NLP for African (Nguni) languages

C. Maria Keet

Department of Computer Science University of Cape Town, South Africa mkeet@cs.uct.ac.za

Guest lecture NLP course Poznan University of Technology, 10 April 2018

1/69

CS@UCT



CS@UCT



2/69

CS@UCT



Outline

Motivation

- Context
- Language 'crash course'

2 Corpus-based spellcheckers

- Error detection for isiZulu and isiXhosa
- Error correction for isiZulu and isiXhosa
- Discussion
- 3 Rule-based NLG
 - What is CNL, NLG?
 - Generating basic sentences
 - Language learning exercises
 - Summary

Outline

Motivation

- Context
- Language 'crash course'

Corpus-based spellcheckers

- Error detection for isiZulu and isiXhosa
- Error correction for isiZulu and isiXhosa
- Discussion

3 Rule-based NLG

- What is CNL, NLG?
- Generating basic sentences
- Language learning exercises

Summary

Motivation

- IsiZulu and isiXhosa most widely spoken languages in South Africa by first language speakers
- 23% or about 11 million people (isiZulu), 8 million (isiXhosa)
- Have very limited ICT support

Motivation

- IsiZulu and isiXhosa most widely spoken languages in South Africa by first language speakers
- 23% or about 11 million people (isiZulu), 8 million (isiXhosa)
- Have very limited ICT support
- They use computers for work, social media... (e.g.: avg. 1.5 mobile connection/pp in SA)
- So, NLP for these languages: searching online, spellcheckers, machine translation, speech, etc.

Basics

- Bantu languages: group of languages spoken in Sub-Saharan Africa
- Bantu means 'human'; bit of a laden term, but still used in linguistics
- Number of languages varies by who counts (> 200 at least)
- Organised in so-called Guthrie zones

Language 'crash course'

Guthrie Zones



7/69

- System of noun classes (besides the I, you (sg.), he/she, we, you (pl.), they)
 - For 3rd pers. sg./pl. nouns (i.e., not pers. pron.): they are classified into a noun class
 - Meinhof identified 23 noun classes (not all of them used)
 - There's some semantics to them: e.g., NC1 for humans, NC9 for animals, NC15 infinitive nouns

- System of noun classes (besides the I, you (sg.), he/she, we, you (pl.), they)
 - For 3rd pers. sg./pl. nouns (i.e., not pers. pron.): they are classified into a noun class
 - Meinhof identified 23 noun classes (not all of them used)
 - There's some semantics to them: e.g., NC1 for humans, NC9 for animals, NC15 infinitive nouns
- Most, but not all, of the languages are agglutinating
 - i.e., what are separate words in, say, English are 'components' of a word
 - Ex: titukakimureeterahoganu 'We have never ever brought it to him'

(Runyankore, Uganda)

- System of noun classes (besides the I, you (sg.), he/she, we, you (pl.), they)
 - For 3rd pers. sg./pl. nouns (i.e., not pers. pron.): they are classified into a noun class
 - Meinhof identified 23 noun classes (not all of them used)
 - There's some semantics to them: e.g., NC1 for humans, NC9 for animals, NC15 infinitive nouns
- Most, but not all, of the languages are agglutinating
 - i.e., what are separate words in, say, English are 'components' of a word
 - Ex: titukakimureeterahoganu (Runyankore, Uganda) 'We have never ever brought it to him' ti tu ka ki mu reet er a ho ga nu neg-(NC2 SC)-RM-(NC7 SC)-(NC1 SC)-VR-App-FV-Loc-Emp-Dec

- System of noun classes (besides the I, you (sg.), he/she, we, you (pl.), they)
 - For 3rd pers. sg./pl. nouns (i.e., not pers. pron.): they are classified into a noun class
 - Meinhof identified 23 noun classes (not all of them used)
 - There's some semantics to them: e.g., NC1 for humans, NC9 for animals, NC15 infinitive nouns
- Most, but not all, of the languages are agglutinating
 - i.e., what are separate words in, say, English are 'components' of a word
 - Ex: titukakimureeterahoganu (Runyankore, Uganda) 'We have never ever brought it to him' ti tu ka ki mu reet er a ho ga nu neg-(NC2 SC)-RM-(NC7 SC)-(NC1 SC)-VR-App-FV-Loc-Emp-Dec
- System of concordial agreement (more about that soon)

NC	AU	PRE	Stem (ex-	Meaning	Example	
			ample)			
1	u-	m(u)-	-fana	humans and other	umfana	boy
2	a-	ba-	-fana	animates	abafana	boys
1a	u-	-	-baba	kinship terms and proper	ubaba	father
2a	o-	-	-baba	names	obaba	fathers
3a	u-	-	-shizi	nonhuman	ushizi	cheese
(2a)	o-	-	-shizi		oshizi	cheeses
3	u-	m(u)-	-fula	trees, plants, non-paired	umfula	river
4	i-	mi-	-fula	body parts	imifula	rivers
5	i-	(li)-	-gama	fruits, paired body parts,	igama	name
6	a-	ma-	-gama	and natural phenomena	amagama	names
7	i-	si-	-hlalo	inanimates and manner/	isihlalo	chair
8	i-	zi-	-hlalo	style	izihlalo	chairs
9a	i-	-	-rabha	nonhuman	irabha	rubber
(6)	a-	ma-	-rabha		amarabha	rubbers
9	i(n)-	-	-ja	animals	inja	dog
10	i-	zi(n)-	-ja		izinja	dogs
11	u-	(lu)-	-thi	inanimates and long thin	uthi	stick
(10)	i-	zi(n)-	-thi	objects	izinthi	sticks
14	u-	bu-	-hle	abstract nouns	ubuhle	beauty
15	u-	ku-	-cula	infinitives	ukucula	to sing
17		ku-		locatives, remote/ general		locative

Concordial agreement—example (isiZulu, South Africa)

Abafana abancane bazozithenga izincwadi ezinkulu **aba**-fana **aba**-ncane **ba**- zo- **zi**- thenga **izi**-ncwadi e-**zi**-nkulu **2**.boy **2**.small **2.SUBJ**-FUT-**10.OBJ**-buy **10**.book REL-**10**.big 'The little boys will buy the big books'

Other illustrative examples (isiZulu)

- 'and', enumerative: na-, phonologically conditioned
 - Ex: milk and butter: ubisi nebhotela
 - Ex: butter and milk: *ibhotela <u>no</u>bisi*

(-a+i-=-e-)(-a+u-=-o-)

Other illustrative examples (isiZulu)

