

Research + Community Organiser



Data Science for Social Impact











UNIVERSITEIT VAN PRETORIA UNIVERSITY OF PRETORIA YUNIBESITHI YA PRETORIA















Vukosi Marivate Principal Investigator and Leader



Herkulaas Combrink Senior Member, Director/Lecturer at



Abiodun Modupe Senior Member



Seani Rananga Senior Member, PhD CS, Lecturer



Idris Abdulmumin Postdoctoral Fellow



Kayode Olaleye Postdoctoral Fellow



Dawit Bogale PhD CS



Temitope Kekere PhD IT



Masikisiki PhD

Miehleketo Mathebula

CS MSc

David Walker

CS MSc



Mpho Mokoatle PhD CS

Rozina Myoya

CS MIT Big Data

Science

Masana Glad Balovi

CS Honours



Michelle Terblanche PhD CS

Mohlatlego Nakeng

CS MIT Big Data

Isheanesu Joseph

Dzingirai

CS Honours



Nombuyiselo Zondi PhD education



Abiola Akinbowale CS MIT Big Data Science



Tsholofelo Gomba CS MSc



Kwanele Radebe CS MSc



Matimba Shingange CS MIT Big Data Science



Kathleen Siminyu CS MSc









Mokoena Lehlohonolo CS Honours





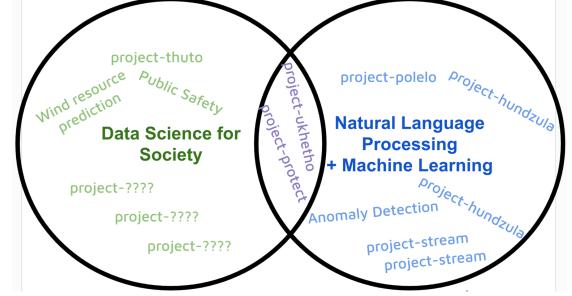
CS Honours



Lesiba Setsiba



+ Alumni





Pitso Khoboko

CS MSc

Sindane CS MSc

Assistant



Christiaan Lombard Keabetswe Madumo CS Honours, Research CS Honours, Research



Assistant

Manyako Manaswe CS Honours



Mulweli Patience Mukwevho CS Honours





CS Honours

CS Honours



Mosa Sethosa CS Honours



CS Honours



Data Science for Social Impact

Let's Get to the End



Al for Mother Tongue Education: Unlocking Potential

Personalized Learning in Local Languages

- Al-driven tutoring systems adapt content to individual student needs.
- Speech-to-text and text-to-speech tools enable accessibility.

Content Creation & Translation

- Machine translation enhances availability of learning materials.
- Large-scale corpus development supports textbook creation in African languages.

Teacher Support & Engagement

- Al-powered tools simplify content generation and assessment.
- Automated grading & feedback in mother tongue languages.



Al for Mother Tongue Education: Unlocking Potential

Preserving Linguistic Diversity

- Al helps document, digitize, and revitalize endangered languages.
- NLP models can foster literacy & language learning from early education stages.

Bridging the Digital Divide with Al-driven Language Inclusion!



Basic Education Opportunities for Al

Data Collection & Corpus Development

- Open Early-grade reading materials can help train language models for foundational understanding.
- Standardized tests and assignments in local languages can provide structured parallel text data.
- Schools can contribute to subject-specific terminology (math, science, etc.) in African languages.
- Collaboration with educators can help refine NLP-based spelling and grammar correction.

Government & Institutional Support

- aking public education materials openly available can assist research in general and NLP research.
- Partnerships between schools, universities, and tech firms can drive AI innovation.
- Empowering educators to integrate AI-driven tools in language teaching.

By integrating NLP into basic education, we can not only improve AI models but also strengthen linguistic diversity and digital literacy in African languages.



Why African Low Resource Languages Now?



Low Resource Natural Language Processing - Languages

Many Speakers, Not Many Resources

Data

- Speech
- Text
- Etc.

Tools for language (digital dictionaries, grammar tools etc.)



Defining the challenges to Low Resource

Languages

 Low availability of resources (Data, Tools, etc.)

Discoverability Q

Reproducibility

Focus

Benchmarks



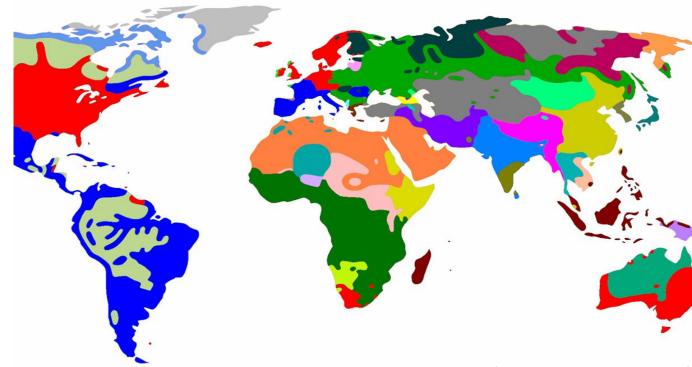
Scale and Complexity

A Focus on Neural Machine Translation for African Languages

Laura Martinus

Jade Abbott

Explore / Johannesburg, South Africa Retro Rabbit / Johannesburg, South Africa laura@explore-ai.net ja@retrorabbit.co.za



Repartition map of the languages over the world (version blank of key) [Wikimedia:User:Industrius]



