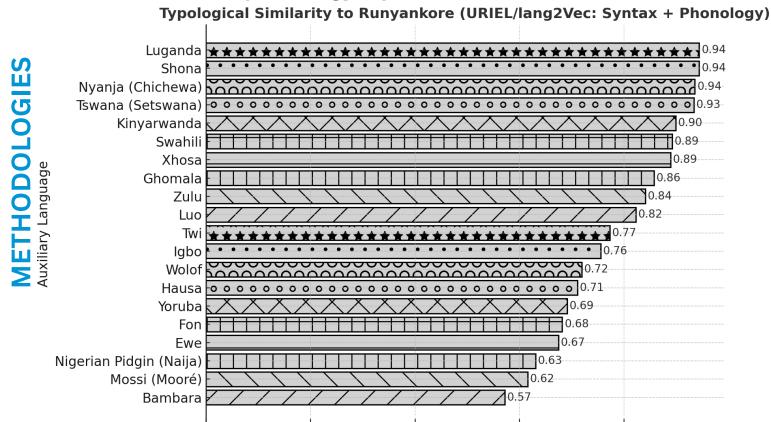
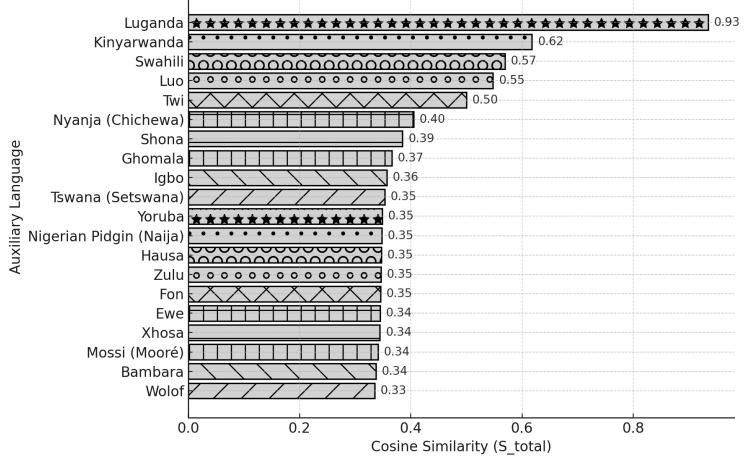
Cross-Lingual Adaptation For Named Entity Recognition in Runyankore

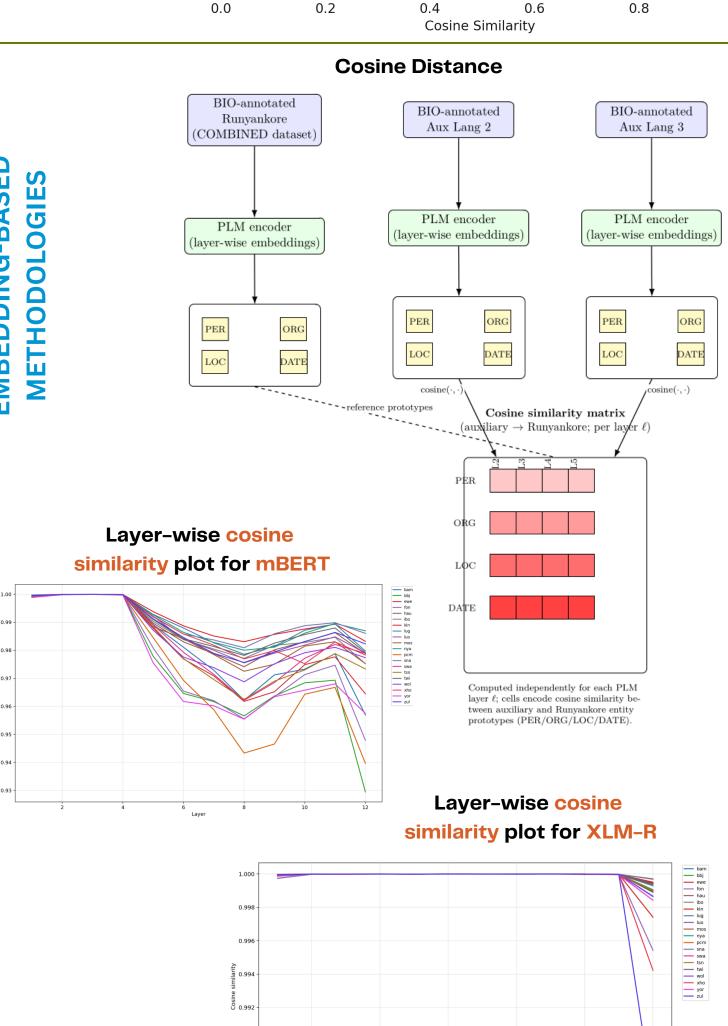
Our Corpus nyn Runyankore [TARGET] Bantu lug Luganda Bantu kin Kinyarwanda Bantu swa Swahili luo Luo (Dholuo) **Nilotic** Kwa twi Twi nya Nyanja (Chichewa) **Bantu** sna Shona Bantu **Grassfields Bantu** bbj Ghomala Volta-Niger ibo Igbo tsn Tswana (Setswana) Bantu yor Yoruba Volta-Niger pcm Nigerian Pidgin (Naija) Chadic hau Hausa **Bantu** zul isiZulu Kwa fon Fon Kwa ewe Ewe **Bantu** xho isiXhosa mos Mossi (Mooré) Gur bam Bambara Mande wol Wolof Atlantic-Congo **Tools and Resources Hugging Face** XLM-RoBERTa, Multilingual BERT (mBERT), AFROXLM-R ANNOTATION TOOL: doccan