- 'and', enumerative: na-, phonologically conditioned
 Ex: milk and butter: ubisi <u>ne</u>bhotela (-a+i-=-e-)
 Ex: butter and milk: ibhotela <u>no</u>bisi (-a+u-=-o-)
- 'is not a': combine NEG SC with PRON, both depend on noun class
 Ex: an animal is not a plant: *isilwane <u>asiwona</u> umuthi* Ex: a plant is not an animal: *umuthi <u>awusona</u> isilwane*

11/69

Other illustrative examples (isiZulu)

- 'and', enumerative: na-, phonologically conditioned
 Ex: milk and butter: ubisi <u>ne</u>bhotela (-a+i-=-e-)
 Ex: butter and milk: ibhotela <u>no</u>bisi (-a+u==-o-)
- 'is not a': combine NEG SC with PRON, both depend on noun class
 Ex: an animal is not a plant: *isilwane <u>asiwona</u> umuthi* Ex: a plant is not an animal: *umuthi <u>awusona</u> isilwane*
- 'all'/'each' (\forall), 'at least one' (\exists): quantifiers depend on noun class
 - Ex: all animals *zonke izilwane* / all plants *yonke imithi*
 - Ex: at least one animal: *isilwane <u>esisodwa</u> /* at least one plant *umuthi* <u>owodwa</u>
- copulative (to be): depends on first letter of noun: *ng* for a-, o-, u-, else *y*-
 - Ex: is a dog: *yinja*
 - Ex: is a grandmother: <u>ngugogo</u>

Outline

Motivation

Context

• Language 'crash course'

2 Corpus-based spellcheckers

- Error detection for isiZulu and isiXhosa
- Error correction for isiZulu and isiXhosa
- Discussion
- 3 Rule-based NLG
 - What is CNL, NLG?
 - Generating basic sentences
 - Language learning exercises

Summary

General issues and aims

- Very limited available spellcheckers for use:
 - Outdated/software not working anymore (OpenOffice plugin)
 - Online one with too many clicks, popups, and ads¹

¹https://www.spellchecker.net/africa_zulu_spell_checker_html > > > > > 0 0 0 13/69

General issues and aims

- Very limited available spellcheckers for use:
 - Outdated/software not working anymore (OpenOffice plugin)
 - Online one with too many clicks, popups, and ads¹
- Investigate development of spellchecker for isiZulu
- Find an approach that can be used across (agglutinating) Bantu languages

¹https://www.spellchecker.net/africa_zulu_spell_checker_html > > > > 0 0 0 13/69

What will work best?

• Dictionary approach won't work due to (theoretically) agglutination and (practically) limited online dictionaries

What will work best?

- Dictionary approach won't work due to (theoretically) agglutination and (practically) limited online dictionaries
- Data-driven statistical model or grammar-based (morphological analyser-based) approach?
 - Try that for both isiZulu and isiXhosa

What will work best?

- Dictionary approach won't work due to (theoretically) agglutination and (practically) limited online dictionaries
- Data-driven statistical model or grammar-based (morphological analyser-based) approach?
 - Try that for both isiZulu and isiXhosa
- ⇒ Rules for the POS categories coded perform better overall (isiXhosa), but data-driven approach faster and more reusable across languages (isiZulu, isiXhosa), despite being underresourced

First iteration [Ndaba et al.(2016)]

- 1. Use corpus for data
- 2. Generate n-gram statistics; tried trigrams and quadrigrams; e.g.:
 - original word ngimbona
 - trigrams: ngi, gim, imb, ...
 - quadrigrams: ngim, gimb, ...
- 3. Compute frequencies of the tri/quadrigrams
- 4. Determine threshold of when a tri/quadrigram would probably be wrong
- 5. Then when a word is given:
 - a. Generate its trigrams
 - b. Check probability each trigram
 - c. If tri/quadrigram below threshold: flag word as incorrectly spelled

Example language model creation

corpus	generate trigrams	unique trigrams	frequency	probability
sivela	siv ive vel ela	ala	1	0.03703704
ngihamba	ngi gih iha ham amb mba	amb	1	0.03703704
uvelaphi	uve vel ela lap aph phi	aph	2	0.07407407
ngivala	ngi giv iva val ala	ela	3	0.11111111
uvelaphi	uve vel ela lap aph phi	gih	1	0.03703704
		giv	1	0.03703704
		ham	1	0.03703704
		iha	1	0.03703704
		iva	1	0.03703704
		ive	1	0.03703704
		lap	2	0.07407407
		mba	1	0.03703704
		ngi	2	0.07407407
		phi	2	0.07407407
		siv	1	0.03703704
		uve	2	0.07407407
		val	1	0.03703704
		vel	3	0.11111111
				1

16/69

Example how this then works in the spellchecker

000	
<u>File Edit H</u> elp	
Elle Edit Help Run Save Copy Paste Clear Help HiZulu orrectly. All trigrams are above the threshold, so it is assumed to have been spelled correctly. A teast one trigram is below the threshold (there is no gni trigram in the language model), so it is probably misspelled, and flagged as such. A word it was not trained on, but all trigrams are above the threshold, so it is assumed to have been spelled correctly. At least one trigram is below the threshold (no vea), so it is probably misspelled. It suggests a correction using more probable trigrams.	Ignore once Ignore all Add to dictionary Suggestions Ingivela Change
	Change all
Double click on errors to make correction one at a time	Exit

Basic approach for testing

- 10-fold cross-validation for training and testing data set
- 3 corpora to test effect of corpus on accuracy
 - Ukwabelana [Spiegler et al.(2010)]; 288 106 words, 87033 unique
 - Section of the isiZulu National Corpus [Khumalo(2015)]; 538 732 words, 33020 unique
 - IsiZulu news items (collected by MK); 21250 words, 9587 unique
- Different thresholds
- 46 know-to-be-incorrect words added
- Use accuracy as measure, with confusion matrix (TP, TN, FP, FN)

Major outcomes

- 89% accuracy (on par with older data [Bosch and Eisele(2005), Prinsloo and de Schryver(2004)])
- The spellchecker performed slightly better with trigrams than with quadrigrams
- Accurate in detecting words that do not occur in the training corpus
- The most updated corpora are preferable



Major outcomes

- 89% accuracy (on par with older data [Bosch and Eisele(2005), Prinsloo and de Schryver(2004)])
- The spellchecker performed slightly better with trigrams than with quadrigrams
- Accurate in detecting words that do not occur in the training corpus
- The most updated corpora are preferable



