# African Language Information Retrieval

*Hussein Suleman*

*University of Cape Town*
*Department of Computer Science*
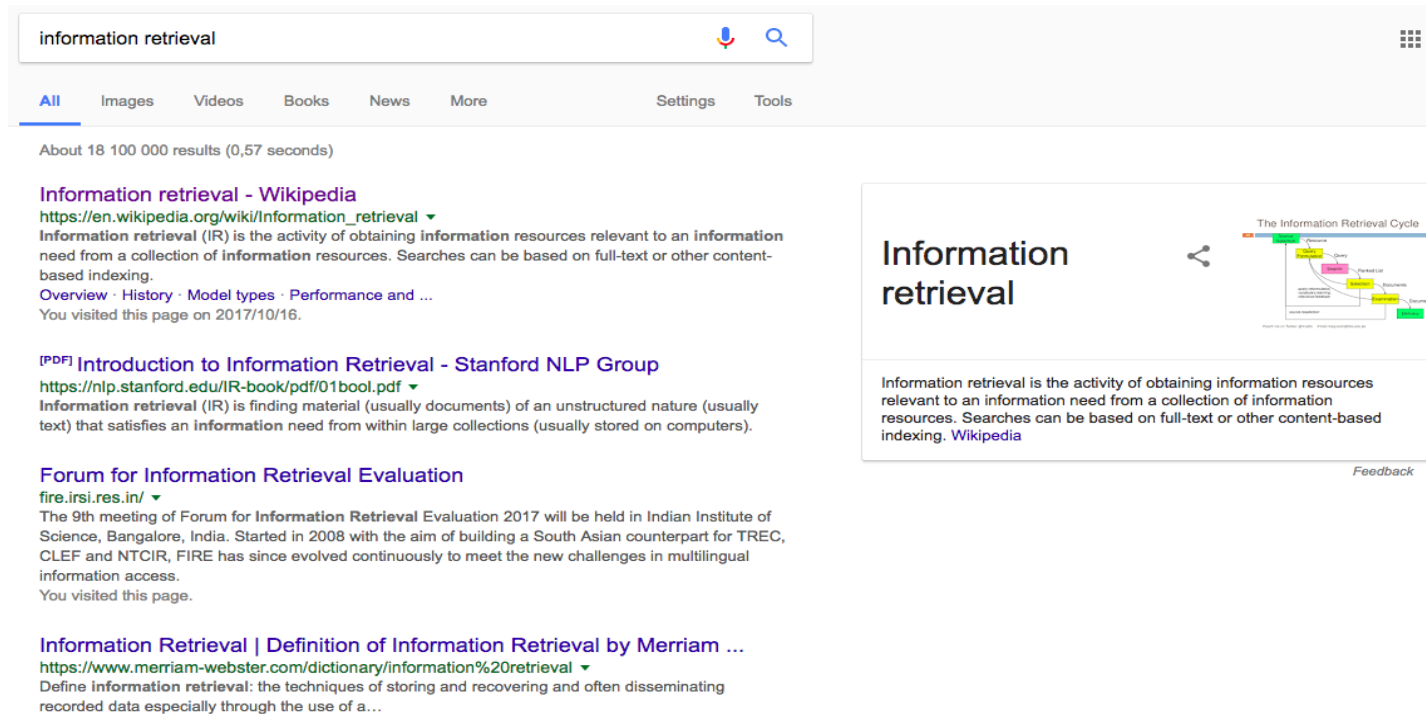*Digital Libraries Laboratory @ Centre for ICT4D*

*December 2017*

# Digital Libraries Lab

- Making information available to people
- Context-sensitive: low resource environment, different languages, different skills/culture, etc.
- 4 areas of interest:
  - Information Retrieval - HS
  - Educational Technology - HS
  - Knowledge Engineering - CMK
  - ICT4D Applications - various
- Some production systems / products

# Information Retrieval (IR)

□ Finding the most relevant information to satisfy an information need.

