

# Natural Language Generation from Ontologies

Thomas Smith  
Electronics and Computer Science  
University of Southampton  
taes1g09@ecs.soton.ac.uk

## ABSTRACT

The field of Natural Language Generation (NLG) deals with transforming structured data input into an understandable natural language output, with the aim of communicating some specific facet of that data. Since the Semantic Web (SW) stores and references semantic data in a structured way, there are many obvious benefits to the use of NLG techniques to help author, understand and summarise SW data. This paper provides an overview of the current approaches to NLG from SW data, and covers both strengths and weaknesses of the state-of-the-art.

## Categories and Subject Descriptors

H.5.2 [Information Interfaces and Presentation]: User Interfaces—*Natural language*; I.2.7 [Artificial Intelligence]: Natural Language Processing—*Language generation*; I.2.4 [Artificial Intelligence]: Knowledge Representation Formalisms and Methods—*Semantic networks*

## General Terms

Algorithms, Design, Documentation, Theory

## Keywords

Semantic web, natural language generation, natural language interface

## 1. INTRODUCTION

The aim of this paper is to provide an overview of the state-of-the-art in approaches to Natural Language Generation (NLG) from Semantic Web (SW) data. Traditional NLG converts a given structured data input into a natural language output, which is dependent on the goal of the communication and the context. The format and characteristics of the input, and the requirements of the output typically govern the architecture of the generator: methods range from simple pre-built templates to complex rule- and projection-based summaries.

Traditional NLG is composed of 5 stages:

1) Content Selection: In order to return only the necessary information, the input is condensed to only relevant data. A bottom-up approach is an informed search of the most relevant nodes, most useful when the goal is to say everything there is to say about an object of the domain [3].

2) Discourse and Information Structuring: It is necessary to derive a coherent discourse structure for the content to be presented. One popular approach is *Rhetorical Structure Theory* (RST) [7], due to the clear distinction between main (nucleus) and supporting (satellite) arguments. This stage can be undertaken at the same time as content selection to help ensure that the nodes selected form a coherent structure, and can be packaged into units of content known as *Elementary Discourse Units* (EDUs)[4].

3) Aggregation: Redundant and repeated information may be removed, and certain concepts merged to provide a more fluent outcome.

4) Lexicalisation: A language-oriented ontology is used to map the semantic representation of the data onto an abstract linguistic representation - one popular method available in OWL-DL is the *Generalised Upper Model* (GUM) [1].

5) Realisation and Linearisation: The final natural language content is arranged, rendered and presented to the user.

## 2. NLG AND THE SEMANTIC WEB

The use of Semantic Web technologies can greatly simplify the traditional NLG process in a number of ways.

A common problem in NLG is the projection of the domain information onto the language information used for the NL rendering. Often solved by using a single bespoke combination knowledge base, however for SW data OWL/RDF provides support for modularisation of knowledge into separate layered ontologies that can represent the connections between domain and communication knowledge whilst still retaining a degree of portability for the communication layer. The communication layer may also contain supplementary information about the context or user profile, for the determination of particularly relevant content during the NLG process. Alternatively, some approaches do not make use of a separate communication layer ontology, and instead annotate the domain classes and individuals directly with lexical information such as noun phrases, gender and plural forms, using the same Semantic Web formalisms as the orig-

inal data. One such system is NaturalOWL [6], which uses traditional reasoning to help annotate OWL ontologies in a multilingual fashion. Finally, the labels of entities have been found by Power et al. [9] to be largely meaningful natural language words, and so domain-independent lexicon creation can be performed cheaply in exchange for a certain degree of fluency. The NLG process can also benefit from the availability of standard SW reasoners in order to perform inference on the domain data and deliver natural language content in a more fluent and natural manner <come back to this with CNL>. The content selection process in many NLG techniques can be made to scale up to the large quantities of data available in SW contexts - many approaches extract fragments of interest based upon a maximum distance from an original node [4]. The increasing standardisation of formats across domains in SW data means that a NLG system developed for one purpose has the potential to be made portable across domains or applications. Task-specific knowledge may be layered separately from the domain knowledge, for example using the GUM [1], and in order to further reduce the amount of task-specific knowledge, the generated language constructions may be reduced to a controlled subset known as a *Controlled Natural Language* (CNL). Generic CNL-grammars are defined in terms of a formal language (i.e., OWL) and linguistic expressions, and for each new domain only lexicons need be defined [9]. NLG techniques are also able to make use of existing linked data concepts in SW-data, for example by using existing provenance models or FOAF information to supplement information available directly from the domain of interest [2].

Natural Language Generation can be used to present Semantic Web concepts and data in a user-friendly and easily understandable manner. Therefore, there are two classes of users that can most benefit from the application of NLG techniques to SW technologies: those who create and populate ontologies, and end-users that make use of data published and communicated via the Semantic Web. The first group would make use of a *Natural Language Interface* (NLI) to present the information in natural language for ontology engineering [11] - typically they might be domain experts but not SW experts. The second group would be seeking specific information from the SW, and NLG techniques would be used to present that information in an appropriate format.