• Tests with more data: 83% accuracy with noisy data, 85% with cleaned data

Try this with isiXhosa

- Use code for isiZulu, but feed it isiXhosa texts to create a language model for isiXhosa
- Determine threshold
- Determine accuracy

Try this with isiXhosa: results

- 20K tokens corpus, mainly medical documents
- Threshold: 0.002 (marginally better than 0.003)
- Accuracy: about 79%
- (current implemented version trained with much more text)

Error correction for isiZulu

• Statistical language model based approach as well

Error correction for isiZulu

- Statistical language model based approach as well
- Insertions, deletions, transpositions, substitutions

Error correction for isiZulu

- Statistical language model based approach as well
- Insertions, deletions, transpositions, substitutions
- Levenshtein distance + probability of alternate trigram

Probabilities of successive trigrams
Error correction for isiZulu

- Statistical language model based approach as well
- Insertions, deletions, transpositions, substitutions
- Levenshtein distance + probability of alternate trigram
 - e.g., with typo: nii
 - intended: ngi
 - distance: 1 (a substitution of i where there should be g)
- Probabilities of successive trigrams

Error correction for isiZulu

- Statistical language model based approach as well
- Insertions, deletions, transpositions, substitutions
- Levenshtein distance + probability of alternate trigram
 - e.g., with typo: nii
 - intended: ngi
 - distance: 1 (a substitution of i where there should be g)
- Probabilities of successive trigrams
- Measures:
 - Can it propose something (C_s) ?
 - Is the intended word among the proposed words (C_v) ? (relevance)

Deletion typo and a suggested correction (running example)



୍ର ବ୍ ୧୦ 23 / 69

Error correction-transposition typo (isiZulu)

000	
<u>File Edit Help</u>	
Run Save Copy Paste Clear Help IsiZulu V	
yebo, <u>mivah</u> eThekwini	Ignore once Ignore all Add to dictionary
	Suggestions
	Change all

Double click on errors to make correction one at a time

Error correction-transposition typo (isiXhosa)

000	
<u>Eile Edit H</u> elp	
Run Save Copy Paste Clear Help IsiXhosa Ebantwini, iinwele zikhula entloko ubukhulu becala, WKAYE Isixa seenwele zomzimba sahlukile kuhlanga nohlanga.	A/
	Ignore once Ignore all Add to dictionary
	kwaye Change
Double click on errors to make correction one at a time	Change all

Results

- It can propose something quite well for each type of typo (> 90% accuracy)
- The relevance varies a lot:
 - Substitutions 59%
 - Insertions 30%
 - Deletions 73%
 - Transpositions 89%

3

26 / 69

Results

- It can propose something quite well for each type of typo (> 90% accuracy)
- The relevance varies a lot:
 - Substitutions 59%
 - Insertions 30%
 - Deletions 73%
 - Transpositions 89%
- Why? We don't know for sure yet

Spelling correction for isiXhosa

- Used the same code as for isiZulu
- But with the isiXhosa language model
- Implemented, but no idea on effectiveness yet

isiZulu text, isiZulu model

000	
<u>File</u> <u>E</u> dit <u>H</u> elp	
Run Save Copy Paste Clear Help isiZulu KUBATSHAZWA Isihluku owesifazane oneminyaka engu-22 ediwengulwa amadoda ayisihlanu eshintshana ngaye eNseleni eMpangeni. Umphatathi wakule ndawo usababaza isigameko okuliwa senzeke ngolwesibili ebusuku. Kuthina owesifazane ubanjwe amadoda ate savapmbelezela esegoluka esuka ejoprimi. Ngokathola kweśolezwe to wesifazane takate <u>rezipibulisa</u> ejoprimin okuliwa ajfižule khona eMseleni. Kuthiwa upiscele umfowabo ngakho uno pomphasica ngoba oberujeka kili. Ngokathola kweśolezwe to wesifazane takate <u>rezipibulisa</u> , ejoprimin okuliwa ajfižule khona eMseleni. Kuthiwa upiscele umfowabo ngakho uno pomphasica ngoba oberujeka kili. Use sifazane kuthiwa upiante nabo.oka okuthe besendleleni bamikenga. Kuthiwa bahambe naje ngenkarhi. Neweifazane kuthiwa upiante nabo.oka okuthe besendleleni bamikenga. Kuthiwa bahambe naje ngenkarhi. Okuhumela amaphoysa KwaZulu-Natal uColonel Jay Naicker ukuginisekisle ukuthi amaphoyisa aseMpangeni aphenya icala lokuthnujwa nelokudiwengulwa abasowa abahlanu emua kokuba benthembise ukuthi bazamoduka, Akekho osabihwe. uphenyo wamaphoysa lokazhutheka. Kuthi ukuthi zingenejazo. "Owenglazane oneminyaka engu-22 uwle icala lokuthnujwa nelokudiwengulwa abasowa abahlanu emua kokuba benthembise ukuthi bazamoduka, Kusho Unakter. "Omenglazane oneminyaka engu-22 uwle icala lokuthnujwa nelokudiwengula. Kuyacara ukuthi zingenejazo. "Simangazwe yishiku sabasowa ngokusininshana ngengane bejdidengula. Kuyacara ukuthi zingenejazo. "Simangazwe yishiku sabasoku ngokusininshana ngengane b	Ignore once Ignore all Add to dictionary Suggestions No suggestions Change Change all
Double click on errors to make correction one at a time	Exit

