# KenCorpus: Kenyan Languages Corpus

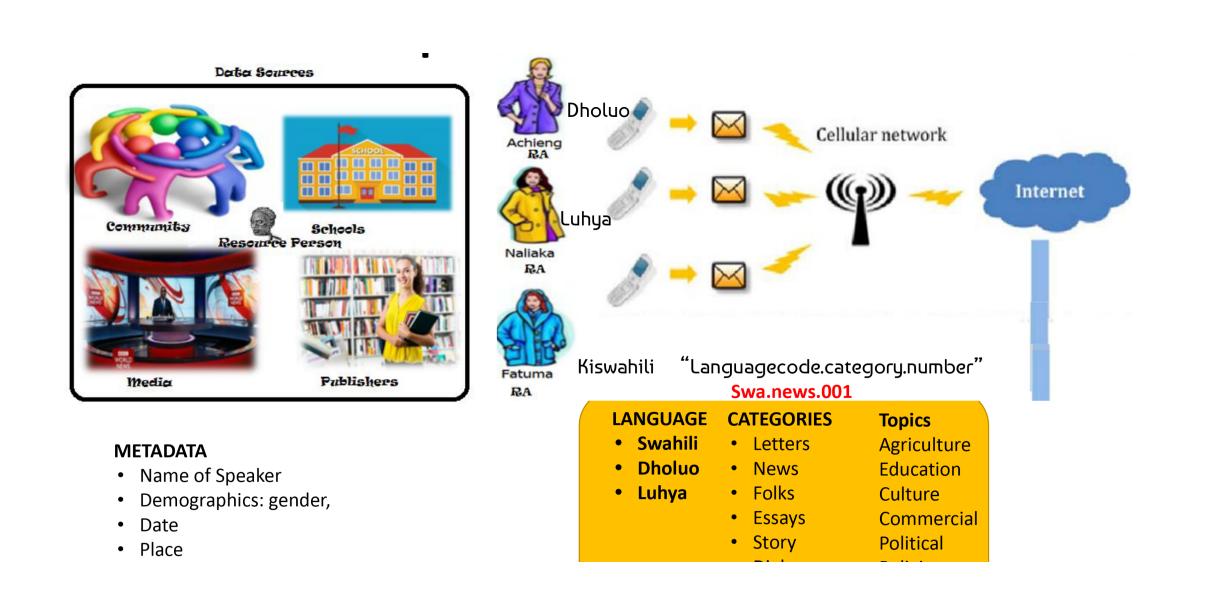
**Lilian Wanzare**, Edward Ombui, Lawrence Muchemi, Barack Wanjawa, Owen McOnyango and Florence Indede

Maseno University, University of Nairobi, Africa Nazarene University

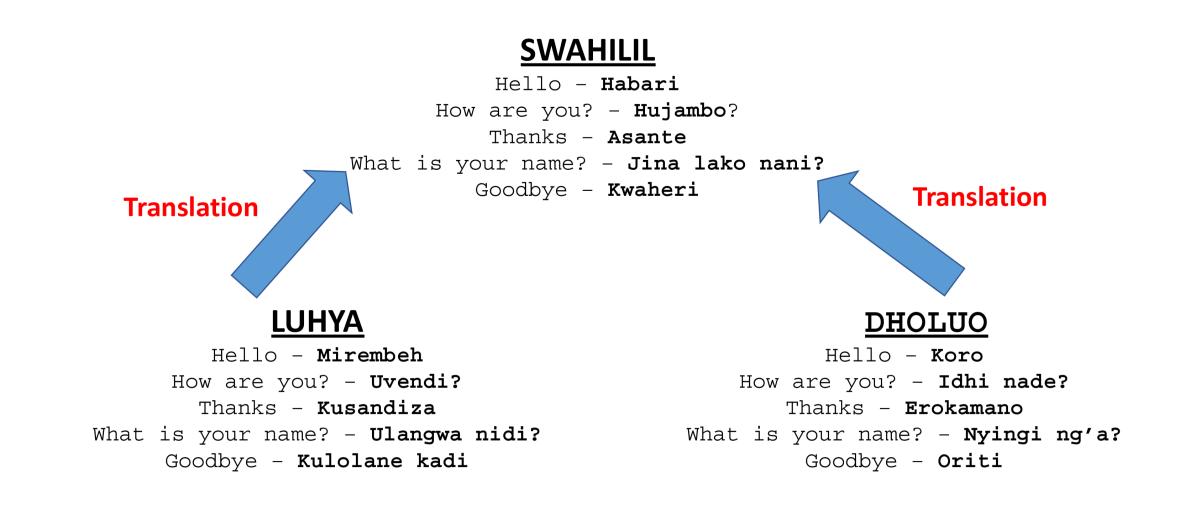
# Kenyan Languages Corpus for Natural Language Processing and Machine Learning