Defining the challenges to Low Resource NLP

- O. The Left-Behinds
- 1. The Scraping-Bys
- 2. The Hopefuls
- 3. The Rising Stars
- 4. The Underdogs
- 5. The Winners

Class	5 Example Languages	#Langs	#Speakers	% of Total Langs
0	Dahalo, Warlpiri, Popoloca, Wallisian, Bora	2191	1.2B	88.38%
1	Cherokee, Fijian, Greenlandic, Bhojpuri, Navajo	222	30M	5.49%
2	Zulu, Konkani, Lao, Maltese, Irish	19	5.7M	0.36%
3	Indonesian, Ukranian, Cebuano, Afrikaans, Hebrew	28	1.8B	4.42%
4	Russian, Hungarian, Vietnamese, Dutch, Korean	18	2.2B	1.07%
5	English, Spanish, German, Japanese, French	7	2.5B	0.28%

Table 1: Number of languages, number of speakers, and percentage of total languages for each language class.

The State and Fate of Linguistic Diversity and Inclusion in the NLP World

Pratik Joshi* Sebastin Santy* Amar Budhiraja* Kalika Bali Monojit Choudhury

Microsoft Research, India {t-prjos, t-sesan, amar.budhiraja, kalikab, monojitc}@microsoft.com



Al current futures: Insatiable Appetite for Data and Compute



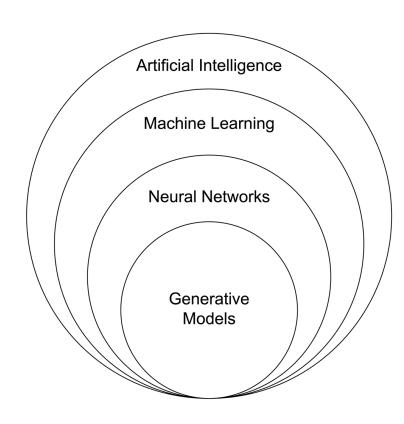
What is old is new again

This is not the first time we are going through an Al/ML hype cycle.

We are now in the age of Generative Al

It won't be the last.

How do we learn from the past to also get a grip on the future?





There will be change

Many ideas that will be floating around.

Most will wither away, some will float to the top.

Predicting the Future of AI with AI: High-Quality link prediction in an exponentially growing knowledge network

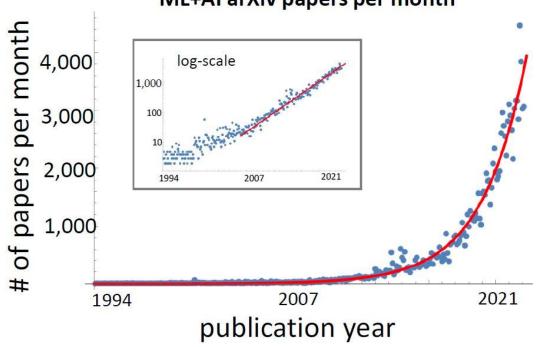
Mario Krenn,^{1,*} Lorenzo Buffoni,² Bruno Coutinho,² Sagi Eppel,³ Jacob Gates Foster,⁴

Independent Researcher, Leiria, Portugal.
 University of Pennsylvania, USA.
 University of California, San Diego, USA.

Andrew Gritsevskiy, 3,5,6 Harlin Lee, 4 Yichao Lu, 7 João P. Moutinho, 2 Nima Sanjabi, 8 Rishi Sonthalia, 4 Ngoc Mai Tran, 9 Francisco Valente, 10 Yangxinyu Xie, 11 Rose Yu, 12 and Michael Kopp 6

1 Max Planck Institute for the Science of Light (MPL), Erlangen, Germany.
2 Instituto de Telecomunicações, Lisbon, Portugal.
3 University of Toronto, Canada.
4 University of California Los Angeles, USA.
5 Cavendish Laboratories, Cavendish, Vermont, USA.
6 Institute of Advanced Research in Artificial Intelligence (IARAI), Vienna, Austria.
7 Layer 6 AI, Toronto, Canada.
8 Independent Researcher, Barcelona, Spain.
9 University of Texas at Austin, USA.

