

Surface realization for Nguni languages by using weather summary production

Zola Mahlaza

Supervisor : Dr C. Maria Keet

July 8, 2016

1 Introduction

The history of South Africa and its general mood towards languages has resulted in human language technologies (HLTs) being seen as a priority by the government [12]. South Africa is a multilingual country with approximately 25 spoken languages, and 11 official languages. According to the country's 2011 census [1], three languages with the most first language speakers are isiZulu, isiXhosa and Afrikaans (in that order). English is only the fourth on the list, yet it is the main language used in public and official discourse for most institutions. This is because South African languages of foreign origin have seen significant investment, often times at the expense of other South African languages. Corrective measures are now being taken hence the investment in HLT for other South African languages. This is important because the country's relatively new constitution states that every South African citizen has the right "to receive education in the official language or languages of their choice in public educational institutions where that education is reasonably practicable" [7]. The same document also states that it is the responsibility of the State to recognize "the historically diminished use and status of the indigenous languages of our people" [8] and therefore it "must take practical and positive measures to elevate the status and advance the use of these languages" [8]. In the quest to fulfill this responsibility, according to Grover et al [12], there have been projects dealing with human languages funded by the South African government through the De-

partment of Arts and Culture (DAC), Department of Science and Technology (DST), and National Research Foundation (NRF). The stipulations from the constitution highlighted above could be achieved through the use of computers or more specifically, machine translation.

Natural language processing (NLP) tools and other HLTs that work with South African languages could have a very positive impact in the country. An area of NLP that is of interest to us is natural language generation (NLG). It is the study of techniques involved in the production of natural language texts from machine representations of knowledge. It has made it possible for computers to explain medical data to patients, summarise statistical data, etc [24, p2]. These are the benefits we would like to see for South African languages. However, there hasn't been a great deal of academic literature on NLG focusing on Nguni languages. This a group comprised of four languages which are spoken by a large population in Southern Africa. Our work is an attempt to lay the groundwork for NLG in Nguni languages by focusing on the methods for auto-generation of weather reports for isiZulu and isiXhosa. This goal can be made much more difficult due to the complex nature of these languages. NLP experts address the complexity of natural languages by first focusing on controlled languages. These are domain specific languages with a simplified grammar and restricted vocabulary. This is the reason our work focuses on weather generation. The use of this domain

restricts the language constructs thus making the scope of the project manageable. In addition to using domain restriction, it is important to note that not all four languages can be dealt with at the same time. We will focus on isiZulu and isiXhosa. These are languages which have similarities, like the other Nguni languages. It is often assumed that the aforesaid similarities can be exploited when building Nguni language technologies. This assumption is the motivating factor to the work done by Pretorius and Bosch [23] on their isiZulu morphological analyser, ZulMorph. Our work will attempt to quantify the similarity between isiXhosa and isiZulu. This will be done by determining whether it is possible to use the same grammar rules for both languages when generating sentences in a specific domain. Finally, we will investigate the improvement that is brought on, if any, by the use of phonological conditioning rules in NL generation.

The remainder of this document is structured as follows; We will discuss existing NLG systems and the overall approaches taken when building them in section 2. We will go deep and look at the models that go into building an NLG system, and also focus on the existing work that looks at similarities between isiZulu and isiXhosa in section 2.1. We will finalize our focus at these systems by highlighting existing systems in the weather domain in section 2.2. In sections 3, 4 and 5 we will discuss the problem we face in respect to the generation of text in Nguni languages. We will also provide the aim of this work and present the research questions we aim to answer. Lastly, in section 6 we will provide the strategy that we will take in trying to answer our research questions.

2 Related work

Natural language generation has uses in a variety of domains. The Bateman & Zock [35] list of NLG systems shows a variation of systems that exist in a number of fields ranging from systems which can produce flight information [2] to systems which can produce biographies [28]. The work done by Reiter and Dale ([24],[25]) details the principles of

building NLG systems. They discuss not only the required material when building NLG systems, but they also go into detail explaining the tasks involved in the process of converting data into text. They present the architecture made up of three major steps. These are document planning, micro-planning and realisation. The pipeline encompasses a large number of steps which are not all compulsory for every NLG systems. These are content determination, discourse planning, sentence aggregation, lexicalisation, referring expression generation, syntax and morphology.

