

CONFIDENCE MEASURES FOR SPOKEN DIALOGUE SYSTEMS

Rubén San-Segundo¹, Bryan Pellom, Kadri Hacioglu, Wayne Ward

Center for Spoken Language Research. University of Colorado
Boulder, Colorado 80309-0594, USA, <http://cslr.colorado.edu>

José M. Pardo

Grupo de Tecnología del Habla. Universidad Politécnica de Madrid
Ciudad Universitaria s/n, 28040 Madrid, Spain, <http://www-gth.die.upm.es>

ABSTRACT

This paper provides improved confidence assessment for detection of word-level speech recognition errors, out of domain utterances and incorrect concepts in the CU Communicator system. New features from the speech understanding component are proposed for confidence annotation at utterance and concept levels. We have considered a neural network to combine all features in each level. Using the data collected from a live telephony system, it is shown that 53.2% of incorrectly recognized words, 53.2% of out of domain utterances and 50.1% of incorrect concepts are detected at a 5% false rejection rate. In addition, the confidence measures are used to improve the word recognition accuracy. Several hypotheses from different speech recognizers are compiled into a word-graph. The word-graph is searched for the hypothesis with the best confidence. We report a 14.0% relative word error rate reduction after this confidence rescaling.

1. INTRODUCTION

In a spoken dialogue system we can define three different levels for confidence measures:

- **Word Level:** in this level, confidence measures provide an idea about the accuracy of each recognized word. For this level we will use decoder and Language Model (LM) features.
- **Utterance Level:** here, the target is the detection of out of domain utterances. In this level we will use acoustic, LM and parsing features.
- **Concept Level:** The end-to-end system performance does not change in cases where the phrases in a sentence that belong to concepts are correctly recognized while the “filler” words or phrases are not correctly recognized. In this level, we focus on parts of phrases that are meaningful to the task. Decoder, LM and parsing features will be used to tag the concepts with the confidence measures.

We use the CU Communicator system as our test-bed for experimentation in this paper. This system is a Hub compliant implementation of DARPA Communicator task [1][2][3]. The

system combines continuous speech recognition, natural language understanding and flexible dialogue control to enable natural conversational interaction by telephone callers to access information about airline flights, hotels and rental cars.

2. DATABASE

The data used for the experiments has been obtained during the telephone data collection from November of 1999 through June of 2000 [3]. Over 900 calls were collected during this period totaling approximately 11,500 utterances. We have randomly split the data into three sets, 60% of the data for training, 20% for evaluation and 20% for testing. We have repeated it six times providing 6-Round Robin data sets to verify the results. The results presented are the average of these experiments. For the experiments in Sec. 3.3, we have used an independent set of utterances from the NIST Multi-Site Data Collection [3].

3. WORD LEVEL

For word level confidence we investigated a subset of the most promising features which were considered in [4][5][6]. We consider decoder and LM features.

Decoder features:

- **Normalized Score:** acoustic score of the word divided by the number of frames that it spans.
- **Count in the Nbest:** percentage of times the word appears in the 100-best hypotheses in similar position.
- **Lattice Density:** number of alternative paths to the word considered in the word-graph generated in the second pass of the recognizer.
- **Phone Perplexity:** average number of phones searched along the frames where the recognized word has been active in the decoding process.

Language Model features [7]:

- **Language Model Back-Off Behavior:** back-off behavior of an N-gram language model along a 5 word context.
- **Language Model Score:** the log-probability for each word in a sequence as computed from a back-off language model along a 5 word context.

¹ This work was performed during a visiting internship at the Center for Spoken Language Research and has been supported by Spanish Education Ministry grant AP97-20257252 and by DARPA through ONR grant #N00014-99-1-0418.