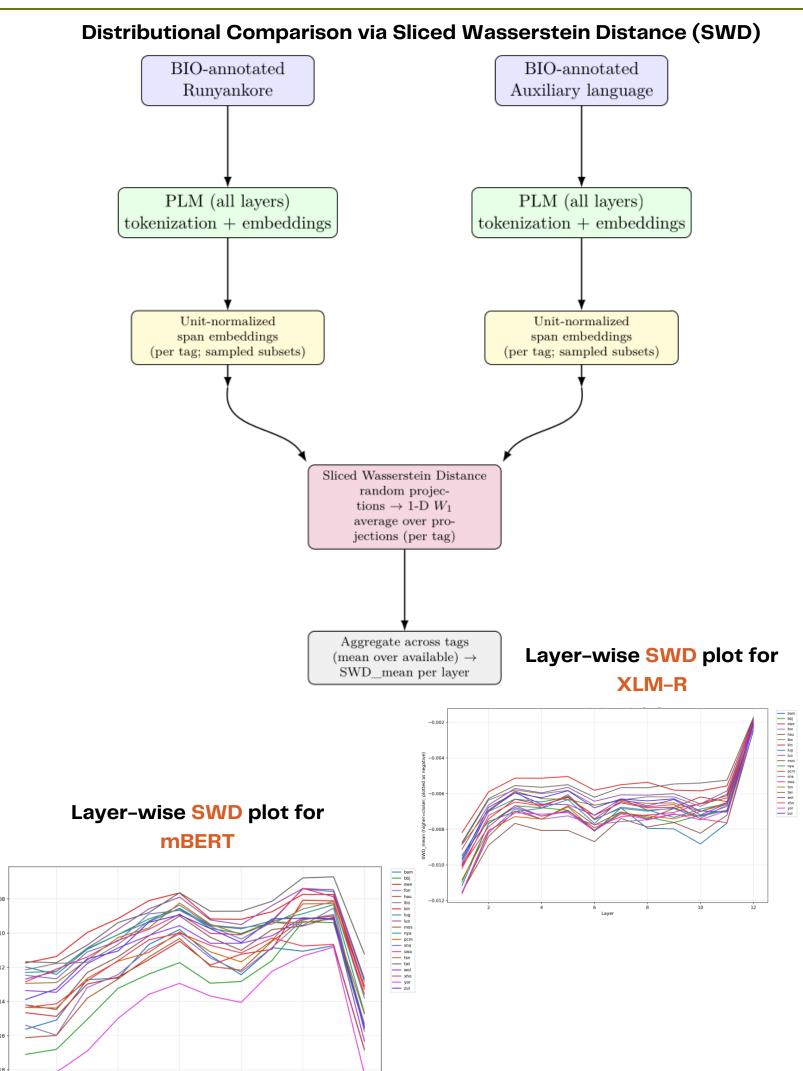
Overall typological similarity to Runyankore computed from URIEL/lang2Vec typological feature vectors combining syntax and phonology representations.



Overall typological similarity to Runyankore (S_total) computed from $\underline{\text{LinguaMeta's}} \text{ script match, shared country/region locale overlap, and}$ geographic proximity based on latitude/longitude distance (τ =1000 KM).