# IR in Africa

- Little understanding of local languages.
- Little relevant application.
- Poor general understanding of concept of search.
- Resource limitations.

# Arabic IR

- Feasibility study in 2006 into IR in isiXhosa, isiZulu, etc.
  - Very few text resources online.
  - Very little research in related areas.
  - Virtually no IR research.
- Arabic as African language!
  - Spoken in most of North Africa.
  - Dialect/usage not the same as Middle East.
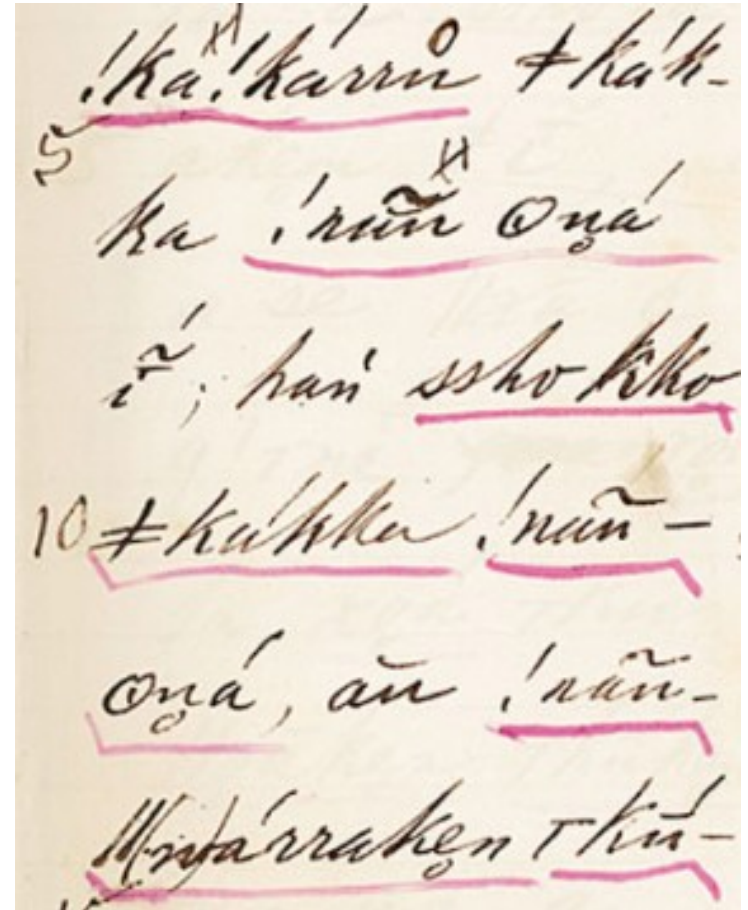  - Fashionable at the time.

# Mixed Language IR

*Mohammed Mustafa Ali, PhD*

- ❑ Noted that Google is language unaware.
- ❑ Poor results for mixed queries – queries in multiple languages.
  - ■ Dominant languages are dominant in results.
  - ■ Mixed language use is very popular in Africa.
- ❑ Solution: Examine queries and rerank based on language-based collection weights.

# |Xam IR

- Extinct Khoisan language.
- Language used in documenting early South African history/culture (25000 pages of stories).
- No Unicode representation.