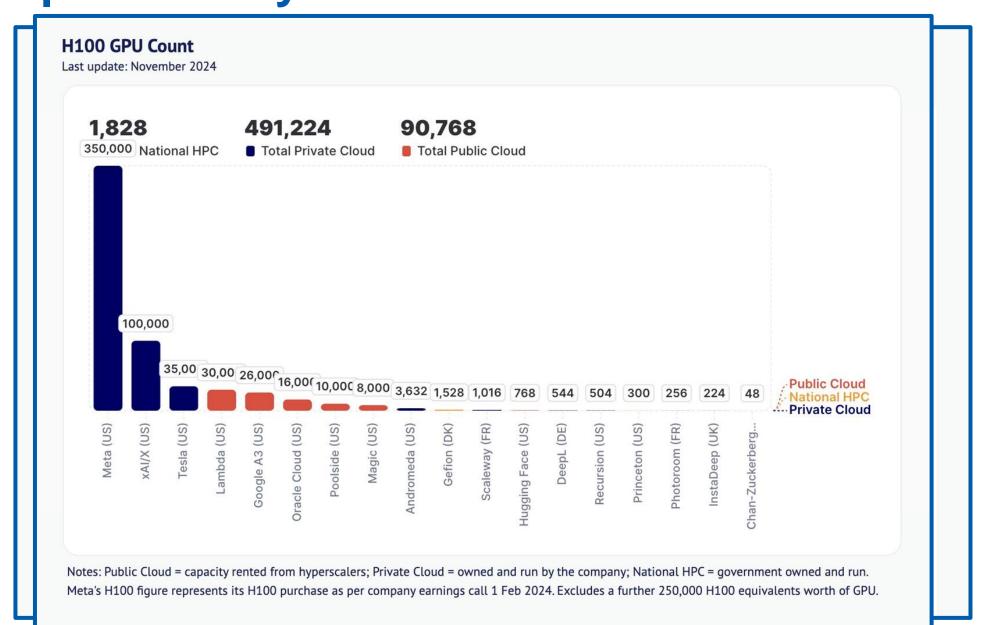
ML+AI arXiv papers per month





Compute is all you need?

https://www.stateof.ai/compute



Effectiveness of Data, At What Cost?



One of the largest gaps in the availability of data is in language

'Brute force' alone is unlikely to achieve parity in data generation in African languages compared to Anglophone, given the extent of the data imbalance.

Language (Nov '23) % Share of Internet content in local language

Global	English	52.60%
Benchmarks	Hindi	0.1%
	Afrikaans	0.003%
	Twi	0.00195%
Top African	Swahili	0.00135%
Languages	Malagasy	0.00022%
	Bambara	0.00025%
	Venda	0.000115%
	Hausa	0.00011%
Average with other African languages*		0.000999%
Sum of African Languages		0.01999%

English 53%

African 0.02%

2,650 x more content in English than African Languages

Lacuna Fund and
Masakhane are key
interventions, but only
scratch the surface



^{*} Afrikaans, Twi, Swahili, Bambara, Malagasy, Hausa, Venda, Haitian, Haitian Creole, Igbo, Luba-Katanga, Ndonga, Rundi, Tokelau, Tswana, Akan, Chichewa, Chewa, Nyanja, Fulah, Ganda, Masai

Data, Data, Data and More Data

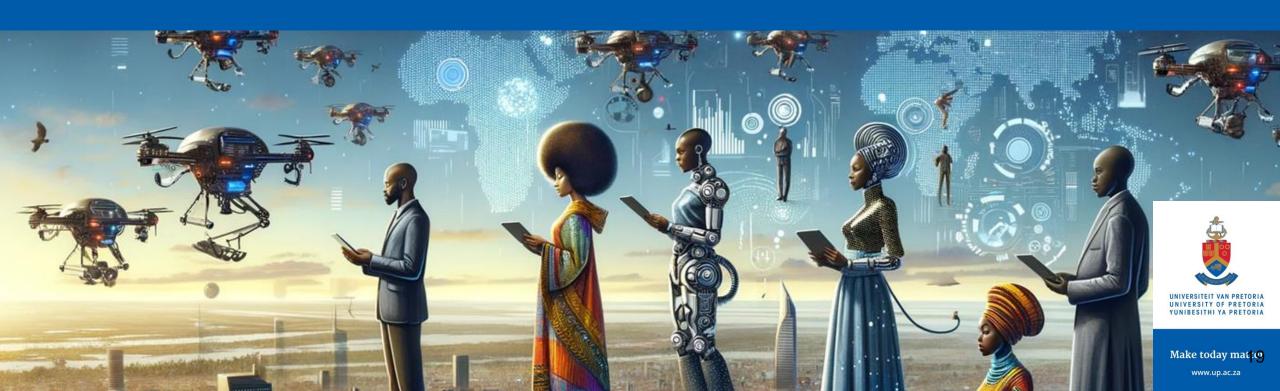
Where is the data coming from to train all these models?

Does it represent us? Bias!!!

How do make sure we protect citizens private information?

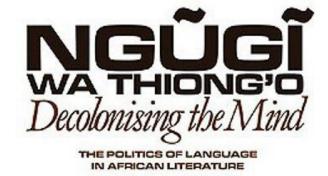
Africa cannot just be a [CHEAP] data market and AI consumer!!!!

Connecting back to our languages



Language is more than symbols

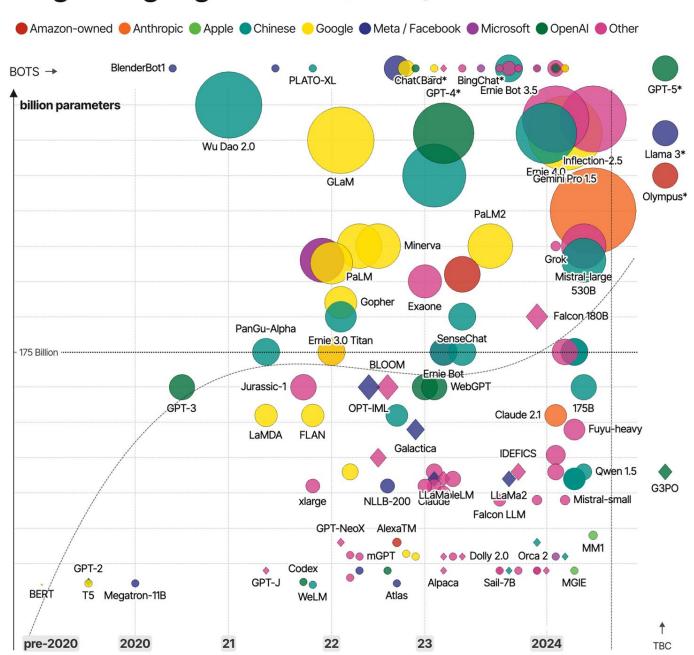
Language as communication and as culture are then products of each other. *Communication creates* culture: culture is a means of communication. Language carries culture, and culture carries, particularly through orature and literature, the entire body of values by which we come to perceive ourselves and our place in the world. How people perceive themselves and affects how they look at their culture, at their places politics and at the social production of wealth, at their entire relationship to nature and to other beings. Language is thus inseparable from ourselves as a community of human beings with a specific form and character, a specific history, a specific relationship to the world — Decolonising the Mind (16)