There has been variation in the approaches of developing NLG systems. Cimiano et al [6] point out that first generation systems used to rely on template-based approaches. They were followed by the use statistical architectures. Current systems, like those developed by Cimiano et al [6] and their contemporaries, use a combination of these two techniques. The work done by Cimiano et al [6] is the development and evaluation of a system capable of converting RDF (Resource Description Framework) data to a natural language. This work depends on the use of an ontology and an ontology lexicon. They follow a three step data-to-text pipeline which comprises of document planning, microplanning and surface realisation. Their hybrid development approach is such that only the last two stages make use of statistical techniques.

The work relies heavily on the ability to represent lexicon data through lemon (Lexicon model for ontologies) and the use of a lexicon database to retrieve inflectional variants. They also have access to a large domain corpus. There are issues with their approach which may be problematic when working with Nguni languages. For instance, in referring expression generation they use a rule based approach. When it comes to the use of pronouns, they use the following rule:

A re-occurring ingredient is replaced by a pronoun if there is no other ingredient mentioned in the previous sentence that has the same number and gender [6, p13]

There are issues with the highlighted basis of their work when it comes to Nguni languages. A simple rule like the one above for referring expression generation will need refinement in order to be able to result in comprehensible Bantu sentences. In isiZulu for instance, Wilkes [32] argued that the notion that pronouns can replace nouns in sentences is misguided. The author pointed out that linguists who believe that this was the case classified pronouns into the following categories; absolute, demonstrative, quantitative, and qualificative pronouns. The much recent work done by Twala [31] captured the dissimilar views on the matter and showed that there were some authors who argued that even though a noun can appear in a sentence as a replacement of the noun, “the basic position of the pronoun is before the noun” [31, p23]. The overall categories pointed out by Twala [31] were; absolute, possessive, demonstrative, quantitative and the demonstrative copulative pronouns. The common thing with the pronouns, regardless of classification, is that they have a concordial nature which makes them more complex than English pronouns. Hence a rule based on swapping the noun with the pronoun becomes much more complicated to put into effect for Bantu languages. Consider the following example, in order to understand the problem with the rule proposed by Cimiano et al [6], that uses three short sentences with two ingredients. Here we see that the two ingredients (onion and oil) satisfy the requirement for replacing the second ingredient with a pronoun. The sentences are in English (1) and isiXhosa (2).

1. Chop one onion. Pour the *oil* in a pan.
Stir *it* until it is warm.
2. Nqunqa itswele elinye. Galela *amafutha* kwipani. [Wa]zamise [*wona*] ade abeshushu.

The introduction of the pronoun in last sentence has not made the sentence much more human-like. In this particular case, the omission of the pronoun is the best recourse. This is because the pronoun is in agreement with verb through its prefix thus the pronoun is implied. The problem, however, is that the rule does not necessarily always hold hence

encoding it would prove to be difficult. Finally, the noun class of the noun “amafutha” (class 6) determines the prefix (wo-) used for the stem (-na) of the pronoun. The gains of encoding the relationship between the prefix and stem for all the pronouns is not worth the effort at this point.

Furthermore, at the moment there is no database to retrieve inflectional variants for South African Bantu languages. Lastly, Chavula et al [5] have shown that lemon, as-is, does not work with Bantu languages. They have proposed the use an additional ontology in order to scaffold lemon in order for it to be able to deal with the noun class systems which are a feature of Bantu languages.

An example of a system built on statistical approaches, and not a hybrid approach like the work of Cimiano et al [6], is MOUNTAIN [18]. This is a language generation system that depends on statistical techniques and the availability of natural corpora. The system is built for dialog systems and according to the authors, such systems have several differences from general NLG systems. These differences are a result of the variety of requirements for dialog and text generation systems. Dialog systems require short utterances and text generation systems require long sentences. Other differences are related to functionalities and the architecture. The given reasons combined with the fact that the building of a corpus for training and testing the proposed system makes the MOUNTAIN approach not feasible.