# Digital Bleek and Lloyd Collection
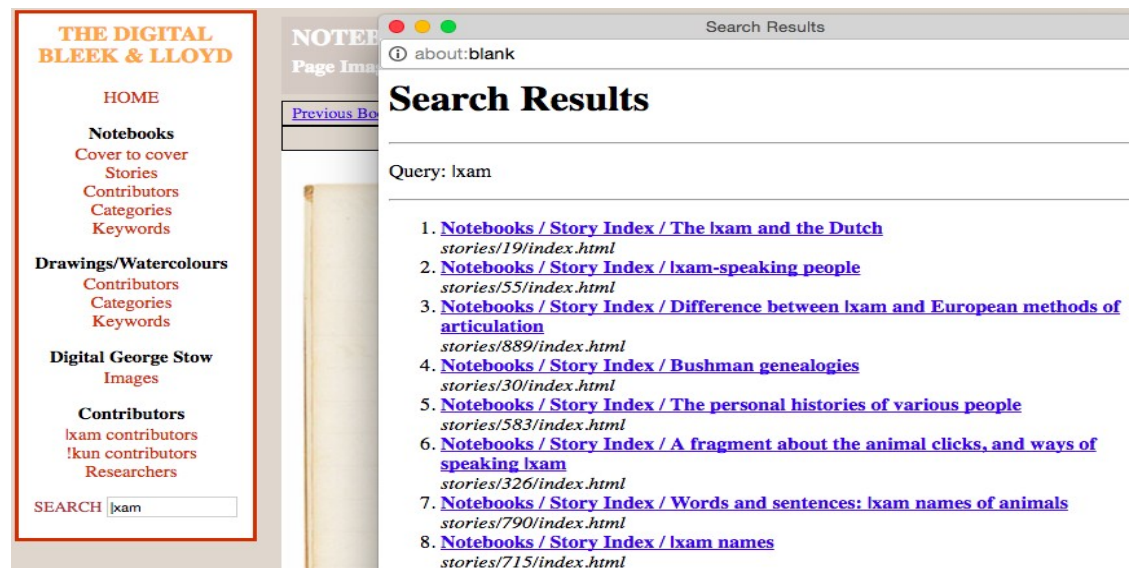
# Bleek and Lloyd: Low Resource IR

- IR engine within the browser – no network needed.
- Only simple transcriptions supported.

# Bleek and Lloyd: Dictionary

*Lebogang Molwantoa, Sanvir Manilal, Kyle Williams, BSc(Hons)*

- □  Visual dictionary – pictures of words.
- □  Find meanings of words in stories by image search.



THE BLEEK AND LLOYD |XAM DICTIONARY

This digital publication is part of a Llarec project to digitise, research and publish the Bleek and Lloyd Archive. Llarec (the Lucy Lloyd Archive, Resource and Exhibition Centre) is a University of Cape Town research centre located at the Michaelis School of Fine Art.

The project has been made possible by funding provided by the Andrew W. Mellon Foundation and De Beers; and is the result of the cooperation of the three curating institutions: University of Cape Town, Iziko South African Museum and The National Library of South Africa.

THE BLEEK LLOYD |XAM DICTIONARY

English–|xam Dictionary

THE BLEEK LLOYD |XAM DICTIONARY

|xam–English Dictionary

# Bleek and Lloyd: Transcription

*Kyle Williams, MSc; Ngoni Munyaradzi, MSc*

- Using machine learning to transcribe |Xam.
- Training data manually generated.
- 45% accuracy at best.

- Crowdsourcing had 10% better performance.
  - Answer determined by agreement among 3 amateur transcribers.

# Bleek and Lloyd: Text Input

*Sunkanmi Olaleye, MSc (current)*

- Inputting |Xam is non-trivial.
- Diacritics above, below and both; single and multiple characters.
- Custom Android keyboards for predictive and directed text entry in |Xam.

# Bantu Language IR

- ❑ Search engines in Bantu languages, especially South African languages (isiZulu, isiXhosa, etc.).

- ❑ Many core IR algorithms are unchanged but some language-specific algorithms needed:
    - ▪ Language identification
    - ▪ Text pre-processing and normalization
    - ▪ Ranking and reranking

# Bantu Language IR: AfriWeb

*Nkosana Malumba, Katlego Moukangwe, BSc(Hons)*

- IsiZulu Search Engine.
- High accuracy in identifying isiZulu vs. English+Italian.
- Simple morphological parser outperformed simple stemmer in IR results.



System Overview

# Bantu Language IR: Speech UI

*Michael Kyeyune, U/G; Morebodi Modise, MSc*

- Speech-driven mobile search interface in isiXhosa.