- 1. Project for the creation, labeling, augmentation and maintenance of datasets for machine learning.
- 2. Main Objective: To build Datasets that enable specific tasks in Natural Language Processing (NLP) and broader research in Machine Learning with the ultimate goal of supporting social impact.
- 3. Objectives include to develop:
  - ا labeled and unlabeled text corpora to support NLP and ML
  - parallel corpora for machine translation
  - corpora to support fundamental NLP tasks: i.e. Part of Speech tagging
  - corpora to support downstream NLP tasks: i.e. question answering
     speech corpora to support Speech synthesis and Text to Speech conversion

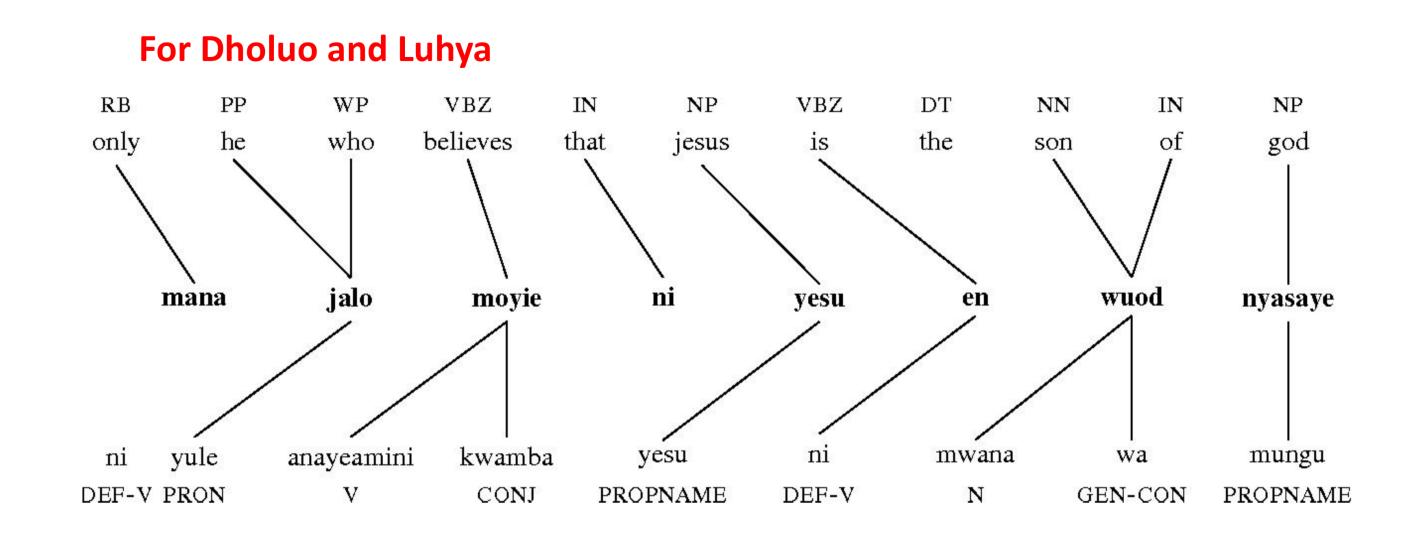