The variety of systems, both in approach and function, means that there is a difference in the inputs expected by each system. For instance, Klein [17] developed a system capable of summarising essay paragraphs. The input to the aforementioned system is a paragraph of text. Davey [9] developed an NLG system whose input is a game of tic-tac-toe and produces English text describing the current state of the game. FOG [11], the weather NLG system which resides in the Forecast Production Assistant (FPA) system is comprised of three steps; data extraction, conceptual processing, and linguistic processing. All

the steps in the work done by Cimiano et al [6] would make up the linguistic processing in FOG. The inputs to FOG are charts which are developed by a forecaster on the FPA. Finally, the BabyTalk project developed by Portet et al [22] takes in 45 minutes' worth of "continuous physiological signals and discrete events" which are taken from hospital equipment to produce summaries for a neonatal intensive care unit. This wide variety shows us that NLG systems work with different types of data hence the forms of representation of said data should be dependant on the system. Common forms include conceptual graphs, RDBMSs and RDF/RDFS.

Bouayad-Agha et al [3] (how to cite, got in press version) say that in spite of the numerous kinds of possible inputs to NLG systems, the 'natural' inputs are semantic/conceptual representations. This is because the said inputs result in more flexible NLG systems. Furthermore, the authors point out that the pipeline architecture presented by Dale and Reiter is not the only architecture in existence. There exists revision-based approaches, optimization approaches and monolithic approaches which map content into text [3, p3]. In all of these approaches, the complexity of the system is determined by the triplet; inputs, context and output. There are other minor issues which contribute to the complexity, such as whether the system exists in isolation or within a larger system. What is relevant, however, is that the complexity brought on by the output is due to the output being affected by variables such as it's size, coherence, fluency, language and modality [3, p2].

2.1 Foundations and Language similarities

A key aspect that cannot be forgotten is that NLG does not only require information pertaining to the application domain, Reiter and Dale [24] point out that it also requires knowledge about the language. It is for this reason that the work done by Twala [31] is relevant. In her dissertation, she discusses the evolution of the grouping of nouns in isiZulu by looking at the groupings presented by numerous authors over the years. The author also provides

a comprehensible overview of the morphological, syntactic and semantic details for each noun class.

There are similarities between isiZulu and isiXhosa, but, be that as it may, the pond of academic literature attempting to quantify and/or document these differences has been moderately dry. An important notion that is brought forward by the likeness of Bantu languages is the generalization of techniques which currently apply to a specific language within this group to other Bantu languages. The work done by L.Pretorius & S. Bosch [23], which falls under natural language understanding (NLU), is evidence to that. The work in question attempted to document some differences pertaining to morphotactics and morphophonological alternations between the two languages. Their goal was to bootstrap the development of an isiXhosa morphological analyser by using their current prototype of an isiZulu morphological analyser, ZulMorph. Their basic approach, they claim, is that they will "use the Zulu morphological structure wherever applicable and only extend the analyser to accommodate differences between the source language (Zulu) and the target language (in this case Xhosa)". This has implicit assumption the differences between the two languages are relatively small. The aspect of their work that is of interest to us is their enquiry into morphotactics.

Morphotactics refers to the ordering of morphemes when forming words. This plays a significant role in isiXhosa and isiZulu due to their agglutinative nature. The authors point out that there are some differences in the workings of affixes between the two languages. For instance, isiXhosa unlike isiZulu has a temporal form for verbs and its role is to indicate when an action occurs [23, p98-p99]. The simple example given by Pretorius & Bosch to illustrate this point, however, is wrong as it violates the juxtaposed vowel prohibition (unless the vowels are the same) in isiZulu and isiXhosa. A better example is given by Ma's thesis which shows future and past tense with the sentence "I have arrived in Grahamstown" [19, p6].

1. Ndi-fik-e e-Rhini (SC-arrive-PST)

LOC-Grahamstown)

2. Ndi-zaku-fik-a e-Rhini (SC-FUT-arrive-FV LOC-Grahamstown)

The sentences above are in isiXhosa and the abbreviations used in the above example are; SC the subject concord, PST is past tense, FUT is future tense, FV is final vowel and LOC is locative.

