

Semantic composition process in a spoken language understanding system

Frédéric Duvert¹, Marie-Jean Meurs¹, Christophe Servan¹, Frédéric Béchet¹,
Fabrice Lefèvre¹, and Renato de Mori^{1*}

LIA - University of Avignon, France
{frederic.duvert, marie-jean.meurs, christophe.servan,
frederic.bechet, fabrice.lefevre, renato.demori}@univ-avignon.fr

Abstract

A knowledge representation formalism for SLU is introduced. It is used for incremental and partially automated annotation of the MEDIA corpus in terms of semantic structures. An automatic interpretation process is described for composing semantic structures from basic semantic constituents using patterns involving constituents and words. The process has procedures for obtaining semantic compositions and for generating frame hypotheses by inference. This process is evaluated on a dialogue corpus manually annotated at the word and semantic constituent levels.

Keywords: Spoken language understanding, semantic structures, frames, conceptual decoding, semantic annotation, semantic inference.

1 Introduction

Semantics deals with the organization of meanings and the relations between signs or symbols and what they denote or mean. Spoken Language Understanding (SLU) is the interpretation of signs conveyed by a speech signal. Relations are represented by Knowledge Sources (KS) and applied by processes using control strategic knowledge. This task is difficult because meaning is mixed with other information, such as speaker identity or noise in the environment.

Natural language sentences are often difficult to parse and spoken messages are often ungrammatical. The knowledge used is often imperfect and the transcription of user utterances in terms of word hypotheses is performed by an Automatic Speech Recognition (ASR) system which makes errors. It was observed that an increase in precision may be achieved by computing a lattice of scored hypotheses of semantic constituents from a lattice of scored word hypotheses Raymond *et al.* (2006). Semantic constituent hypotheses are generated using stochastic finite state machines (FSM) along the line of research presented in Riccardi and Gorin (2000).

*This work is supported by the 6th Framework Research Programme of the European Union (EU), Project LUNA, IST contract no 33549. For more information about the LUNA project, please visit <http://www.ist-luna.eu/>.

This paper describes a novel semantic composition and evaluation process which composes semantic constituents into semantic structures. Constituents are generated by a translation process from word lattices. Constituents and words have links to patterns. When patterns match with features based on constituent and word hypotheses, structure building procedures are executed. Confidence values based on probabilities are used for selecting hypotheses. The approach has been tested on the fairly complex French MEDIA corpus, available through the ELDA corpus distribution agency.

2 The MEDIA corpus and the generation of basic constituent hypotheses

2.1 Corpus description

The MEDIA corpus Bonneau-Maynard *et al.* (2005) has been recorded using a *Wizard of Oz* system simulating a telephone server for tourist information and hotel booking. Eight scenario categories were defined with different levels of complexity. The corpus accounts 1257 dialogs from 250 speakers and contains about 70 hours of dialogs. The training portion of the corpus is conceptually rich with more than 80 basic concepts manually transcribed and annotated.

2.2 Conceptual decoding for generating basic constituents

The MEDIA corpus is annotated with basic semantic constituents but not with semantic structures. Basic semantic constituents are hypothesized and scored following the approach described in Raymond *et al.* (2006).

The conceptual decoding process is seen as a translation process in which stochastic Language Models are implemented by Finite State Machines (FSM) which output labels for semantic constituents. There is an FSM for each elementary conceptual constituent. These FSMs are transducers that take words at the input and output the concept tag conveyed by the accepted phrase. An HMM tagger, also encoded as an FSM is used to rescore every path in the word/concept graph. This HMM tagger is trained on the MEDIA training corpus. This approach is called an *integrated* decoding approach as the ASR and SLU processes are done together by looking at the same time for the best sequence of words and concepts. The result of the translation process is a *structured* n-best list of interpretations that can be seen as an abstraction of all the possible interpretations of an utterance.

3 Composing semantic relations into structures

Semantic structures can be derived from semantic knowledge obtained with a semantic theory. Examples are semantic networks to represent entities and their relations Woods (1975) or function/argument structures Jackendoff (1990). A convenient way for representing and reasoning about semantic knowledge is to represent it as a set of *logic formulae* from which computational structures such

as frames can be derived. A frame is a model for representing semantic entities and their properties.

Part of a frame is a data structure which describes the properties of a semantic structure, the constraints which should be respected by the values the property can assume, and procedures for obtaining property values from signs coded in the speech signal. By filling slots, frame instances are generated. Acceptable frames for the semantic representation of a domain can be characterized by a *frame grammar*.