isiZulu text, isiXhosa model

000	
Eile Edit Help	
Run Save Copy Paste Clear Help IsiXhosa KUBATSHAZMA IsiNkuu owesifazane oneminyaka engu-22 edivergulva amadoda ayisiNanu eshintshana ngaye eNseleni eMpangeni. Umphalani walue dnawo uzubabaza inginyako odurhan esareka ngalues piliti ehusuku. Kuthan owesifazane ubanjwe amadoda athe ayonphelezela asego ndu esuka sizimili ehusuku. Waluku owesifazane ubanjwe amadoda athe ayonphelezela asego ndu esuka sizimili ehusuku. Npokutola kwesifazane kuthwa ohunyuku okuthwa amenzakalisile kuvela ukuthi aye kuye amtshela ukuthi ayamazi umfowabo ngakho azomphelezela ngoba ubengakafiki. Lo wesifazane kuthwa ubandhe nabasolwa okuthe besendleleni bamikeda. Kuthiwa bahambe naye ngenkani bayomfaka endini abafike bamikengi kuvona. Kusen kuthwa ubandhe nabasolwa okuthe besendleleni bamikeda. Kuthiwa bahambe nabasolwa okuthe besendleleni bamikeda. Okukumela maphoyisa KwaZulu-Maalu (Colonel lay Naicker ukuqiniseksile ukuthi amaphoyisa aseMpangeni aphenya icala lokuthonyaka menzekalisile kusuku kuthi kanaboka asalwa abahanu emuva kokuba bemthembise ukuthi hazongodukwa kowesifazane. "Owesifazane oneminyaka engu-22 uule icala lokuthunjwa nelokudhwengulwa abaolwa abahanu emuva kokuba bemthembise ukuthi hazongodukwa onejukanisma negnape perkyidhengulak. Kutivaca aukuthi lingane zettu saziphenbile", kusho umthombo ongathandanga ukudaluku uthe lolu daba lukhulunyelwa phansi endawen injengoba luthatwa njengehlazo. Ommore umtombo uthe kuyadumaza ukuthi lokhu kwenzeka ngenyanga yabesifazane u-Agasti okumele engabe imbokodo <u>yavikelwa. Kukhulunye</u> kakhulu ngokuwkekwa kwabesifazane selokhu kuqale u-Agasti. Kuyadumaza ukuthi kukhona abesiisa abasahlukumeza abesifazane ngale <u>miela</u> kusho	Ignore once Ignore all Add to dictionary Suggestions No suggestions Change all Change all
Double click on errors to make correction one at a time	Exit

isiXhosa text, isiXhosa model

000	
<u>File Edit Help</u>	
Elle Edit Help Run Save Copy Paste Clear Help IsiXhosa ▼ Waggibela nini ukuza kwenza uvavanyo lomzimba? Matagibela ukwenza uvavanyo lomzimba? Matagibela ukwenza uvavanyo lomzimba? Ukusebrna kwegazi, JESC okanye luttra-soand Natagibela ukwenza uvavanyo kusha nje? Ukusebrna kwegazi, JESC okanye luttra-soand Ukusebrna kwegazi, JESC okanye luttra-soand Natagibela ukwenza uvavanyo kusha nje? Ukusebrna kwegazi, JESC okanye luttra-soand Ukusebrna kwegazi, JESC okanye luttra-soand Natagibela ukusebra kwegazi, JESC okanye luttra-soand Natagibela ukusebra kwegazi, JESC okanye luttra-soand Jeunakuke muse umkhono wasekhoho? Natagaberzu kama-80 Auwbonakali ngathi uyebe kakhulu, intle loo nto Ingaba uziboga njalo? Ukuba dimlyuka ndibaleka kazibenzusi, kundithatha ixesha ukuba ndifumane ukuphefumla kwam kwakhona.	Ignore once Ignore all
<u>Kudingek</u> a ndiphume ngokongezelelweyo. Ingalicebo elihile elo. Injani yona indiela otya ngayo? Ndicinga nditya ngendlela efanelekileyo. Ubazi ndiba mehamburgere du sesha nelo useha. Kodwa ngokuthe iikelele nditva ukutha okufanelekileyo.	Add to dictionary
Ngoku, ndiza kumamela intiiziyo yakho. Yinu, yabanda! Musa ukukhathazeka, <u>vistethoscope</u> yam je kuphela. Ngoku, phefumlela ngaphakathi uze <u>uwubambe</u> umphefumio wakho. Ndicela unyuse isheti yakho, upoferumile kakhulu. Yonke into ivakala kakuhle. Masikhe sijonge <u>umada</u> wakho. Ndicela uvule kakhulu uze uthi <u>'ah'.</u>	No suggestions
	Change
	Change all

Double click on errors to make correction one at a time

Exit

isiXhosa text, isiZulu model

Eile Edit Help	
Run Save Copy Paste Clear Help isiZulu •	Δ.
Ndaggibela ukwenza uzwaznyo lwam ikwiminyaka emibini edlulileyo. Ubuke wanalo olunye uzwaznyo kutsha nje? Ukusebenza kwegazi, <u>IEKC</u> okanye <u>uktra-sound</u> Akingatsho, <u>Pendikhe</u> ndenza <u>I-K-taray</u> ezimbalwa <u>kwaqqirha</u> wamazinyo. Ubuheli uzika njani ngokuthe jikelele? Akukho zikahaza, mahari	A/
Usunalukhe umuse umkhono wasekhohlo? Ndingwenela ukuhatha futhe lakho legazi. ngama_12_2 angaphezu (wama_60, Awubonakali ngathi utwebe kakhulu, intle loo nto Ingaba uziloonga njalo?	Ignore once
nazi, anungesno Ukuba animyka <u>anibaleka izitepusi, kundithatha ixesha</u> ukuba <u>ndifumane</u> ukuphefumla kwam kwakhona. Kudingeka <u>ndiphume ngokongezelekweyo.</u> Inglaicebo elihe elo. Injani yona indiela <u>atya</u> ngayo? <u>Ndicinga ndiya</u> ngendela e fanekkileyo.	Add to dictionary
Uyazi, <u>ndiba nehamburger</u> elo <u>xesha</u> , nelo <u>xesha</u> , kodwa ngokuthe jikelele, <u>nditya ukutva</u> okufanelekileyo. Ngoku, <u>ndrak</u> kumamela <u>inditzyo</u> yakho. <u>Yhu</u> , yabanda! Musa ukukhathazeka, <u>vistethoscope</u> yam je kuphela. Ngoku, <u>nbefumlela</u> ngaphakathi uze uwubambe umphefumlo wakho. <u>Ndicela umuse</u> isheti yakho, <u>uphefumle</u> kakhulu. Yonke into bakala kakuhie.	Suggestions No suggestions
Masikhe <u>slionge</u> umqala wakho. <u>Ndicela</u> uvule kakhulu uze uthi <u>'ah'.</u>	Change

Double click on errors to make correction one at a time

31/69

Exit

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?
- Corpus size? Genre?
- Timeliness of the text

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?
- Corpus size? Genre?
- Timeliness of the text
- Lowercase vs upper case (e.g., eGoli)

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?
- Corpus size? Genre?
- Timeliness of the text
- Lowercase vs upper case (e.g., eGoli)
- Sociolinguistics, if dialects have words written differently

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?
- Corpus size? Genre?
- Timeliness of the text
- Lowercase vs upper case (e.g., eGoli)
- Sociolinguistics, if dialects have words written differently
- Room for improvement on the corrector

- Setting the threshold is difficult
 - Cleaned data or noisy data?
 - Trigrams on text proper or on non-punctuation-marks?
 - Cleaned trigrams or not?
- Corpus size? Genre?
- Timeliness of the text
- Lowercase vs upper case (e.g., eGoli)
- Sociolinguistics, if dialects have words written differently
- Room for improvement on the corrector
- How much do the isiZulu and isiXhosa language models differ?