There are no methods for verbalising concepts from machine representations in isiXhosa, to our knowledge. This is why the work done by Keet and Khumalo [15] becomes a foundational aspect of our work. Keet and Khumalo [15] investigated the formation of methods which would allow one to be able to create an isiZulu controlled natural language (CNL). A CNL is a subset of a specific natural language. The difference is that the grammar and vocabulary is restricted. They have shown that a template based NLG system will not work with isiZulu, and correspondingly other Nguni languages, due to the complex noun class system. They developed verbalization patterns in isiZulu for logic constructs such as subsumption, negation, universal and existential quantification. Furthermore, they have also shown that a template based system that makes use of the developed verbalisation patterns will also not work - a “full-fledged grammar engine” [15, p23] is required. The transferability of these patterns to other Nguni languages such as isiXhosa seems like a possibility. This is due to similarities between the Nguni languages.

There already has been work done to take the methods developed for isiZulu and reuse them for another Bantu language, Runyankore. This is a language spoken in Uganda and other Central/East African countries such as Burundi and Kenya. Runyaronke shares similarities with isiZulu, and that can be deduced by virtue of these languages being in the same language group. The two languages also have some differences, for instance, isiZulu has five distinct tenses whereas Runyankore has fourteen [4, p2]. The verbalization patterns faced similar issues and according to Byagumisha et al these are due to factors such as “the noun class of the name of the concept, the category of the concept,

whether the concept is atomic or an expression, the quantifier use in the axiom, and the position of the concept in the axiom” [4, p7]. Nonetheless, their work provided more evidence that the bootstrapping approach when building human language technologies for Bantu languages significantly reduces development time and requires less resources.

Context free grammars (CFG) are entities which come from formal language theory. They are a set of rules which determine how to form sentences/words in a formal language. Formal languages are not necessarily equivalent to natural languages. CFGs can be used, however, to model a controlled natural language. A controlled natural language (CNL), as briefly mentioned above, is a subset of a natural language. Controlled natural languages not only restrict vocabulary, but according to Wyner et al [34], they also restrict morphological forms, grammatical constructions, semantic interpretations and pragmatics. The benefits of CNLs are that they can be realized through the use of grammar formalisms such as CFGs. There are also other forms of grammars which can be used towards this goal. Other examples include context sensitive grammars (CSG), probabilistic CFGs (PCFG), etc. CFGs and their variants, to our knowledge, see more use in NLG compared to CSGs. The scarcity of NLG systems that makes use of CSGs might be due to a number of issues. The main reason is the observation which was made by Simmons and Yu. They argued that context sensitive grammars were not attractive because they are conceptually and computationally difficult to deal with [29, p392]. It is for this reason that the class of grammars, mildly context-sensitive grammars, exists. This is a group of grammars which are more powerful than CFGs as they include the notion of context, however, do not face the same challenges as a CSG. Formalisms within this group, to name a few, include the tree-adjoining grammar (TAG), head grammar (HG) and the combinatory categorial grammar (CCG). The last of which was investigated by Karagol-Aya [14] who attempted using it to model the morphotactics and syntax of Turkish. CCGs are generally used when mapping natural language to a logic form in NLU. However,

Karagol-Aya [14] uses an adapted version of the “semantic head-driven bottom-up generation” [14, p139] algorithm to generate a natural language.

There have been alternative approaches towards surface realization. The building of simpleNLG introduced the idea of a realization engine. This is a tool whose function is to create the lexicon and syntactical representations along with actions to allow a developer full control over the realization process. The important distinction between it and a traditional surface realizer is that, according to Gatt and Reiter, it’s key feature is the ability to move the responsibility of “making appropriate linguistic choices given the semantic input” [10, p90] out of the realization tool and into the hands of the developer. This means that simpleNLG is therefore only responsible for mechanical tasks such as mapping the semantic input into a syntactic structure and then linearizing it. A consequential benefit of this is that the engine does not require a strict input formalism and this means that the developer can decide on any suitable representation for the realization process. This tool is made exclusively for English. It has been adapted to other languages such as French, German and Brazilian Portuguese.