4 Progressive annotation of the corpus in terms of semantic structures

A frame based KS was manually composed to describe the semantic composition knowledge of the MEDIA domain. Some frames describe generic knowledge like spatial relations, some others are application specific. These frames were defined according to the *Berkeley FrameNet* paradigm adopted in lun.

The MEDIA KS is composed of 21 basic frames with a total of 85 slots. The meaning representation language (MRL) contains conceptual constituents and semantic structure building procedures. These procedures are part of the semantics of the MRL. Semantic constituents and some words have links to patterns π_j . When a pattern matches with the incoming data, frame instantiations are created. Based on frame instances, inferences are performed. Different frames linked by relations may be instantiated by a single pattern.

An initial set of 463 turns from 15 dialogues was manually annotated. The *FrameNet* Lowe *et al.* (1997) annotation format was used. A frame visualization tool, called FriZ, dedicated to process speech dialogues was developed to support manual annotation and verification of subsequent automatic annotations. The average manual annotation time per dialogue is around 2 hours.

Patterns were generalized by progressively annotating data with available knowledge, evaluating confidence of the results and manually annotating samples with low confidence.

Each word/concept sequence is analyzed thanks to the logical rules developed on the MEDIA training corpus. These rules use the attributes, the values and the specifiers obtained in the first decoding phase in order to infer the frames. This operation could also benefit from information related to other speech events, for example to the speaker pitch or to the hypotheses generated in the previous dialogue turns (stored in an agenda). These sources of information are not yet integrated in the work described in this paper.

5 Experimental results

Tests were performed on a corpus of 1249 dialog turns for a total of 2938 constituents. For a word error rate of 30.3%, the attribute error rate is about 25%. Each further information (specifiers and normalized values) add roughly an extra 6% to the error rates. The Oracle error rates, obtained by manually selecting the best hypotheses in the n -best list of interpretations (with $n = 20$), are lower by an absolute 8% than the 1-best error rates.

The frame hypotheses obtained on the output of the interpretation process has also been evaluated in view of. Since manual frame annotations were not available for the test corpus, the manual annotations of words and concepts were used to derive a reference frame annotation. After the composition and inference knowledge described in the previous section has been applied, a random sampling on the test user turns was performed by two human experts to manually assessing the accuracy of the automatic structure annotation. An F-measure of 0.90 (0.96 precision and 0.85 recall) was measured on 100 turns when comparing manual annotations and automatic frame annotations of exact transcriptions. This high accuracy allows to use the automatically-derived annotations as reference annotations.

6 Conclusion

A knowledge representation formalism for SLU has been introduced. It has been used for incremental and partially automated annotation of the MEDIA corpus in terms of semantic structures. Automatic annotations were evaluated and submitted to a human expert where confidence was low. An automatic interpretation process has been introduced for composing semantic structures from basic semantic constituents using patterns involving constituents and words. The process has procedures for obtaining semantic compositions and for generating frame hypotheses by inference.

Results in terms of F-measures are presented showing that the knowledge and the process have good capabilities for producing semantic structure hypotheses. This research will be pursued by using structural semantic knowledge for selecting possible constituents beyond the 1-best hypothesis in the whole lattice of concept hypotheses.

References

- (), Project LUNA : www.ist-luna.eu.
- Hélène BONNEAU-MAYNARD, Sophie ROSSET, Christelle AYACHE, Anne KUHN, and Djamel MOSTEFA (2005), Semantic annotation of the French Media dialog corpus, in *Eurospeech*, Lisboa, Portugal.
- R. JACKENDOFF (1990), Semantic Structures, *The MIT Press, Cambridge Mass.*
- J.B. LOWE, C.F. BAKER, and C.J. FILLMORE (1997), A frame-semantic approach to semantic annotation, in *Proceedings of the SIGLEX Workshop on Tagging Text with Lexical Semantics: Why, What, and How?*, Washington D.C., USA.
- Christian RAYMOND, Frederic BECHET, Renato De MORI, and Geraldine DAMNATI (2006), On the use of finite state transducers for semantic interpretation, *Speech Communication*, 48(3-4):288–304.
- Giuseppe RICCARDI and Al GORIN (2000), Stochastic language adaptation over time and state in natural spoken dialogue systems, *IEEE Trans. on Speech and Audio Processing*, 8(1):3–10.
- W.A. WOODS (1975), *What's in a Link: Foundations for Semantic Networks*, Bolt, Beranek and Newman.