## Data Collection



#### Translation



# Part of Speech Tagging

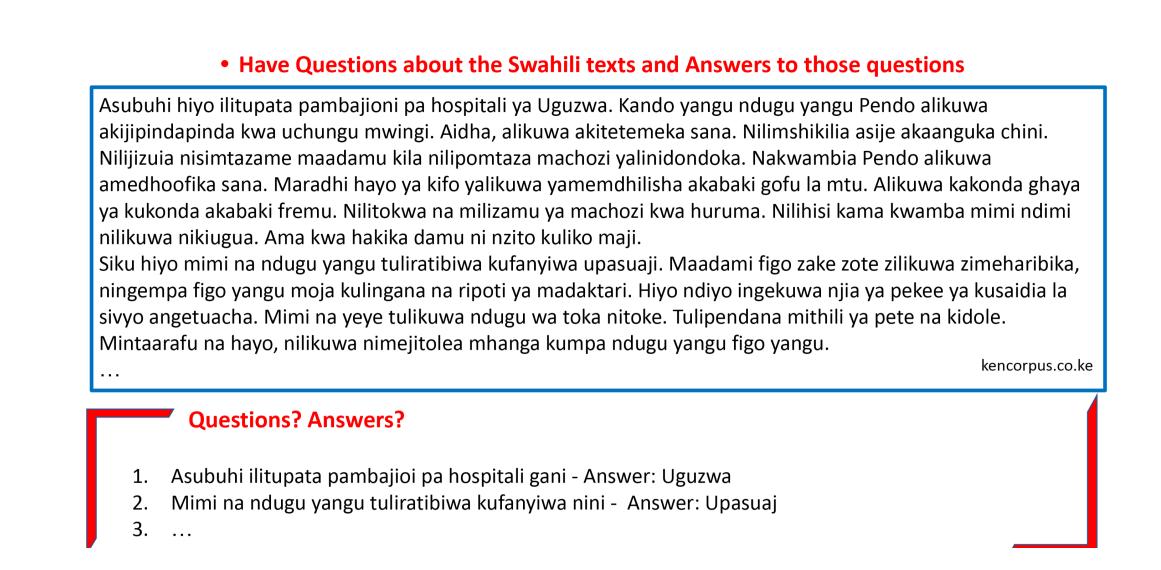


#### Link to Datasets

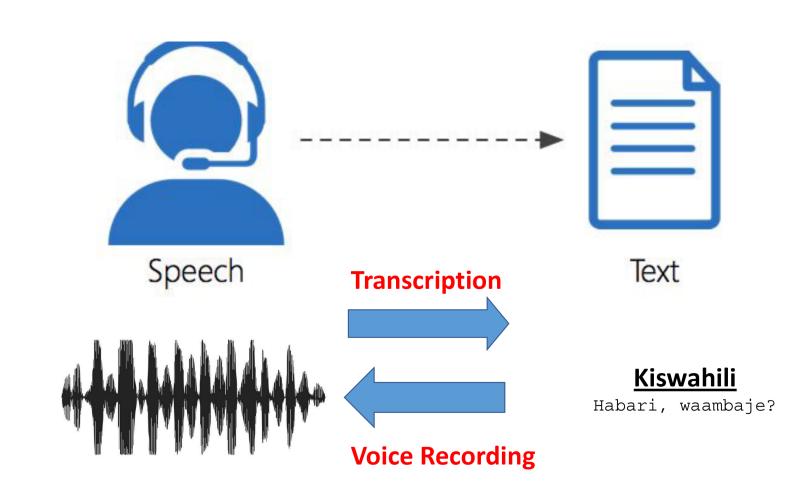
https://kencorpus.co.ke/

https://dataverse.harvard.edu/dataverse/Kencorpus

## Question Answering



### Transcription



# KenCorpus Datasets

Lang.	Texts	Audio Hrs	Data	No.	No.
Swahili	2585	19	MT	12,400(sents)	_
Dholuo	546	99	POS	143,000(words)	_
Luhya	977	57	QA	1445(texts)	7526(QA)
Table 1. Text and speech			Trans.	27(Hrs)	30,000(words)
data collected.				Table 2. Datasets.	

# Acknolegement

- LACUNA Funds
- Language research assistants and resource persons
- Maseno University
- University of Nairobi
- Africa Nazarene University

kencorpus@maseno.ac,ke Poster Number: #