2.2 Weather text generation

According to the Bateman and Zock [35] list, the automation of the production of weather summary text is the second most favourite application of NLG systems. It follows behind health-care/medicine. This is an observation that is confirmed by Sripada et al [30]. A joint initiative by the University of Montreal and the Canadian government’s Environment Canada have detailed the use of NLP in producing weather forecasts [11]. They have developed the Forecast Generator (FOG), a system capable of creating forecast test summaries from weather maps. FOG exists within a bigger system called the Forecast Production Assistant (FPA). The goal of the FPA is the automation of routing aspects of weather reporting in order to allow forecasters to focus on “scientific questions”[11, pg45]. The key aspect about the FOG is that it is bilingual. It produces

texts in English and French. FOG has three steps, and these are data extraction, conceptual processing and linguistic processing. The data extraction step is not of interest to us.

The authors of FOG faced the same two lexical challenges we currently face. These are deciding which professional words to use when describing weather concepts. The second is how to generate text in two languages from the same input. They dealt with the first challenge by using words which were decided upon by the forecasters. The second challenge was dealt with the introduction of an abstract interlingua that will capture the syntax irrespective of language. This interlingua will be used when generating a deep syntactic representation for each language. The deep syntactic representation is further used when determining the surface syntax. The introduction of different surface syntax representations is done because FOG will make use of independent language grammars to map surface syntax into text. The authors are quick to point out that the deep syntactic representation is capable of modelling English and French because these languages have a similar communication/semantic style. This representation is not guaranteed to be passable for other languages. The level of abstractness in the representation is completely in the control of the developer. An example can be in the representation of the two concepts; “lower” and “slower”. These two concepts are used when describing two different concepts like temperature and wind speed. They can both be modelled in a more abstract level with the concept “diminish”. This abstractness allows a common representation for two different concept types and therefore two different languages. It is not clear whether such a technique is necessary for isiZulu and isiXhosa. It would only be necessary in the event that isiXhosa and isiZulu cannot be realized using the same grammar rules.

They use a language model which is flexible in order to “accommodate a variety of forecast types, and regional needs and tastes”[11, pg50] and it allows better maintenance of the system. They also chose to use a model that supported

speech synthesis. The reason for this was that they needed the system to work with telephone-answering systems. This requirement is great as it leads to accessibility of weather services to people who are visually impaired. They made use of Meaning-Text Theory (MTT). MTT is not a new theory/model. It was first proposed by Igor Mel'čuk in his 1970 and 80's work ([20], [21]), as referenced by Kittredge et al [16]. A number of its qualities are not necessary when dealing with controlled natural languages. We will not be investigating it due to its complexity, our limited time frame and lack of literature pertaining to the use of MTT with Nguni languages. It might be, in the future, worth investigating the use of MTT for the generation of text in Nguni languages.

We currently tell the weather through predictions. The accuracy of said predictions sometimes varies. This is why the NLG system should be able account for this. This is evident in the work done by Sripada et al [30] for the UK's national weather service. This consideration to the varying accuracy is important because their work generates weather summaries for a number of days in the future. This means that the system needs to account for the loss in accuracy in the text it generates to makes sure that users of the system are not mislead and thus lose trust forecast generating system. Sripada et al [30] also faced the issue of a lack of a corpus when building the system. They dealt with this by obtaining text samples from experts and supplementing them by a domain language; weatherese. It is not clear whether or not their evaluation method verified that the respondents of their survey were not weather experts. Nonetheless, of 35 respondents, they reported that 97% were satisfied with the understandability of the of the forecast. This shows that the absence of a corpus does not always result in failure. The internal details and techniques used in the system are not given. However, we know that it is based on the Arria NLG engine [27]. The engine is a commercial product that is used in areas such as financial services, advertising and marketing, oil and gas, etc. The technical overview does claim that it is language-independent. The engine is based on the traditional techniques which have already been

discussed here. It has are other aspects/modules which are incorporated for commercial appeal. Unfortunately, there is no evidence to support the idea that it could be used with Bantu languages.