3.1 Feature combination

We have considered a Multi-Layer Perceptron (MLP) to combine all the features. In this study, the features were quantized with 115 binary inputs. 10 bits per feature were used except for the 5-context LM back-offs where it was necessary to use only 5 bits to code all possible situations. We have coded the features considering more resolution in ranks with more training data. The hidden layer consisted of 30 units and one output node was used to model the word-level confidence. During weight estimation, a target value of 1 is assigned when the decoder correctly recognizes the word and a value of 0, is assigned during incorrect recognition (e.g., substitutions and insertions).

3.2 Experiments

Table 1 summarizes the correct detection rates for word-level recognition errors at false rejection rates of 2.5% and 5.0%. In this table we also present the minimum classification error and the baseline error that corresponds with the recognition rate. It can be seen that LM features provide better indicators for word-level confidence than decoder features. For example, using LM features alone, 42.0% of mis-recognized words were detected at a false rejection rate of 5%, similar to [7]. Using decoder features we only reject 28.5% of the mis-recognized words.

Word Level	Correct Detection Rates		Classification Error Baseline 19.0%
	2.5% FR	5.0% FR	
Decoder Features	16.9%	28.5%	17.5%
LM Features	28.3%	42.0%	15.0%
ALL Features	39.0%	53.2%	12.8%

Table 1. Correct detection of mis-recognized words at a 2.5% and 5.0% false rejection rate (FR) rate. Minimum classification error and Baseline Error (Substitutions and Insertions) are also shown.

The best results are obtained by combining all features. In this case we reject more than half of the incorrect words at 5% false rejection. These features reduce the classification error by 6.2%.

3.3 Combining multiple hypotheses

The CU Communicator utilizes parallel banks of recognizers to obtain hypotheses from speaker-independent and female adapted telephone acoustic models. In the next sections we present several methods to use confidence measures for combining hypotheses from different decoders to improve speech recognition accuracy.



Figure 2: Word-graph generated for the example in Figure 1. The bold arrows indicate the best hypothesis.

3.3.1 Flat List Confidence Rescoring (FLCR)

For each hypothesis output from the bank of decoders, we calculate the average confidence along the whole sentence. The hypothesis with the highest average confidence value is selected as the best hypothesis.

3.3.2 Word Graph Confidence Rescoring (WGCR)

In this case the key is to build a Word-Graph with all hypotheses and find the path along the graph with the highest average confidence value. This path can produce a new hypothesis, different from those used to build the graph. The idea is to pick up the best parts from different sentences.

Word-Graph Generation: For each decoder hypothesis we tag each word with its confidence value. For example, consider Fig.1 that shows three hypotheses with confidence values shown underneath. In this figure we represent each word as an edge and we put a node between two consecutive words.

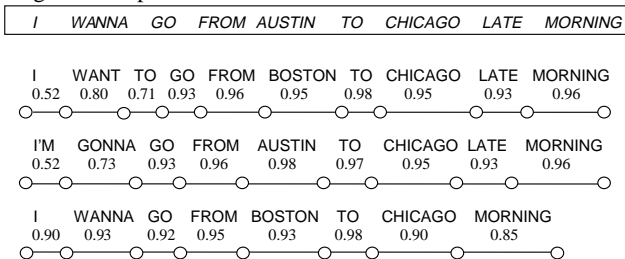


Figure 1. Reference utterance and three alternatives.

Next we join nodes from different hypotheses to build the graph. We join the beginning and end nodes of all the hypotheses and initial nodes of all words from different hypotheses situated in similar positions in the phrase. Finally we prune parallel transitions by picking the transition with the maximum confidence. The final word-graph is shown in Fig. 2.

Best Path Calculation: Using dynamic programming, the best path along the graph can be calculated. In this case our heuristic has been the *Accumulated Average Confidence*, i.e. the average confidence from the initial node to the current node. In Fig. 2 we can see the sentence finally obtained. Following the bolded edges we can see how each word was obtained from each hypothesis. In this case the resulting phrase matches perfectly with the reference and it is different from each individual hypothesis.