Outline



Context

• Language 'crash course'

Corpus-based spellcheckers

- Error detection for isiZulu and isiXhosa
- Error correction for isiZulu and isiXhosa
- Discussion

Rule-based NLG

- What is CNL, NLG?
- Generating basic sentences
- Language learning exercises

Short answer

- Ccontrolled Naural Language: constrain the grammar/vocabulary of a natural language
- Natural Language Generation: generate natural language text from structured data, information, or knowledge

Natural language interfaces with some CNL or NLG

- Many tools, webpages, etc. with some natural language component
- Querying of information in natural language (cf. a query language SQL, SPARQL)
- Business rules typically specified in a natural language
- etc.

Example: iCal calendar entry with canned text

	my collo	quium
	location	None
	all-day	0
	from	12/06/2014 01:00 PM
	to	12/06/2014 02:00 PM
	repeat	None ‡
	show as	Busy ‡
	calendar	Work ‡
	alarm	Message with Sound ‡ ⊲) Basso ‡
		1 hours before ‡
	alarm	None ‡
×	invitees	Add Invitees

Example: Saadiq Moolla's mobile healthcare app



Chest Pain

Have you had any recent pain in your chest? - Uke waba nobuhlungu esifubeni maduzane?

Does the pain radiate to your jaw, neck or arm? - Engabe ubuhlungu bakho bujikeleza emihlathini, emqaleni noma nasezingalweni?

Does anything precipitate or relieve the pain? - Ingabe ikhona into eyenza ubuhlungu buqhubeke noma eyehlisa ubuhlungu?

Dyspnoea



Home » History » Cardiovascular History

Chest Pain

Have you had any recent pain in your chest? - Ingaba kutshanje ukhe weva iintlungu esifubeni?

Does the pain radiate to your jaw, neck or arm? - Ingaba iintlungu zinwenwela emhlathini, entanyeni okanye engalweni?

Does anything precipitate or relieve the pain? - Ingaba ikhon@7 / 69

Example: Query formulation with Quelo [Franconi et al.(2010)]

I am looking for a car dealer. It should sell a new car. The body style of the new car should	
be an off-road car). The new car should run on a diesel. Its model should be a Range Rover.	



I am looking for a car	. It should run on a diesel.						
	∇	it should be equipped with an equipment	V	with an engine	F	V	with a diesel engine
	∇	it should be located in a country	V	with an optional feature	۲		· · · · · · · · · · · · · · · · · · ·
Scramble Clear Exe	∇	it should be produced by something	∣⊽	with a transmission system	F		Ready.

Example: Business rules and conceptual data models



Each Course is taught by at least one Professor Each Professor teaches at least one Course

The 'NLG pipeline'



 What structured data/info/ knowledge do you want to put into NL sentences?
 In what order should it be presented? 3. Which messages to put together into a sentence?

4. Which words and phrases will it use for each domain concept and relation?

5. Which words or phrases to select to identify domain entities?

6. Use grammar rules to produce syntactically, morphologically, and orthographically correct (and is also meaningful)

NLG, principal approaches to generate the text

- Canned text
- Templates
 - Notably for English [Fuchs et al.(2010), Schwitter et al.(2008), Third et al.(2011), Curland and Halpin(2007)],
 - but also other languages [Jarrar et al.(2006)]
- Grammar engines, such as [Kuhn(2013)], Grammatical Framework (http://www.grammaticalframework.org/), SimpleNLG

NLG, principal approaches to generate the text

- Canned text
- Templates
 - Notably for English [Fuchs et al.(2010), Schwitter et al.(2008), Third et al.(2011), Curland and Halpin(2007)],
 - but also other languages [Jarrar et al.(2006)]
- Grammar engines, such as [Kuhn(2013)], Grammatical Framework (http://www.grammaticalframework.org/), SimpleNLG
- \Rightarrow CNL, NLG

Business rules/conceptual data models and logic reconstruction

BR: **Each** Course is taught by **at least one** Professor FOL: $\forall x \text{ (Course}(x) \rightarrow \exists y \text{ (is_taught_by}(x, y) \land \text{Professor}(y)))$ DL: Course $\sqsubseteq \exists \text{ is_taught_by}.\text{Professor}$

```
<Constraint xsi:type="Mandatory"> <Constraint xsi:type="Mandatory">
<Text> -[Mandatory] Cada</Text>
<Description of the constraint xsi:type="Mandatory">
<Description of the constraint xsi:type="Mandatory">
<Text> -[Mandatory] Cada</Text>
<Description of the constraint xsi:type="Mandatory">
<Description of the constraint xsi:type="Mandatory">
<Text> -[Mandatory] Cada</Text>
<Description of the constraint xsi:type="Mandatory">
<Description of the constraint
```







NL Grammars, illustration

 $egin{array}{ccc} \textit{Noun} & \longrightarrow & \textit{car} \mid \textit{train} \ \textit{Adjective} & \longrightarrow & \textit{big} \mid \textit{broken} \end{array}$

. . .

. . .

(and complexity of the grammar)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三日 - シック

44 / 69



• Can the template-based approach be used also for isiZulu?

Question

• Can the template-based approach be used also for isiZulu?

- If so, create those templates
- If not, start with basics for a grammar engine
Question

- Can the template-based approach be used also for isiZulu?
 - If so, create those templates
 - If not, start with basics for a grammar engine
- Use a practically useful language to benefit both ICT and linguists and, possibly, some subject domain (e.g., medicine)
- Details in [Keet and Khumalo(2014b), Keet and Khumalo(2014a), Keet and Khumalo(2017)]

A logic foundation for isiZulu knowledge-to-text

- Roughly OWL 2 EL
- OWL 2 EL is a W3C-standardised profile of OWL 2
- Tools, ontologies in OWL 2 (notably SNOMED CT)

\mathcal{ALC} syntax

- Concepts denoting entity types/classes/unary predicates/universals, including top ⊤ and bottom ⊥;
- Roles denoting relationships/associations/n-ary predicates/properties;
- Constructors: and $\sqcap,$ or $\sqcup,$ and not $\neg;$ quantifications forall \forall and exists \exists
- *Complex concepts* using constructors: Let *C* and *D* be concept names, *R* a role name, then
 - $\neg C$, $C \sqcap D$, and $C \sqcup D$ are concepts, and
 - $\forall R.C$ and $\exists R.C$ are concepts
- Individuals