The Rise and Rise of A.I. Size = no. of parameters open-at Large Language Models (LLMs) & their associated bots like ChatGPT



Source https://informationisbeautiful.net/visualizations/therise-of-generative-ai-large-language-models-llms-likechatgpt/

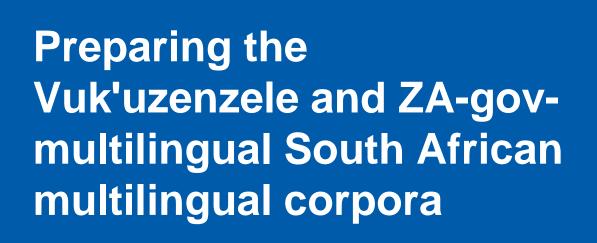
All Wikipedias ordered by number of articles

The languages listed here are Wikipedias that have been created as also be sorted by other columns. It excludes closed Wikipedias. History, after 15 October 2022, in the history of the current source of the

These statistics are updated four times a day. See commons:Dat doesn't update the numbers to what is seen at Commons, try pur

Nº ÷	Language +	Language (local) +	Wiki +	Articles +
1	English	English	en	6,815,681
2	Cebuano	Cebuano	ceb	6,119,222
3	German	Deutsch	de	2,903,386
4	French	français	fr	2,607,022
5	Swedish	svenska	sv	2,583,112
6	Dutch	Nederlands	nl	2,156,687
7	Russian	русский	ru	1,975,846
8	Spanish	español	es	1,947,807
9	Italian	italiano	it	1,860,435
10	Egyptian Arabic	مصرى	arz	1,623,024
82	Swahili	Kiswahili	SW	80,02

110 Yoruba	Yorùbá	yo	33,971
162 Zulu	isiZulu	zu	11,311
Southern	Sanatha	at	040



Richard Lastrucci¹, Isheanesu Dzingirai¹, Jenalea Rajab², Andani Madodonga¹, Matimba Shingange¹, Daniel Njini¹, Vukosi Marivate^{1,3}

¹Department of Computer Science, University of Pretoria ²School of Computer Science and Applied Mathematics, University of the Witwatersrand ³Lelapa AI

richard.lastrucci@tuks.co.za,ishe.dzingirai@gmail.com,jenalea.rajab@gmail.com,andanim412@gmail.com,mrosslyns@gmail.com,vukosi.marivate@cs.up.ac.za



Make today matter

www.up.ac.za

Setting the scene

Low Resource Language data is available, but may *not be easily accessible* for use in Natural Language Processing

Providing as much *metadata* as possible that identifies the source of the data, who might have written it and other translation information.

Some languages like *isiNdebele* has few resources and can be enhanced.

We build upon prior work such as that on *Autshumato* (Groenewald and Fourie, 2009; Groenewald and du Plooy, 2010)

- Hendrik J Groenewald and Liza du Plooy. 2010. Processing parallel text corpora for three south african language pairs in the autshumato project.
- Hendrik Johannes Groenewald and Wildrich Fourie. 2009. Introducing the autshumato integrated translation environment.



Goals of the Project

Liberate and prepare textual multilingual data from Government Communication Information System [GCIS] of the South African government.

Extract and align sentence pairs for all South African languages within the data.

Automate as much of the process as possible.

Provide *benchmarks for Machine Translation* using this data.

Release the cleaned and machine readable data openly for other researchers to use.



Creating the datasets



Cabinet Statements

What is this source?

- Statement of Cabinet Meetings of the South African government.
- A few weeks after the English statement is published, translated statements are made available.

Need

- Extract statements and make them easily accessible for NLP.
- Automate the process with monitoring of website changes and grow the dataset.



Cabinet Statements

You can use the filters to show only results that match your interests

Title

Speech Date

Statement on the Cabinet Meeting of 29 March 2023

Statement on the Cabinet Meeting of 15 March 2023

Statement on the Cabinet Meeting of 1 March 2023

Statement on the Cabinet Meeting of 1 March 2023

Statement on the Cabinet Meeting of 15 February 2023

Statement on the Cabinet Meeting of 15 February 2023

Statement on the Cabinet Meeting of 30 November 2022

Statement on the Cabinet Meeting of 16 November 2022

Statement on the Cabinet Meeting of 16 November 2022



Make today matter

Extracting Cabinet Statements

We extract the data from the HTML and prepare a JSON payload.