3.3.3 Experiments

We conducted experiments with four system configurations. In the baseline system that was evaluated by NIST [3], the CU Communicator initially runs two decoders in parallel. One decoder uses speaker independent models while the other uses female adapted models.

After 500 frames of input (5 sec.), the system selects one “best” decoder to use for the remainder of the telephone dialog. In this work, we have considered running decoders continuously in parallel throughout the dialog.

Word Error Rates for Four System Configurations	
Method	WER
Baseline	27.2%
Best Path Score	26.2%
FLCR	24.2%
WGCR	23.4%

Table 2. Word Error Rate Results.

In the second configuration, we select the hypothesis with the best path score as output from the decoder search. Finally, we consider using the proposed FLCR and WGCR methods.

Results are shown in Table 2.

From these results we can affirm that confidence measures provide an important role in reducing the WER by combining hypotheses from different decoders. The WGCR method proposed in this paper produces a reduction of 3.8 % in WER (14% relative) and performs better than the FLCR method. In these experiments the difference between FLCR and WGCR is low because the number of decoders run in parallel is small, only two, and the average utterance length is also small: 2.4 words per utterance.

4. UTTERANCE LEVEL

In this level we use all sources of information: decoder, LM and parsing features: Decoder and LM features:

- **Average Word Confidence:** this is the average word confidence along the sentence calculated in Sec. 3.

Parsing features:

- **Number of words parsed in the sentence:** number of words from the sentence belonging to a concept or a rule used to parse a concept.

- **Number of words that can be parsed:** number of words from the sentence belonging to a concept or any rule in the task grammar.

- **Number of Concepts:** number of concepts obtained in the sentence.

- **Average Count in the 100-best:** average percentage of times that a concept appears in the 100-best hypotheses.

- **Percentage of hypotheses in the 100-best with any concept:** with this feature we want to represent how many hypotheses parse with at least one extracted concept.

For combining all features we have considered an MLP. In this study, the features were not quantized and they were used as input to the MLP. Because of this, a preprocessing is required to limit the dynamic range of each feature to the (0,1) interval. Here, the normalization consists of scaling utilizing the minimum and maximum value obtained for each feature in the training set.

4.1 Experiments

Table 3 summarizes the correct detection rates for out of domain utterances at false rejection rates of 2.5% and 5.0%. It can be seen that parsing features provide better indicators for utterance-level confidence than decoder and LM features. These results are better than those obtained in [7].

Utterance Level	Correct Detection Rates		Classification Error Baseline 4.8%
	2.5% FR	5.0% FR	
Decoder + LM Features	41.1%	49.8%	4.2%
Parsing Features	43.8%	52.8%	4.1%
ALL Features	45.7%	53.2%	4.0%

Table 3. Correct detection of out of domain utterances at a 2.5% and 5.0% false rejection rate (FR) rate. Minimum classification error and Baseline Error are also shown.

5. CONCEPT LEVEL

We have calculated correct and incorrect concepts considering hypotheses and references, passing both through the parser and comparing them using a dynamic programming algorithm. Similar to the WER it is possible to calculate insertions, deletions and substitutions for concepts. In our system, we have a Concept Error Rate (CER) of 27.9%. The amount of incorrect concepts obtained is 16.5% (substitutions and insertions). In [8] we can see similar work in this level.

In this level we use all sources of information: decoder, LM and parsing features:

- **Average word confidence in the words belonging to the rule used to get the concept:** in this case we calculate the average word confidence obtained following the definition in Sec. 3 along the rule applied to get the concept.

- **Average word confidence for the value of the concept:** average word confidence for the words that can be considered as the value of the concept. Considering the sentence: “I wanna go to Chicago.”: the phrase “go to Chicago” contain the words belonging to the rule applied to get the concept: [Arrival City], and “Chicago” is the value for this concept.

- **Number of words in the rule**

- **Number of words in the concept value**

- **Concept count in the 100-Best:** Each hypothesis in the 100-best are parsed and the percentage of times that a concept appears in the hypotheses are counted.