## 2.1 NLG for Ontology Engineering

In a NLI for ontology engineering, NLG techniques may be used to support and guide an author through the development of ontologies and the construction of valid queries. Providing natural language paraphrases of input definitions can help to highlight trivially satisfiable and other incorrect restrictions [10]. CNLs are useful in this context as they can provide unambiguous translations between SW concepts and natural language. As the NLI is presented during the creation of the SW resource, lexical annotations are often unavailable and so linguistic resources at this stage are typically derived automatically from the entity and relation labelling information available in the ontology. The content selection and discourse structuring steps typically present in NLG are largely unnecessary for the NLI, as it can be assumed that the user will want access to all of the data in the same structure as the source ontology. Some basic aggrega-

tion is generally possible in the form of redundancy elimination and restriction combination (e.g., min-cardinality and max-cardinality may be replaced by “between ... and ...”). Finally the ontology may be lexicalised and realised via a CNL such as *Attempto Controlled English* (ACE), using the Attempto Parsing Engine [5].

## 2.2 NLG for Publication

The use of NLG for publication of SW data is a more complex problem than the generation of an authoring NLI, as the scope of communication goals is far broader. Depending on the context and requirements of the end user, the desired output may be anything from the answer for a specific question to a comprehensive summary of the complete ontology. For example, NaturalOWL [6] generates personalised multilingual descriptions of individuals and classes from the domain of interest. The NLG tasks rely on annotations to the input OWL ontology that specify, among other things, and item’s interest for a given type of user. During content selection, the system orders information about items according to the properties of the current user profile, and then constructs a discourse structure using only a number of the most interesting facts. Items are annotated with the appropriate lexical information, and after aggregation into suitable complex sentences presented to the end user. Not all SW systems rely on prior annotation of the data however. ArtEquAKT is a system that harvests biographical information about artists in order to provide a personalised dynamic biography in response to user queries [12]. Candidate fragments are selected and analysed, knowledge triples are extracted, consolidated and stored, and a personalised biography is generated according to the information gathered, the core ontology and a set of human-generated templates.

## 3. FUTURE WORK

There are a number of issues that still impede efficient integration of NLG with SW technologies. Though work on the Generalised Upper Model [1] and its various successors have provided popular mechanisms for codifying NLG-relevant knowledge alongside SW data, there is still no approach-independent way to provide all relevant information, which hampers the portability of individual generators. There is also no standardisation for the interfaces of individual models for NLG. Efforts have been made to formally specify consensus input or output representations, e.g. *Reference Architecture for Generation Systems* (RAGS) [8], however the complexity of the proposed framework hindered adoption by the NLG community. The resulting lack of portability of NLG solutions for SW data prevents them from fully taking advantage of the opportunities provided by linked data [3].

## 4. CONCLUSIONS

The goal of the Semantic Web is to allow agents to perform complex data retrieval and processing tasks and present understandable results to the end user. Natural Language Generation provides a mechanism by which results may be clearly presented, and in some fields is already a viable option. Further work will improve the portability, efficiency and robustness of existing approaches, enabling Natural Language output for all domains and applications of the Semantic Web.

## 5. REFERENCES

- [1] J. A. Bateman, B. Magnini, and G. Fabris. The generalized upper model knowledge base: Organization and use. *Towards Very Large Knowledge Bases*, pages 60–72, 1995.
- [2] C. Bizer, T. Heath, and T. Berners-Lee. Linked data-the story so far. *International Journal on Semantic Web and Information Systems (IJSWIS)*, 5(3):1–22, 2009.
- [3] N. Bouayad-Agha, G. Casamayor, L. Wanner, and C. Mellish. Content selection from semantic web data. In *Proceedings of the Seventh International Natural Language Generation Conference*, pages 146–149. Association for Computational Linguistics, 2012.
- [4] Y. Dai, S. Zhang, J. Chen, T. Chen, and W. Zhang. Semantic network language generation based on a semantic networks serialization grammar. *World Wide Web*, 13(3):307–341, 2010.
- [5] J. L. De Coi, N. E. Fuchs, K. Kaljurand, and T. Kuhn. Controlled english for reasoning on the semantic web. In *Semantic techniques for the web*, pages 276–308. Springer, 2009.
- [6] D. Galanis, G. Karakatsiotis, G. Lampouras, and I. Androutsopoulos. An open-source natural language generator for owl ontologies and its use in protégé and second life. In *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics: Demonstrations Session*, pages 17–20. Association for Computational Linguistics, 2009.
- [7] W. Mann and S. Thompson. *Rhetorical Structure Theory: A theory of text organization*. Ablex Publishing Corporation, Norwood, New Jersey, 1987.
- [8] C. Mellish. Using semantic web technology to support nlg case study: Owl finds rags. In *Proceedings of the 6th International Natural Language Generation Conference*, pages 85–93. Association for Computational Linguistics, 2010.
- [9] R. Power and A. Third. Expressing owl axioms by english sentences: dubious in theory, feasible in practice. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 1006–1013. Association for Computational Linguistics, 2010.
- [10] A. Rector, N. Drummond, M. Horridge, J. Rogers, H. Knublauch, R. Stevens, H. Wang, and C. Wroe. Owl pizzas: Practical experience of teaching OWL-DL: Common errors & common patterns. In *Engineering Knowledge in the Age of the Semantic Web*, pages 63–81. Springer, 2004.
- [11] P. R. Smart. Controlled natural languages and the semantic web. 2008.
- [12] M. J. Weal, H. Alani, S. Kim, P. H. Lewis, D. E. Millard, P. A. Sinclair, D. C. D. Roure, and N. R. Shadbolt. Ontologies as facilitators for repurposing web documents. *International Journal of Human-Computer Studies*, 65(6):537 – 562, 2007.