There also has been variation in the approaches in NLG systems built within the weather domain. The work done by Winkler et al [33] uses the idea of a catalogue-based system. They mention that the idea has been used in generating severe weather warnings before but has never been used with a complex sentence type and a bigger domain. They built a system for the Swiss avalanche bulletin capable of producing text in German, French, Italian and English. The system uses a collection of sentence templates where each sentence is split into at most 10 segments. This approach, however, is a more advance form of templates. We can therefore deduce that it will suffer the same constraints as templates when it comes to Bantu languages

3 Problem statement

In our examination of the current state and use of Nguni languages, we have observed that there is no fast and large scale producer, automated or otherwise, of weather summaries in said languages. This is due to several factors such as (1) There are multiple Nguni languages, each of which has numerous dialects and hiring human authors to interpret weather data and produce these summaries is expensive and inefficient, (2) There is no automated system to achieve the stated goal because of, among other things, the complexity of Nguni languages (which is due to their noun class systems and the setup of concordial agreement) and the shortage of computer scientists working with Nguni languages.

Furthermore, the syntactic similarity between isiZulu and isiXhosa has never been formally quantified. The qualification would us allow to determine whether the same grammatical constructs can be reused between the two languages.

4 Aim

The first aim is to attempt to eliminate the dependence on a specific language for NLG systems in order to focus on the dependence on a group of closely related languages, starting with isiZulu and isiXhosa. This will be done by investigating the similarity in the syntax of the languages and how it can be exploited. This will partially address the multiple language problem highlighted in the previous section. Additionally, we investigate the degree to which a set of phonological conditioning rules can improve generated isiXhosa sentences hence bringing us a step closer to an automatic natural language generation system for isiZulu/isiXhosa.

5 Research questions

1. How syntactically similar are isiZulu and isiXhosa? Can the same grammar rules be used for both languages to produce comprehensible sentences?
2. What is the degree of improvement in comprehensibility on the generated sentences can be brought on by the introduction of phonological conditioning rules?

6 Methods and anticipated outcomes

We will investigate technologies which are capable of sentence derivation using some context aware grammar. We will then pick a suitable grammar formalism and tools for generating sentences. This will be followed by the incremental development of a grammar for each language. This will be done until we're capable of generating $\geq 75\%$ correctness for the generated sentences. The posed first research question put forth the notion of similarity. This necessitates the use of a metric to determine this similarity. We will devise metrics which capture the effect of tense in the grammar, the order of certain sentence constituents and therefore semantic style. Furthermore, we will attempt to create a unified

grammar should our metrics show that there is a similarity. The scale of deciding whether a similarity exists will depend on the developed metrics.

The focus of our is in surface realization, a step which does not exist in isolation and as such, we will assume the pipeline NLG architecture which has three major steps; document planning, micro-planning and realisation. We will not focus our attention on building a system capable of generating text for different audiences. It is for this reason and others that are not mentioned that will see document planning and micro-planning will not be given much attention. Nonetheless, the construction of entities and relations to be used is still required and it will be done manually.

The most popular approach for the analysis of target text is corpus-based. We will make no attempt to use experts from organizations such as the South African Weather Service (SAWS). This is because the dependence on a second party to produce the text/terms is not guaranteed to be finished within a reasonable timeline. The provisional recourse is to make use of English sources of weather terms and language. We will study the work done by Reiter et al [26] which details the techniques they have used when choosing words to use for their generated weather summaries. We will manually translate the words by hand and supplement the translation with the use of literature [13] and, should it be necessary, use University of Cape Town (UCT) students who speak isiXhosa/isiZulu as a first language as translators.

We are aware that translations may sometimes distort meaning and fail to capture the essence of the original sentence. Language is not a fixed entity. The same language may sometimes vary across a group of people who have different social attributes such as ethnicity, religion, education, etc. This is a factor that will need to be considered when evaluating the correctness for the generated sentences. This is the reason why linguists at the School of African Languages and Literature at UCT will be used to determine correctness of the generated sentences.

The metrics which we will develop should not and will not be specific to the weather generation domain. The expected result should also reveal the importance and therefore the priority of phonological conditioning in generating Nguni sentences. The limitation, however, is that we will not quantify the lexical similarities between isiZulu and isiXhosa. Therefore, there will be no coefficient highlighting that similarity. Furthermore, the metrics we produce may take different forms from lexical similarity coefficients and could possibly have different ranges. Their nature cannot be predetermined as they are dependant on the grammar formalism, its representation and other details.