The JSON payload for each speech records:

- Date,
- Datetime,
- Title (in English),
- Url (top url for speech),
- Language payload for each language (eng, afr, nbl, xho, zul, nso, sep, tsn, ssw, ven, tso).
 - Title (in language),
 - Text (in language),
 - Url (for the translation).





government communications

Department:
Government Communication and Information System
REPUBLIC OF SOUTH AFRICA



Prepared Dataset

We have 162 cabinet statements spanning 2 May 2013 to 1 December 2022.

The dataset will update automatically when new, translated, statements are available on the gov.za website.

The dataset, code and automated scrapers are available at at https://github.com/dsfsi/gov-za-multilingual and Zenodo





government communications

Department:
Government Communication and Information System
REPUBLIC OF SOUTH AFRICA



Vuk'uzenzele Gov Newspaper

What is this source?

- "Vuk'uzenzele is a free Government Newspaper, committed to making a difference in the lives of South Africans"
- Published in 2 editions per month. Available physically and digitally [PDF] in all 11 South African languages.

Need

- Extract text from PDFS. Clean and make them easily accessible for NLP.
- Make the extraction process repeatable with humans in the loop.

k'iizenze

Ku endla swimakiwa swa m na swo tisasekisa hi mpepo

Health • Rural Development • Employment • Safety & Security • Education Viik'iizenze

Producing hair and beauty products with impepho



usiness based in Egotvibeni vil age in Tsolo, in the Eastern Cape

med by harsh sunlight on he



Make today matter

www.up.ac.za

https://www.vukuzenzele.gov.za/

Extracting Vuk'uzenzele PDFs

To clean it, a team member goes through each extracted text file and formats it as follows:

- Line 1: Title of article (in language)
- Line 2: empty line
- Line 3: Author of article (if available. If not, defaults to Vukuzenzele Unnamed)
- Line 4: empty line
- Line 5-end: body of article



nelerisiwile hi: Vufambisi bya Mfumo bya Yuhlanganisi na Mahungu (GCIS)

English/Xitsonga

Dzivamisoko 2023 Nkandziviso w

Ku endla swimakiwa swa misisi na swo tisasekisa hi mpepo



Sihle Mand

oko a karhatiwile hi xivundaza na ku pfaleleka eka fuluele ya yene eBorobeni ra Kapa hi madyambu man'wana hi nakarhi wa ntungu wa COVID-19, Nomhah Detwana a ri na xinkadyana xa rivoni ra giulupu leri ri nga yisa eka ku tswariwa ka masungulo ya bindru lera conga. Dotwana, wa 31 wa malembe, musunguri wa Namhla Coliection, a hlamusela riendora yene aka Yukuzuction, a hlamusela riendora yene aka Yukuzu-

> ula Collection i nongoloko wa bindzu wa swo rshlonge swa ntumbuluko leri ri kumekaka eka aya ra Egotybené Giolo, eKapa-Wuhumadyambu. to ra bindzu leri ri sungutile loko Dotwana a nakambe khandilhele ra khale leri ri nga rushehi sirasmono leri a ri goviwile hi dyambu lera aya mona fasitereni ra yena ra rasimbhi.

nhluvukiso wa vaaki ya ku va a ri mutivi wa swanangalasamahungunyingi yi hela hi ku hela ka 2020, utwana u pakile itbege ta yena ku ya sungula vhengele eren egirachini ya manana wa yena elgotyibeni. Vediz endla khomietiki na switirhisiwa swa sihungulo wa miri ku nga ri na mirhi ekaya hi ku

Swi ya emahlweni eka pheji 2



government communications

Government Communication and Information System REPUBLIC OF SOUTH AFRICA



Prepared Dataset

We have 53 editions of the newspaper spanning January 2020 to July 2022. (More are being added constantly)

Automations have been built to download and archive the PDFs, however manual effort is still required to extract and identify translated articles.

The dataset, code and automated scrapers are available at at https://github.com/dsfsi/vukuzenzele-nlp and Zenodo



Ku endla swimakiwa swa misisi na swo tisasekisa hi mpepo



Sihle Manda

oko a karhatiwile hi xivundza na k pfaleleka eka fulete ya yena eCbrobo mi ra Kapa hi madyambu man'wana ika mi ra Kapa hi madyambu man'wana ika miha Dotwana a ri na xinkadyana xa rivor ra glulupu leri ri nga yisa eka ku tswariwa k masungulo ya binduzi lero xonga. Dotwana, wa 31 wa malembe, musunguri wa Namhia Colle ction, a hilamusela riendora yene aka Vuk'uze

ole.

amhla Collection i nongoloko wa bindzu wa swo
yisa nklonge swa ntambuluko leri ri kumekala eka
okudaya ng Egoybe efi Colo, ekpa-ye-kumadyambu.
radzo ra hindzu leri ri sungurile loko Dotwana a
sikas nakamek kandhler si kalale teri ri nga ruhewa hi sinamono leri a ri govirule bi dyambu lera
mitanya mona Sisteneira ayan ra raminbil.
koko kordinka ya yena na Dorobe ra Kapa eka siynnge
hinburukoko wa va sidy sa ku za ri muthivi wa swa-

nhangalasamahungunyingi yi hela hi ku hela ka 2020, otwana u pakile libege ta yena ku ya sungula vhengele yena egirachini ya manana wa yena e Egotyibeni. Ndzi endla khosimetiki na switirhisiwa swa utshungulo wa miri ku nga ri na mirhi ekaya hi ku

Curi va amahhuani aka nhaji 3



Government Communication and Information System
REPUBLIC OF SOUTH AFRICA



Sentence Alignment and MT Benchmarks



Sentence Alignment: Output

The output of the sentence alignment task is 55 distinct parallel corpora for both the Vukuzenzele & Cabinet Statements.