${\cal ALC}$ semantics

- domain of interpretation, and an interpretation, where:
 - Domain Δ is a non-empty set of objects
 - Interpretation: ${}^{\mathcal{I}}$ is the interpretation function, domain $\Delta^{\mathcal{I}}$
 - $\cdot^{\mathcal{I}}$ maps every concept name A to a subset $A^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$
 - ${}^{\mathcal{I}}$ maps every role name R to a subset $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} imes \Delta^{\mathcal{I}}$
 - $\cdot^{\mathcal{I}}$ maps every individual name *a* to elements of $\Delta^{\mathcal{I}}$: $a^{\mathcal{I}} \in \Delta^{\mathcal{I}}$

• Note:
$$\top^{\mathcal{I}} = \Delta^{\mathcal{I}}$$
 and $\perp^{\mathcal{I}} = \emptyset$

•
$$(\neg C)^{\mathcal{I}} = \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}}$$

- $(C \sqcap D)^{\mathcal{I}} = C^{\mathcal{I}} \cap D^{\mathcal{I}}$
- $(C \sqcup D)^{\mathcal{I}} = C^{\mathcal{I}} \cup D^{\mathcal{I}}$
- $(\forall R.C)^{\mathcal{I}} = \{x \mid \forall y.R^{\mathcal{I}}(x,y) \to C^{\mathcal{I}}(y)\}$
- $(\exists R.C)^{\mathcal{I}} = \{x \mid \exists y.R^{\mathcal{I}}(x,y) \land C^{\mathcal{I}}(y)\}$

Universal Quantification

- Consider here only the universal quantification at the start of the concept inclusion axiom ('nominal head')
- 'all'/'each' uses -onke, prefixed with the oral prefix of the noun class of that first noun (OWL class/DL concept) on lhs of ⊑

```
(U1) Boy ⊑ ...
wonke umfana ...
bonke abafana ...
('all boys...'; u- + -onke)
(U2) Phone ⊑ ...
lonke ifoni ...
onke amafoni ...
('all phones...'; i- + -onke)
```

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC}_{\mathbf{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	SO-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC_{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	SO-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a-onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i-onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC		\mathbf{QC} (all)		NEG SC	PRON	RC	QCdwa	EC
	$\mathbf{QC}_{\mathbf{oral}}$	-onke	QC_{nke}					
1	u-onke –	wonke	wo-	aka-	yena	0-	ye-	mu-
2	ba-onke	\rightarrow bonke	bo-	aba-	bona	aba-	bo-	ba-
1a	u-onke –	wonke	wo-	aka-	yena	0-	ye-	mu-
2a	ba-onke	\rightarrow bonke	bo-	aba-	bona	aba-	bo-	ba-
3a	u-onke —	wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	ba-onke ·	\rightarrow bonke	bo-	aba-	bona	aba-	bo-	ba-
3	u-onke —	• wonke	wo-	awu-	wona	0-	wo-	mu-
4	i-onke \rightarrow	yonke	yo-	ayi-	yona	e-	yo-	mi-
5	li-onke –	· lonke	lo-	ali-	lona	eli-	lo-	li-
6	a-onke —	onke	0-	awa-	wona	a-	wo-	ma-
7	si-onke –	≻ sonke	SO-	asi-	sona	esi-	so-	si-
8	zi-onke –	> zonke	zo-	azi-	zona	ezi	zo-	zi-
9a	i-onke \rightarrow	yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a-onke —	onke	0-	awa-	wona	a-	wo-	ma-
9	i-onke \rightarrow	yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi-onke –	> zonke	zo-	azi-	zona	ezi-	zo-	zi-
11	lu-onke -	→ lonke	lo-	alu-	lona	olu-	lo-	lu-
(10)	zi-onke –	> zonke	zo-	azi-	zona	ezi-	zo-	zi-
14	ba-onke ·	\rightarrow bonke	bo-	abu-	bona	obu-	bo-	bu-
15	ku-onke	\rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

Subsumption

- Two different ways of carving up the nouns to determine which rules apply: semantic and syntactic
- Need to choose between
 - singular and plural
 - with or without the universal quantification voiced
 - generic or determinate
 - (S1) MedicinalHerb \Box Plant

ikhambi <u>ng</u>umuthi amakhambi <u>yi</u>mithi <u>wonke</u> amakhambi ngumuthi

- (S3) Cellphone ⊑ Phone Umakhalekhukhwini uyifoni

('medicinal herb is a plant')

('medicinal herbs are plants')

('all medicinal herbs are a plant')

('giraffes are animals'; generic)

('cellphone <u>is a</u> phone'; determ.)

Possible subsumption patterns

- a. N_1 <copulative ng/y depending on first letter of $N_2 > N_2$.
- b. <plural of N_1 > <copulative ng/y depending on first letter of plural of N_2 ><plural of N_2 >.
- c. <All-concord for NC_x>onke <plural of N_1 , being of NC_x> <copulative ng/y depending on first letter of $N_2 > N_2$.

Subsumption: adding negation

- Need to choose between
 - singular and plural, and with or without the universal quantification voiced
- Copulative is omitted
- Combines the negative subject concord (NEG SC) of the noun class of the first noun (*aku*-) with the pronomial (PRON) of the noun class of second noun (*-yona*)

```
(SN1) Cup \sqsubseteq \negGlass
```

indebe <u>akuyona</u> ingilazi

zonke izindebe aziyona ingilazi

('cup not a glass')

('all cups not a glass')

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC}_{\mathbf{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	so-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i-onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i-onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC	QC (all)	QC (all)		PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC_{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	SO-	asi-	sona	esi-	so-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC_{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	D -	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	D-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	SO-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

Possible negation (disjointness) patterns

- a. $<N_1$ of NC_x> <NEG SC of NC_x><PRON of NC_y> $<N_2$ of NC_y>.
- b. <All-concord for NC_x>onke <plural N₁, being of NC_x> <NEG SC of NC_x><PRON of NC_y> <N₂ with NC_y>.

Existential Quantification

(E1) Giraffe ⊑ ∃eats.Twig

yonke indlulamithi idla ihlamvana <u>elilodwa</u> zonke izindlulamithi zidla ihlamvana <u>elilodwa</u> ('each giraffe eats <u>at least one</u> twig') ('all giraffes eat <u>at least one</u> twig')

a. <All-concord for NC_x>onke <pl. N_1 , is in NC_x> <conjugated verb> < N_2 of NC_y> <RC for NC_y><QC for NC_y>dwa.