(a) Query submission interface

(b) Detecting voice queries

(c) Detected voice query with list of voice results

Fig. 3: Mobile voice interface

# Bantu Language IR: Similar Language IR

*Catherine Chavula, PhD (current);*
*Sinead Urisohn, Andre Lopes, BSc(Hons)*

- Exploit language similarity for those who can read multiple languages.
  - Reranking to emphasize language similarity in addition to relevance.
  - Universal language group text pre-processing, such as stemming.

# Bantu Language IR: Latest Projects

*Nyasha Katemauswa, U/G*

- ChiShona Search Engine.
  - Can we adapt the isiZulu framework to get better results in chiShona?

- KiSwahili IR.
  - Coming soon ...

# Corpora

- Corpora for African Language IR are rare.
  - There are limited corpora for speech recognition, speech synthesis, MT, etc.

- Very few documents online.

- Wikipedia has about 1000 (poor quality) pages in a large Bantu language collection!

- Lots of OOV, loan words, mixed texts, etc.

# Corpora: Crowdsourcing

*Sean Packham, MSc*
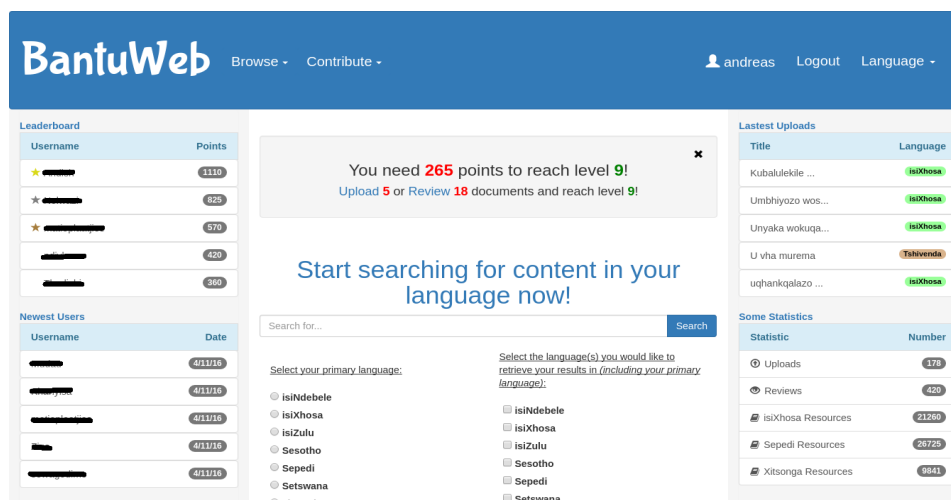
- Parallel corpus in isiXhosa-English.
- Will people contribute if money paid is varied or there is no money but only gamification?
  - Payment is only criterion!

# Corpora: SALANG

*Andreas von Holy, Osher Shuman, Alon Bresler, Bsc(Hons)*

◻ Create a central portal for documents in any SA Bantu language, with gamification, multilingual search, etc.

# Corpora: Long-term effects

*Jackson Moji, MSc (current)*

- Does gamification for corpus creation work in the long term?
  - Will people lose interest?
  - Will they continue to contribute?
  - How is intrinsic motivation affected by time?

- Extension of SALang project.

# IR for Development

*Gina Paihama, PhD (current)*

- How can we give users directed results to address unemployment?

*Selvas Mwanza, PhD (current)*

- Can we use Twitter data to evaluate developmental measures in society (e.g., level of free speech)?

# Other IR: Time-sensitive search

*Jivashi Nagar, PhD (current)*

- Can we exploit periodic patterns in search behaviour to provide contextually-relevant results?
  - Java in working hours leads to programming
  - Java over weekends leads to coffee

# Other IR: ETD Metadata Search

# questions, comments, ...

http://dl.cs.uct.ac.za/

enkosi
hamba kakuhle
thank you and go well