

RECENT ADVANCES IN JANUS: A SPEECH TRANSLATION SYSTEM

*M. Woszczyna, N. Coccaro, A. Eisele, A. Lavie, A. M. Nir, T. Polzin, I. Rogina,
C.P. Rose, T. Sloboda, M. Timita, J. Tsutsumi, N. Aki-Wibel, A. Wibel, W. Wörd*

Carnegie Mellon University
University of Karlsruhe

ABSTRACT

Present recent advances from our efforts in increasing coverage, robustness, generality and speed of JANUS, CMU's German-to-speech translation system. JANUS is a speaker-independent system which translates spoken utterances in German into one of German, English or French. The system has been designed around the task of continuous translation (CT). It has initially been built around a database of 12 read dialogs, encompassing around 500 words. We have since been working along several dimensions to improve coverage and to move toward spontaneous

INTRODUCTION

In this paper we describe recent improvements of the German-to-speech translation system. Improvements have been made mainly along the following dimensions: 1.) better context-dependent modeling in the speech recognition module, 2.) improved language models, smoothing, and word equivalence classes improve coverage and robustness of the sentences that the system accepts, 3.) an improved N-best search reduces run-time from several minutes to now real time, 4.) trigram and parser rescoring improves selection of suitable hypotheses from the N-best list for subsequent translation. On the machine translation side, 5.) a cleaner interlingua was designed and syntactic and domain-specific analysis were separated for reusability of components and improved translation, 6.) a semantic analysis module for semantic analysis.

The

pendent segment weights.

Error rates using context dependent phonemes are lower by a factor 2 to 3 for English (1.5 to 2 for German) than using context independent phonemes. Results are shown in table 1.

language model	English		German				
	PP	WA	PP	WA			
none	40.0	58.2	42.5	63.0			
word-pairs		28.9	83.4	20.8	89.1		
bigrams			16.2	92.6	18.3	93.7	
smoothed bigrams				18.1	91.5	28.1	93.7
after resorting						—	—

Table 1: Word Accuracy

The performance on German is significantly better than on English in all cases.

When the standard GLR parser fails on all sentence candidates, this robust GLR parser is applied to the best sentence candidate.

3.2 The Interlingua

The output of the parser, known as "syntactic structure", is then fed into a mapper to produce an Interlingua representation. For the mapper, we use a software tool known as Transformation Kit [10]. A mapping grammar with about 300 rules is written for the Conference Registration domain of English.

```

((PREV-UTTERANCE ((SPEECH-ACT*ACKNOWL) (VALUE*HELLO))
  (TIME*PRESENT)
  (PRIY
  ((DEFINITE+) (NUMBER*SG)
  (AN M-)
  (TYPE*CONFERENCE)
  (CONCEPT*OFFICE))
(SPEECH-ACT*IDENTIFY-OTHER)

```

Figure 2: Example: Interlingua Output

Figure 2 is an example of Interlingua representation produced from the sentence "Hello is this the conference office". In the example, "Hello" is represented as speech-act *ACKNOWLEDGEMENT, and the rest as speech-act *IDENTIFY-OTHER.

3.3 The Generator

The generation of target language from an Interlingua representation involves two steps. First, with the same Transformation Kit used in the analysis phase, Interlingua representation is mapped into syntactic structure of the target language. There are about 300 rules for the generation mapping grammar for Japanese. The generation mapping grammar is used to generate the target sentence given the Interlingua representation.

side there is a “built-in” robustness against these phenomena in a connectionist system

The connectionist parsing process is able to combine symbolic information (e.g. syntactic features of words) with non-symbolic information (e.g. statistical likelihood of sentence types). Moreover, the system can easily integrate different knowledge sources. For example, instead of just training on the symbolic information, we trained PARSEC on both the symbolic information and the pitch contour. After training, the system was able to use the pitch contour to determine the sentence structure. These results were