenges from code-switching and non-standardized orthography further highlight the limitations of current automatic metrics. Future research should explore hybrid models, context-aware embeddings, and expanded datasets to enhance scientific Nigerian Pidgin translation. Addressing these challenges will improve machine translation for scientific literacy, education, and accessibility for Nigerian Pidgin speakers.

REFERENCES

- Ife Adebara, AbdelRahim Elmadany, and Muhammad Abdul-Mageed. Cheetah: Natural language generation for 517 African languages. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12798–12823, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.691>.
- David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajuddeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulmumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. A few thousand translations go a long way! leveraging pre-trained models for African news translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3053–3070, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.223. URL <https://aclanthology.org/2022.naacl-main.223>.
- Ammar, G. Mulcaire, Y. Tsvetkov, G. Lample, C. Dyer, and N. A. Smith. Massively multilingual word embeddings. *arXiv preprint*, 2016.
- Artetxe, G. Labaka, and E. Agirre. Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2289–2294, 2016.
- Artetxe, G. Labaka, and E. Agirre. Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 451–462, 2017a.
- M. Artetxe, G. Labaka, E. Agirre, and K. Cho. Unsupervised neural machine translation. *arXiv preprint*, 2017b.
- Zehila Babaci-WilHITE. The use of local languages for effective science literacy as a human right. In *Human rights in language and STEM education*, 2016.
- Bojanowski, E. Grave, A. Joulin, and T. Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- Conneau, G. Lample, M. A. Ranzato, L. Denoyer, and H. Jégou. Word translation without parallel data. *arXiv preprint*, 2017.
- Abdelrahim Elmadany, Ife Adebara, and Muhammad Abdul-Mageed. Toucan: Many-to-many translation for 150 african language pairs. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 13189–13206, 2024.
- Joulin, P. Bojanowski, T. Mikolov, H. Jégou, and E. Grave. Loss in translation: Learning bilingual word mapping with a retrieval criterion. *arXiv preprint*, 2018.
- Najiba Kasraee. Bbc blogs - academy - working towards a standard pidgin, 2017. URL <https://www.bbc.co.uk/blogs/academy/entries/70f2a30c-40a5-463c-9480-1d63e7d5f44a>. Accessed: 2024-05-21.
- G. Lample, A. Conneau, L. Denoyer, and M. A. Ranzato. Unsupervised machine translation using monolingual corpora only. *arXiv preprint*, 2017.
- G. Lample, M. Ott, A. Conneau, L. Denoyer, and M. A. Ranzato. Phrase-based & neural unsupervised machine translation. *arXiv preprint*, 2018.

- Mikolov, Q. V. Le, and I. Sutskever. Exploiting similarities among languages for machine translation. *arXiv preprint*, 2013a.
- Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pp. 3111–3119, 2013b.
- Arijit Nag, Soumen Chakrabarti, Animesh Mukherjee, and Niloy Ganguly. Efficient continual pre-training of llms for low-resource languages, 2024. URL <https://arxiv.org/abs/2412.10244>.
- Ebelechukwu Nwafor. A survey of machine translation tasks on nigerian languages. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 6480–6486, 2022. URL <https://aclanthology.org/2022.lrec-1.695>.
- Kelechi Ogueji and Orevaoghene Ahia. Pidginunmt: Unsupervised neural machine translation from west african pidgin to english. 2019.
- Akintunde Oladipo, Mofetoluwa Adeyemi, Orevaoghene Ahia, Abraham Owodunni, Odunayo Ogundepo, David Adelani, and Jimmy Lin. Better quality pre-training data and t5 models for African languages. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 158–168, Singapore, December 2023. Association for Computational Linguistics. URL <https://aclanthology.org/2023.emnlp-main.11>.
- Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543, 2014.
- Peter and H. G. Wolf. A comparison of the varieties of west african pidgin english. *World Englishes*, 26(1):3–21, 2007.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL <https://arxiv.org/abs/1910.10683>.
- Ravi and K. Knight. Deciphering foreign language. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 12–21, 2011.
- Ruder, I. Vulić, and A. Søgaard. A survey of cross-lingual word embedding models. *arXiv preprint*, 2017.
- S. L. Smith, D. H. Turban, S. Hamblin, and N. Y. Hammerla. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *arXiv preprint*, 2017.
- Statista. Population in nigeria by languages spoken, 2021. URL <https://www.statista.com/statistics/1285383/population-in-nigeria-by-languages-spoken/>. Accessed: Feb 4, 2025.