References

- [1] Statistics South Africa. 2011 census : Census in brief. http://www.statssa.gov.za/census/census_2011/census_products/Census_2011_Census_in_brief.pdf. Accessed: 03 May 2016.
- [2] Scott Axelrod. Natural language generation in the ibm flight information system. In *Proceedings of the ANLP-NAACL 2000 Workshop on Conversational Systems*, pages 21–26. Association for Computational Linguistics, 2000.
- [3] Nadjat Bouayad-Agha, Gerard Casamayor, and Leo Wanner. Natural language generation in the context of the semantic web. *Semantic Web*, 5(6):493–513, 2014.
- [4] Keet C.M. DeRenzi B. Byamugisha, J. Bootstrapping a runyankore cnl from an isizulu cnl. In *Proceedings of the 5th Workshop on Controlled Natural Language (CNL’16), Aberdeen, UK, 25-27 July 2016*. Springer, in press 2016.
- [5] Catherine Chavula and C. Maria Keet. Is lemon sufficient for building multilingual ontologies for bantu languages? In *Proceedings of the 11th International Workshop on OWL: Experiences and Directions (OWLED 2014) co-located with 13th International Semantic Web Conference on (ISWC 2014), Riva del Garda, Italy, October 17-18, 2014.*, pages 61–72, 2014.
- [6] Philipp Cimiano, Janna Lüker, David Nagel, and Christina Unger. Exploiting ontology lexica for generating natural language texts from rdf data. In *Proceedings of the 14th European Workshop on Natural Language Generation*, pages 10–19, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- [7] South African Constitution. Bill of rights, chapter 2, section 7-39. <http://www.justice.gov.za/legislation/constitution/chp02.html#sthash.zjzzdCaW.dpuf>. Accessed: 19 April 2016.
- [8] South African Constitution. Founding provisions, chapter 1, section 1-6. <http://www.justice.gov.za/legislation/constitution/chp01.html#sthash.DthAk1WL.dpuf>. Accessed: 19 April 2016.
- [9] Anthony Davey. *The formalisation of discourse production*. PhD thesis, The University of Edinburgh, 1974.
- [10] Albert Gatt and Ehud Reiter. Simplenlg: A realisation engine for practical applications. In *Proceedings of the 12th European Workshop on Natural Language Generation, ENLG ’09*, pages 90–93, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.
- [11] Eli Goldberg, Norbert Driedger, and Richard I Kittredge. Using natural-language processing to produce weather forecasts. *IEEE Expert*, 9(2):45–53, 1994.
- [12] Aditi Sharma Grover, Gerhard B Van Huyssteen, and Marthinus W Pretorius. The south african human language technology audit. *Language resources and evaluation*, 45(3):271–288, 2011.
- [13] Language Inc. Weather terminology (multilingual south africa). [https://www.language-inc.org/en/resources/Resources/Weather\%20Terminology \ %20\(Multilingual \](https://www.language-inc.org/en/resources/Resources/Weather\%20Terminology\%20(Multilingual\%20South\%20Africa))