Introduction / Motivation

DATASETS:

Runyankore datasets:

Languages (MPTC)

Auxiliary languages: MasakhaNER 2.0

1. Sunbird African Language Technology (SALT)

2. Multilingual Parallel Text Corpora for East African

African languages remain underrepresented in NLP, with most multilingual models performing best on high-resource languages. **Runyankore**, a major Bantu language of **Uganda**, lacks any NER dataset, excluding it from key benchmarks like MasakhaNER 2.0.

This work develops the first annotated

Runyankore NER dataset and investigates crosslingual transfer from related African languages.
Using both typological features
(URIEL/Lang2Vec, LinguaMeta) and embeddingbased similarity measures (cosine distance,
Sliced Wasserstein Distance (SWD)) on the
DATE, PERSON, LOCATION, and
ORGANIZATION entities, the study identifies
which linguistic relationships best support

transfer learning.

Dataset Creation

The Runyankore NER corpus was annotated with a semiautomated approach, where XLM-R was trained with Luganda NER data, and used to identify initial entity suggestions in the Runyankore dataset (Zero-shot transfer), following MasakhaNER 2.0 guidelines, covering PERSON, ORG, LOC, DATE, and O for non-entity types. Over 5,000 sentences were double-annotated and cross-checked for quality assurance.

Example (BIO format):

Besigye B-PER Entity distribution across SALT, MPTC, and COMBINED datasets. akeetaba O Columns B and I represent Beginning and Inside tags, respectively. omu O Entity Type Train \mathbf{Dev} Test Dataset mirimo O DATE 328y'eby'obutegyeki O LOC 620omuri O SALT ORG 136191215Uganda B-LOC PER 113omu O DATE LOC 33 mwaka B-DATE MPTC ORG gwa I-DATE PER rukumi I-DATE DATE 534 437 - 588rwenda I-DATE 714 175LOC 185COMBINED ORG 478154245 $161 \quad 232$ nshanju I-DATE PER 200 107 87 43 114 56

Methodology

This study combines data creation, language similarity analysis, and cross-lingual NER experiments to evaluate transfer into Runyankore.

- Dataset Development: A new Runyankore NER corpus was annotated using MasakhaNER 2.0 guidelines, covering Person, Organization, Location, and Date entities, initially identified in raw Runyankore data by applying a zero-shot transfer from an XLM-R PLM trained on Luganda NER data, then manually verified and corrected by a human.
- Model Evaluation: Three multilingual pre-trained models: mBERT, XLM-R, and AfroXLM-R, were tested under zero-shot and cross-lingual setups.
- Typological Similarity: Computed from URIEL/lang2Vec and LinguaMeta features, focusing on syntax, phonology, and geographic proximity to identify linguistically related source languages.
- Embedding-Based Similarity: Calculated from contextual embeddings using cosine distance and the Sliced Wasserstein Distance (SWD) to capture representational closeness across languages.
- Transfer Comparison: The two similarity strategies: typological vs embedding-based guided auxiliary language selection to determine which better predicts cross-lingual transfer effectiveness.

Together, these techniques provide a structured framework for analyzing how linguistic and representation-level similarities influence the success of transfer learning for low-resource African languages.

Research Objectives

isiNdebelesh

This study investigates two key questions:

(1) How do multilingual model pretraining choices influence NER performance in Runyankore?

(2) Which language similarity strategy, typological or embedding-based, more effectively guides auxiliary language selection for cross-lingual transfer?

Findings indicate that embedding-based similarity shows a modest but consistent correlation with improved transfer, particularly among closely related Bantu languages such as Luganda and Kinyarwanda.

Overall, embedding-driven selection achieved slight yet measurable gains over typology-based approaches.

Conclusion & Future Work

This study introduced the first Runyankore NER dataset and demonstrated the advantages of embedding-based auxiliary language selection for cross-lingual transfer over typological-based selection.

Future work will extend to other Runyakitara languages: Rukiga, Runyoro, and Rutoro, and integrate findings into MasakhaNER pipelines to strengthen African-language representation in multilingual NLP.

Author: Prosper ARINEITWE ASIIMWE
Msc. Computer Science (Artificial Intelligence)
Supervisors: Dr. Jan BUYS

Dr. Francois MEYER

arnari002@myuct.ac.za

jan.buys@uct.ac.za francois.meyer@uct.ac.za