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC}_{\mathbf{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	so-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	QC _{oral+onke}	$\mathbf{QC_{nke}}$					
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	SO-	asi-	sona	esi-	SO-	si-
8	$ ext{zi-onke} ightarrow ext{zonke}$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke \rightarrow konke	zo-	aku-	khona	oku-	zo-	ku-

NC	QC (all)		NEG SC	PRON	RC	QCdwa	EC
	$QC_{oral+onke}$	$\mathbf{QC}_{\mathbf{nke}}$				unu	
1	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
1a	u -onke \rightarrow wonke	wo-	aka-	yena	0-	ye-	mu-
2a	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3a	u -onke \rightarrow wonke	wo-	aka-	wona	0-	ye-	mu-
(2a)	$ba-onke \rightarrow bonke$	bo-	aba-	bona	aba-	bo-	ba-
3	u -onke \rightarrow wonke	wo-	awu-	wona	0-	wo-	mu-
4	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	mi-
5	$li-onke \rightarrow lonke$	lo-	ali-	lona	eli-	lo-	li-
6	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
7	$si-onke \rightarrow sonke$	so-	asi-	sona	esi-	SO-	si-
8	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi	zo-	zi-
9a	i -onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
(6)	a -onke \rightarrow onke	0-	awa-	wona	a-	wo-	ma-
9	i-onke \rightarrow yonke	yo-	ayi-	yona	e-	yo-	yi-
10	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
11	$lu-onke \rightarrow lonke$	lo-	alu-	lona	olu-	lo-	lu-
(10)	zi -onke $\rightarrow zonke$	zo-	azi-	zona	ezi-	zo-	zi-
14	$ba-onke \rightarrow bonke$	bo-	abu-	bona	obu-	bo-	bu-
15	ku -onke $\rightarrow konke$	zo-	aku-	khona	oku-	zo-	ku-

Example

- $\forall x \ (\operatorname{Professor}(x) \to \exists y \ (\operatorname{teaches}(x, y) \land \operatorname{Course}(y)))$
- Professor $\sqsubseteq \exists$ teaches.Course
- Each Professor teaches at least one Course

Example

- $\forall x (uSolwazi(x) \rightarrow \exists y (ufundisa(x, y) \land lsifundo(y)))$
- uSolwazi ⊑ ∃ ufundisa.lsifundo
- ?

$\forall x \text{ (uSolwazi}(x) \rightarrow \exists y \text{ (ufundisa}(x, y) \land \text{ lsifundo}(y))) \\ \text{uSolwazi} \sqsubseteq \exists \text{ ufundisa.lsifundo}$

$\forall x (uSolwazi(x) \rightarrow$	NC	AU	PRE	x,	$\frac{v}{v}$	lsifundo(v)))
	1	u-	m(u)-	ŗ.	NC	QC (all)
	2	a-	ba-			QC _{oral+onke}
	la	u-	-		1	u-onke \rightarrow wonke
look-up NC	$_{2a}$	0-	-		2	$ba-onke \rightarrow bonke$
pluralise	3a	u-	-		1a	u-onke \rightarrow wonke
	(2a)	0-	-		2a	ba-onke → bonke
for-all ———	3	u-	m(u)-		3a	u -onke \rightarrow wonke
	4	i-	mi-		(2a)	$ba-onke \rightarrow bonke$
	5	i-	(li)-		3	u-onke \rightarrow wonke
	6	a-	ma-	e	4	i-onke \rightarrow yonke
	7	i-	si-		5	$li-onke \rightarrow lonke$
	8	i-	zi-		6	a-onke \rightarrow onke
	9a	i-	-	-	7	$si-onke \rightarrow sonke$
	(6)	a-	ma-		8	$zi-onke \rightarrow zonke$
	9	i(n)-	-		9a	i-onke \rightarrow yonke
	10	i-	zi(n)-	·	(6)	a -onke \rightarrow onke
	11	u-	(lu)-		9	i-onke \rightarrow yonke
	(10)	i-	zi(n)-	ŀ	10	zi -onke \rightarrow zonke
	14	u-	bu-	•	11	$lu-onke \rightarrow lonke$
	15	u-	ku-	Ŀ	(10)	zi -onke \rightarrow zonke :
	17		ku-	L	14	$ba-onke \rightarrow bonke$
Bonke oSolwa	ızi				15	ku -onke \rightarrow konke
					< D	▶ ★課 ▶ ★注 ▶ ★注 ▶

≣ •⁄ < (~ 59 / 69

$$\forall x \text{ (uSolwazi}(x) \rightarrow \exists y \text{ (ufundisa}(x, y) \land \text{ lsifundo}(y)))}$$

$$\textbf{uSolwazi} \sqsubseteq \exists \text{(ufundisa)} \qquad \text{for relevant NC. Here:}$$

$$ngi-$$

$$u-$$

$$u-$$

$$u-$$

$$u-$$

$$si-$$

$$ni-$$

$$ba-$$



A PERSONAL A PERSON PERSON

$\forall x \text{ (uSolwazi}(x) \rightarrow \exists y \text{ (ufundisa}(x, y) \land \text{ lsifundo}(y)))$ uSolwazi $\sqsubseteq \exists$ ufundisa (lsifundo)





Bonke oSolwazi bafundisa Isifundo esisodwa

How to evaluate?

- Typical way of evaluating: ask linguists and/or intended target group
- Questions depend on what you want to know; e.g.,
 - Does the text capture the semantics adequately?
 - Must it really be grammatically correct or is understandable also acceptable?
 - Compared against alternate representation (figures, tables) or human-authored text?

How to evaluate?

- Typical way of evaluating: ask linguists and/or intended target group
- Questions depend on what you want to know; e.g.,
 - Does the text capture the semantics adequately?
 - Must it really be grammatically correct or is understandable also acceptable?
 - Compared against alternate representation (figures, tables) or human-authored text?
- Survey, asked linguists and non-linguists for their preferences
- 10 questions pitting the patterns against each other
- Online, with isiZulu-localised version of Limesurvey

The NLG algorithms can be used elsewhere

- Paper-based language learning exercises
- Exercise books have a lot of exercises on 'give plural noun', 'complete verb' etc

The NLG algorithms can be used elsewhere

- Paper-based language learning exercises
- Exercise books have a lot of exercises on 'give plural noun', 'complete verb' etc
- Our algorithms already can do that!
- Reuse the algorithms to pluralise and conjugate
- Proof of concept tool, tried to use both NLP (corpus, POS tagger) and the grammar engine of NLG

Examples of the CNL it uses

- Pluralise subject
 - Q: * Umfowethu bayaphuza
 - A: Abafowethu bayaphuza [prefixSG+stem] [PLSC+VerbRoot+FV] [prefixPL+stem] [PLSC+VerbRoot+FV]