Table 1: Language List with ISO 639-2 codes

Name	Code
isiZulu	zul
isiXhosa	xho
Afrikaans	afr
English	eng
Sepedi	nso
Setswana	tsn
Xitsonga	tso
Sesotho	sot
siSwati	ssw
Tshivenda	ven
isiNdebele	nbl

Table 2: Top ten datasets with the most observations with a cosine score greater than or equal to 0.65 in Vuk'uzenzele.

Language	No. of observations in		
pair	Vuk'uzenzele		
ssw-xho	2,202		
ssw-zul	2,183		
xho-zul	2,102		
nso-xho	2,081		
nso-tso	2,071		
ssw-tso	2,034		
nso-ssw	2,021		
tsn-tso	2,020		
tsn-xho	2,009		
tso-xho	2,009		

Table 3: Top ten datasets with the most observations with a cosine score greater than or equal to 0.65

Language	No. of observations
pair	in ZAgov Multilin-
	gual
nbl-ven	18,984
nso-ssw	18,697
zul-ssw	18,563
xho-ssw	18,387
xho-zul	18,145
xho-nso	18,110
xho-tso	17,954
ssw-tso	17,880
zul-tso	17,789
zul-nso	17,630





Data Characteristics for Aligned Translations

Table 5: Characteristics of the translation data for the Vuk'uzenzele (Vuk.) datasets

Size	
#sents (#src / #trg tokens)	
136 (3.4k / 3.9k)	
1715 (53.7k / 41.6k)	
1588 (29.9k / 37.6k)	
260 (9.9k / 7.5k)	
1366 (49.8k / 31.7k)	
1998 (58.9k/ 46.6k)	
230 (9.1k / 7k)	
1338 (25.8k / 31.5k)	
1874 (34.1k / 43k)	

Table 6: Characteristics of the translation data for the ZA-gov-multilingual (Gov.) datasets

Translation	Size		
Direction	#sents (#src / #trg tokens)		
nbl→eng	3513 (63.9k / 107k)		
nso→eng	14742 (460.9k / 375k)		
ssw→eng	15139 (291k/ 377.8k)		
sot→eng	4995 (145.9k / 153.5k)		
tsn→eng	14068 (493.1k / 362.2k)		
tso→eng	15393 (466.4k / 381.2k)		
ven→eng	3404 (68.2k / 96.6k)		
xho→eng	15853 (318.2k / 389.5k)		
zul→eng	15503 (327.5k / 384.1k)		



NMT Benchmarks

Training translation models on a small amount of data can hinder the quality

Therefore we leverage a Massively Multilingual Model (M2M100)(Fan et al , 2021) and fine-tune it on our relatively small datasets (Adelani etal, 2022)

To provide our results in context and for comparison we also fine-tune the M2M100 model on subsets of the existing **Autshumato** corpora (Which exist only in the 'eng-xxx')

We focus our efforts on providing NMT benchmarks for the low resource African languages

- David Adelani et al. 2022. A few thousand translations go a long way! leveraging pre-trained models for African news translation
- Angela Fan et al. 2021. Beyond english-centric multilingual machine translation.



MT Results

The highest BLEU scores are distributed across the ZA-gov-multilingual and Autshumato NMT models

ZA-gov-multilingual models achieving a higher score for **Setswana**, **Xitsonga and isiZulu**

Highest benchmark result for **Xitsonga** across all datasets, demonstrating the effectiveness of transfer learning for new low-resource language datasets.

Our contributions:

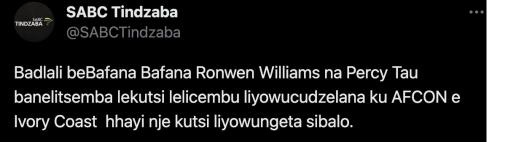
- Extend the benchmark translation resources (in the government data domain) to isiNdebele, isiXhosa, siSwati and Tshivenda
- Broaden the translation direction beyond English as the source language.

Table 4: BLEU scores for Massively Multilingual Transfer on xxx-eng translations using the Vuk'uzenzele (Vuk.), ZA-gov-multilingual (Gov.) and subsets of the available Autshumato datasets (Aut.)

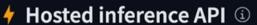
Translation	BLEU			
Direction	Vuk.	Gov.	Aut.	
nbl→eng	7.33	8.04	12.24	
nso→eng	9.29	26.50	-	
ssw→eng	4.80	28.72	-	
sot→eng	4.55	10.21	14.83	
tsn→eng	2.80	29.68	28.04	
tso→eng	13.86	35.40	32.10	
ven→eng	2.32	9.68	17.24	
xho→eng	6.05	26.81	-	
zul→eng	9.97	30.03	25.90	

It is noted that variations in the Autsumato subset selections will yield different results and an in-depth analysis is left for future work.

[•] Source languages which were not including in the original M2M100 pretraining are highlighted



#SABCNews
#SABCTindzaba



☐ Text2Text Generation

Badlali beBafana Bafana Ronwen Williams na Percy Tau banelitsemba lekutsi lelicembu liyowucudzelana ku AFCON e Ivory Coast hhayi nje kutsi liyowungeta sibalo.