- %20South\%20Africa).pdf. Accessed: 14 June 2016.
- [14] Burcu Karagol-Ayan. Morphosyntactic generation of turkish surface forms. In *ESSLLI Student Papers*, pages 137–144, 1999.
 - [15] C Maria Keet and Langa Khumalo. Toward a knowledge-to-text controlled natural language of isizulu. *Language Resources and Evaluation*, pages 1–27, 2016.
 - [16] Richard Kittredge, Lidija Iordanskaja, and Alain Polguère. Multilingual text generation and the meaning-text theory. *Proceedings of TMI-88, Pittsburgh, PA, June, 1988*.
 - [17] Sheldon Klein. Automatic paraphrasing in essay format. *Mechanical Translation and Computational Linguistics*, 8:68–83, 1965.
 - [18] Brian Langner, Stephan Vogel, and Alan W Black. Evaluating a dialog language generation system: comparing the mountain system to other nlg approaches. In *11th Annual Conference of the International Speech Communication Association 2010, Makuhari, Chiba, Japan, 26-30 September, 2010*, pages 1109–1112, 2010.
 - [19] Xiujie Ma. An analysis of temporal relations in languages : a comparative study of mandarin and isixhosa. Master’s thesis, Rhodes University, Faculty of Humanities, English Language and Linguistics, 2013.
 - [20] Igor A. Mel’čuk. Meaning-text Models - a Recent Trend in Soviet Linguistics. *Annual Review of Anthropology*, 10:27–62, 1981.
 - [21] Igor A. Mel’čuk and Aleksandr K. Žolkovskij. Towards a functioning ‘Meaning-Text’ model of language. *Linguistics : An Interdisciplinary Journal of the Language Sciences*, 8:10–47, 1970.
 - [22] François Portet, Ehud Reiter, Jim Hunter, and Somayajulu Sripada. Automatic generation of textual summaries from neonatal intensive care data. In *Artificial Intelligence in Medicine, 11th Conference on Artificial Intelligence in Medicine, AIME 2007, Amsterdam, The Netherlands, July 7-11, 2007, Proceedings*, pages 227–236, 2007.
 - [23] Laurette Pretorius and Sonja Bosch. Exploiting cross-linguistic similarities in zulu and xhosa computational morphology. In *Proceedings of the First Workshop on Language Technologies for African Languages*, pages 96–103. Association for Computational Linguistics, 2009.
 - [24] Ehud Reiter and Robert Dale. Building applied natural language generation systems. *Natural Language Engineering*, 3(01):57–87, 1997.
 - [25] Ehud Reiter, Robert Dale, and Zhiwei Feng. *Building natural language generation systems*, volume 33. MIT Press, 2000.
 - [26] Ehud Reiter, Somayajulu Sripada, Jim Hunter, Jin Yu, and Ian Davy. Choosing words in computer-generated weather forecasts. *Artificial Intelligence*, 167(1):137–169, 2005.
 - [27] Ehud Reiter, Y Sripada, and R Dale. Technical overview: The arria nlg engine. Technical report, Arria NLG, 2015. Accessed: 13 June 2016.
 - [28] Wendy Hall Paul H. Lewis David E. Millard Nigel R. Shadbolt Mark J. Weal Sanghee Kim, Harith Alani. Artequakt: Generating tailored biographies with automatically annotated fragments from the web. In *Proceedings of the ECAI 2002 Workshop on Semantic Authoring, Annotation & Knowledge Markup, Lyon, July 22-26, 2002*, 2002.
 - [29] Robert F. Simmons and Yeong-Ho Yu. The acquisition and use of context-dependent grammars for english. *Comput. Linguist.*, 18(4):391–418, December 1992.
 - [30] Somayajulu Sripada, Neil Burnett, Ross Turner, John Mastin, and Dave Evans. *Proceedings of the 8th International Natural Language Generation Conference (INLG)*, chapter A Case Study: NLG meeting Weather Industry Demand for

- Quality and Quantity of Textual Weather Forecasts, pages 1–5. Association for Computational Linguistics, 2014.
- [31] Edith Khanyisile Twala. The noun class system of isizulu. Master’s thesis, University of Johannesburg, 1992.
 - [32] Arnett Wilkes. Comments on the function of the abbreviated absolute pronouns in zulu grammar. *South African Journal of African Languages*, 7(4):137–142, 1987.
 - [33] Kurt Winkler, Tobias Kuhn, and Martin Volk. Evaluating the fully automatic multi-language translation of the swiss avalanche bulletin. In *Proceedings of the 4th International Workshop, CNL 2014, Galway, Ireland, August 20-22, 2014*, pages 44–54, 2014.
 - [34] Adam Wyner, Krasimir Angelov, Guntis Barzdins, Danica Damljanovic, Brian Davis, Norbert Fuchs, Stefan Hoefler, Ken Jones, Kaarel Kaljurand, Tobias Kuhn, Martin Luts, Jonathan Pool, Mike Rosner, Rolf Schwitter, and John Sowa. On controlled natural languages: Properties and prospects. In Norbert E. Fuchs, editor, *Controlled Natural Language: Workshop on Controlled Natural Language, CNL 2009, Marettimo Island, Italy, June 8-10, 2009. Revised Papers*, pages 281–289. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010.
 - [35] John Bateman & Michael Zock. Bateman/zock list of nlg systems. <http://www.nlg-wiki.org/systems/>. Accessed: 14 April 2016.