Examples of the CNL it uses

- Pluralise subject
 - Q: * Umfowethu bayaphuza
 - A: Abafowethu bayaphuza [prefixSG+stem] [PLSC+VerbRoot+FV] [prefixPL+stem] [PLSC+VerbRoot+FV]
- Negate the verb
 - Q: Batotoba
 - A: Abatotobi

[PLSC+VerbRoot+FV] [PLNEGSC+VerbRoot+NEGFV]

Examples of the CNL it uses

- Pluralise subject
 - Q: * Umfowethu bayaphuza
 - A: Abafowethu bayaphuza [prefixSG+stem] [PLSC+VerbRoot+FV] [prefixPL+stem] [PLSC+VerbRoot+FV]
- Negate the verb
 - Q: Batotoba
 - A: Abatotobi
 - [PLSC+VerbRoot+FV] [PLNEGSC+VerbRoot+NEGFV]
- Possible to combine components for new exercises
 - [prefixSG+stem] [SGSC+VerbRoot+FV] [prefixSG+stem] [prefixPL+stem] [PLNEGSC+VerbRoot+NEGFV] [prefixPL+stem]
 - Q: umfowethu usula inkomishi '(my) brother washes the cup'
 - A: abafowethu abasuli izinkomishi '(my) brothers do not wash the cups'

Outline

- Motivation
 - Context
 - Language 'crash course'
 - Corpus-based spellcheckers
 - Error detection for isiZulu and isiXhosa
 - Error correction for isiZulu and isiXhosa
 - Discussion
- 3 Rule-based NLG
 - What is CNL, NLG?
 - Generating basic sentences
 - Language learning exercises



Summary

- Some natural language understanding, some generation
- N-grams and learning from a corpus (spellchecker)
- Corpus affects how well the tool performs (spellchecker, isiZulu CALL)
- Templates inapplicable to isiZulu due to its grammar (OWL verbalisation), hence a tailor-made grammar engine
- NLG algorithms generic and modularised in the sense that they can be reused in other tools (CALL exercises)
- Not addressed much now, but no less important: underresourced language

Summary

References I



Sonia E. Bosch and Roald Eisele.

The effectiveness of morphological rules for an isiZulu spelling checker. South African Journal of African Languages, 25(1):25-36, 2005.



M. Curland and T. Halpin.

Model driven development with NORMA.

In Proceedings of the 40th International Conference on System Sciences (HICSS-40), pages 286a–286a, IEEE Computer Society, 2007. Los Alamitos, Hawaii,



Enrico Franconi, Paolo Guagliardo, and Marco Trevisan.

An intelligent query interface based on ontology navigation. In Workshop on Visual Interfaces to the Social and Semantic Web (VISSW'10), 2010.

Hong Kong, February 2010.



Norbert E. Fuchs, Kaarel Kaliurand, and Tobias Kuhn.

Discourse Representation Structures for ACE 6.6.

Technical Report ifi-2010.0010, Department of Informatics, University of Zurich, Zurich, Switzerland, 2010.



Mustafa Jarrar, C. Maria Keet, and Paolo Dongilli.

Multilingual verbalization of ORM conceptual models and axiomatized ontologies. Starlab technical report, Vrije Universiteit Brussel, Belgium, February 2006. URL http://www.meteck.org/files/ORMmultiverb JKD.pdf.

C. M. Keet and L. Khumalo.

Toward a knowledge-to-text controlled natural language of isiZulu. Language Resources and Evaluation, 51(1):131-157, 2017. doi: 10.1007/s10579-016-9340-0.
References II



C. Maria Keet and Langa Khumalo.

Toward verbalizing logical theories in isiZulu.

In B. Davis, T. Kuhn, and K. Kaljurand, editors, Proceedings of the 4th Workshop on Controlled Natural Language (CNL'14), volume 8625 of LNAI, pages 78–89. Springer, 2014a. 20-22 August 2014, Galway, Ireland.

C. Maria Keet and Langa Khumalo.

Basics for a grammar engine to verbalize logical theories in isiZulu.

In A. Bikakis et al., editors, Proceedings of the 8th International Web Rule Symposium (RuleML'14), volume 8620 of LNCS, pages 216–225. Springer, 2014b. August 18-20, 2014, Prague, Czech Republic.



Langa Khumalo.

Advances in developing corpora in African languages. *Kuwala*, 1(2):21–30, 2015.



Tobias Kuhn.

A principled approach to grammars for controlled natural languages and predictive editors. Journal of Logic, Language and Information, 22(1):33–70, 2013.

B. Ndaba, H. Suleman, C. M. Keet, and L. Khumalo.

The effects of a corpus on isizulu spellcheckers based on n-grams. In Paul Cunningham and Miriam Cunningham, editors, *IST-Africa 2016*. IIMC International Information Management Corporation, 2016. 11-13 May, 2016. Durban, South Africa.

References III



D. J. Prinsloo and G.-M. de Schryver.

Spellcheckers for the south african languages, part 2: the utilisation of clusters of circumfixes. South African Journal of African Languages, :83–94, 2004.



R. Schwitter, K. Kaljurand, A. Cregan, C. Dolbear, and G. Hart.

A comparison of three controlled natural languages for OWL 1.1. In *Proc. of OWLED 2008 DC*, 2008. Washington, DC, USA metropolitan area, on 1-2 April 2008.

Sebastian Spiegler, Andrew van der Spuy, and Peter A. Flach.

Ukwabelana – an open-source morphological zulu corpus. In Proceedings of the 23rd International Conference on Computational Linguistics (COLING'10), pages 1020–1028. Association for Computational Linguistics, 2010. Beijing.

Allan Third, Sandra Williams, and Richard Power.

OWL to English: a tool for generating organised easily-navigated hypertexts from ontologies. poster/demo paper, Open Unversity UK, 2011.
10th International Semantic Web Conference (ISWC'11), 23-27 Oct 2011, Bonn, Germany.

Thank you!

Questions?

Spellchecker with contributions from:

- IsiZulu Linguist: Langa Khumalo
- IsiZulu spellchecker
 - Hussein Suleman, Balone Ndaba, Langa Khumalo, Norman Pilusa, Frida Mjaria
- IsiXhosa spellchecker
 - Nthabiseng Mashiane, Siseko Neti, Norman Pilusa
- Participants in the evaluations from ULPDO@UKZN and Linguistics@UCT
- Additional texts from INC, Mantoa Motinyane-Masoko, MeMaT translators, publicly available texts