Compute

∺+Enter

Computation time on Intel Xeon 3rd Gen Scalable cpu: 3.522 s

Bafana Bafana Ronwen Williams and Percy Tau are confident that the team is competing to AFCON on Ivory Coast and not only raise awards.

X

0.0





Breaking News



Mo, 11 sep. \bigcirc 30° Tu, 12 sep. -0- 24° We, 13 sep. -0- 23° Th, 14 sep. -0- 27°







News /

Dikgang /

Ba lebisitswe molato wa petelelo le go kgothosa

BA LEBISITSWE MOLATO WA PETELELO LE GO KGOTHOSA

10 Sep 2023

Borre ba le babedi ba ba tlholegang kwa lefatsheng la Zimbabwe ba itshupile fa pele ga lekgotla ka Labone mabapi le molato wa go tsaya ka dikgoka le go thubetsa mme yo o dingwaga tse di masome a matlhano le bobedi (52) wa kgotla ya Kanamo kwa Mahalapye.

Borre bao, wa dingwaga tse di masome a mararo le bobedi (32) le wa tse di masome a mararo le borataro (36) mo bekeng e e fetileng ba ne ba tsenelela mme yo mo ntlong ba mo amoga megala ya letheka e le mebedi ba bo ba mmetelela.

Mookamela Mapodisi a Mahalapye Central Supt. Boitshepho Mudongo o kaile fa ka lesego ba kgonne go ba tshwara mme megala yotlhe ka bobedi e bonwe.

A re bobedi jo bo amanngwa le dikgetsi tse dingwe tsa go kgothosa di le pedi ka ba fitlhetswe ka dilwana tse di neng di begilwe fa di

X

Hosted inference API ①

👯 Token Classification

Borre ba le babedi ba ba tlholegang kwa lefatsheng la Zimbabwe ba itshupile fa pele ga lekgotla ka Labone mabapi le molato wa go tsaya ka dikgoka le go thubetsa mme yo o dingwaga tse di masome a matlhano le bobedi (52) wa kgotla ya Kanamo kwa Mahalapye.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: cached

Borre ba le babedi ba ba tlholegang kwa lefatsheng la Zimbabwe Loc ba itshupile fa pele ga lekgotla ka Labone DATE mabapi le molato wa go tsaya ka dikgoka le go thubetsa mme yo o dingwaga tse di masome a matlhano DATE le bobedi DATE (52) wa kgotla ya Loc Kana PER mo kwa Loc Mahalapye Loc.

♦ Hosted inference API ③

👯 Token Classification

Borre ba le babedi ba ba tlholegang kwa lefatsheng la Zimbabwe ba itshupile fa pele ga lekgotla ka Labone mabapi le molato wa go tsaya ka dikgoka le go thubetsa mme yo o dingwaga tse di masome a matlhano le bobedi (52) wa kgotla ya Kanamo kwa Mahalapye.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.042 s

```
Borre Noun
                    babedi NUM
                                ba DET
                                                tlholegang VERB kwa DET lefatsheng PROPN la DET
                                                                                                Zimbabwe PROPN
           ba le DET
                                        ba PRON
                                                                                                                ba PRON
                                      lekgotla Noun
itshupile VERB
             fa PUNCT
                      pele ADV
                                                    ka ADP Labone mabapi Noun le AUX molato Noun
                               ga ADP
                                                                                                   wa DET go tsaya VERB
                                                                                                                         ka ADP
                                      mme ADV yo DET o dingwaga NUM tse DET di masome NUM
dikgoka Noun
             le cconj go thubetsa VERB
                                                                                                       matlhano NUM
                                                                                                                     le cconj
                                                                                              a PROPN
bobedi NUM
             PUNCT 52 NUM ) PUNCT WA DET kgotla NOUN ya DET Kana PROPN mo kwa DET
                                                                                    Mahalapye PROPN . PUNCT
```

X

Mavito: South African Terminology, Lexicon and Glossary Project

https://dsfsi.github.io/za-mavito/

Vukosi Marivate (PI)* Fiskani Banda Richard Lastrucci Mohlatlego Nakeng Kayode Olaleye Thapelo Sindane

Data Science for Social Impact

University of Pretoria





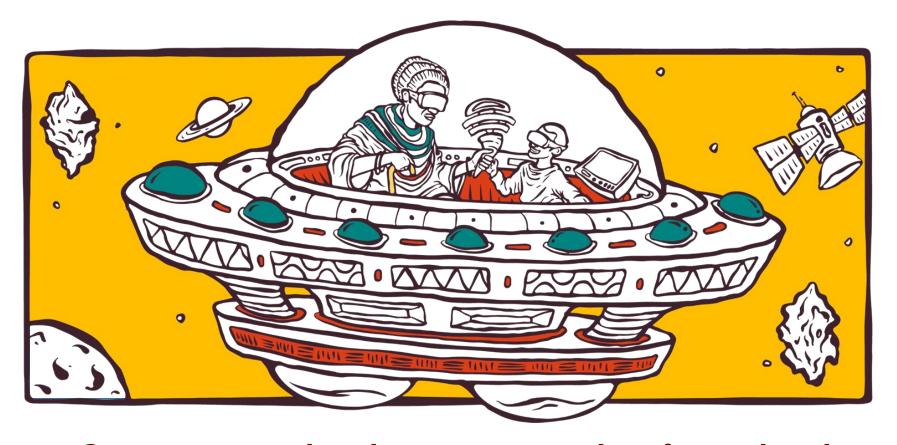
Make today matter

We need to get tools in people hands!

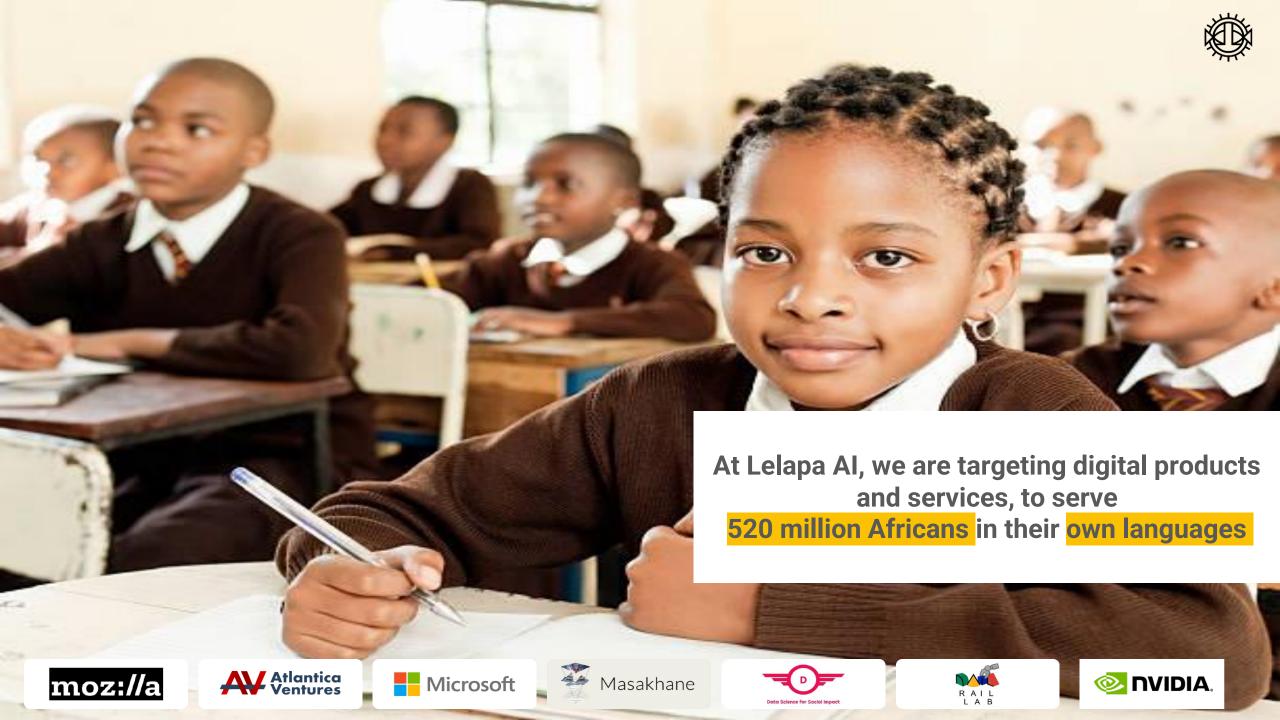




LELAPA AI



Grow your market, know your market, foster loyal connections



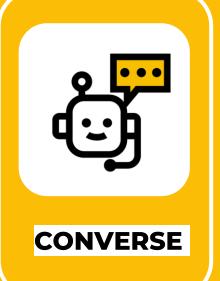
VULAVULA FEATURES



Voice to text conversion



Text Analysis functions



NLU capability to power chat



Text to text translation



Text to speech conversion

Solutions

EGRA

Vi

Support the enhancement of EGRA-Al for home languages

Teacher support

Virtual Teaching
Assistant:
Support in creating
multilingual lesson
plans and chat for live
teacher support.

Education language model

Train an education
(literacy & numeracy)
specific model with
robust data for use
across multiple literary
programmes.

Partner org AI enablem ent

Enable AI capabilities of partners & their current programmesto prioritise based on most feasible over a period



Collaborations









MEDIA MONITORING International Research Centre AFRICA of Artificial Inteligence under the auspices of UNESCO











Network for Artificial Intelligence, Knowledge and SUStainable development a nexus and central meeting point between AI and SDGs









AND INFORMATION TECHNOLOGY LAW



Make today matre: www.up.ac.za

DSFSI Support - Thank You

















FUTURE PROFESSORS PROGRAMME (PHASE 01)

MaSS





Meta

facebook









JPMORGAN CHASE & CO.



AFRICA

science & innovation

Department: Science and Innovation REPUBLIC OF SOUTH AFRICA



National Research **Foundation**



TensorFlow









Make today matre www.up.ac.za

Thank you

Questions



Prof. Vukosi Marivate

vukosi.marivate@cs.up.ac.za

https://dsfsi.github.io

- @vukosi
- @DSFSI_Research

Keep in touch

Join our research group newsletter

https://dsup.substack.com/

Made with ♥□ in Tshwane



Faculty of Engineering, Built Environment and Information Technology

Fakulteit Ingenieurswese, Bou-omgewing en Inligtingtegnologie / Lefapha la Boetšenere, Tikologo ya Kago le Theknolotši ya Tshedimošo