- **Concepts and value count in the 100-Best:** in this case we consider the number of times that a concept appears with the same value along the 100-best hypotheses.

These two features are useful when a confusable pair appears in the hypotheses. For example, consider two cities with high confusion between them: Austin and Boston. When we have the sentence: “I wanna go from Denver to Austin” in the 100-

best we observe many times how the word “Austin” is substituted by “Boston”. In this case the concept [Departing City] has high confidence because it appears in almost all the hypotheses but its value changes considerably, so this is a good measure of the concept value confidence. Obtaining large differences between these values means that the concept will probably be right but its value is confused.

The next two features are obtained considering a Concept Language Model. In our case we have trained a 3-gram LM considering the concepts obtained from the references in the training set. Similar to the case for word level, we have considered:

- **Language Model Back-Off Behavior:** back-off behavior of an N-gram language model along a 5-concept context.
- **Language Model Score:** the log-probability for each concept in a sequence as computed from a back-off language model along a 5-concept context.

For combining all features we have considered a MLP with the same characteristics described in the previous section.

5.1 Results.

Table 4 summarizes the correct detection rates of incorrect concepts at false rejection rates of 2.5% and 5.0%. We have run three experiments: considering the Average Word Confidence features (AWCs) separated, the rest of the features and all together. It can be seen that using only the AWC features we get slightly improved results than considering the remaining proposed features. For example, 47.1% of the incorrect concepts were detected at a false rejection rate of 5%. Using the rest features we reject 40.1% for the same false rejection. The best results are obtained combining all features. In this case we reject more than 50% of incorrect concepts at 5% false rejection. Considering that we have 16.5% of the incorrect concepts, these features reduce the classification error in 4.5%.

Concept Level	Correct Detection Rates		Classification Error Baseline 16.5%
	2.5% FR	5.0% FR	
AWCs	31.0%	47.1%	12.8%
Remaining Features	29.3%	40.1%	13.5%
ALL Features	35.9%	50.1%	12.0%

Table 4. Wrong concepts correct detection at a 2.5% and 5.0% false rejection rate (FR) rate. Minimum classification error and Baseline Error are also shown.

6. CONCLUSIONS

In this paper we present an analysis for confidence annotation for spoken dialogue systems at different levels: word, utterance and concept levels. All results are reported with respect to data collected from the CU Communicator during a seven month period (over 900 calls). At the word level we detect 53.2% of mis-recognized words at a 5% false rejection rate reducing the

classification error by 6.2%. From the experimental observations, features coming from the LM perform better than decoder features. The best results are obtained by combining all features together. We propose the use of confidence measures as heuristic to combine several hypotheses from different recognizers. We analyze two options for combining the hypotheses: the Flat List Confidence Rescoring (FLCR) method and the Word-Graph Confidence Rescoring (WGCR) method. The word-graph generation algorithm is described and results with this method are reported. Using the WGCR method we reach a 14% relative word error rate reduction. For the utterance level combining all features, 53.2% of out of domain utterances are detected at 5% false rejection rate. At the concept level a new set of features are proposed detecting more than 50% of incorrect concepts at a 5% false rejection rate. New features from the speech understanding component are proposed for confidence annotation at utterance and concept levels

7. FUTURE WORK

In a future work, we will analyze separately all features proposed from the speech understanding component, in order to obtain each contribution. For the utterance level we will consider the detection of, not only out of domain, but also misunderstanding utterances because of poor recognition. We also will consider the compilation of N-best lists from different coders into a single word graph to get a richer collection of hypotheses. Finally we will use the concept confidence to reduce the semantic errors that have more impact on the system's end-to-end performance.

8. REFERENCES

- [1] W. Ward, B. Pellom. “The CU Communicator System,” *Proc. IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, Keystone Colorado, 1999.
- [2] <http://fofoca.mitre.org>
- [3] B. Pellom, W. Ward, S. Pradhan, "The CU Communicator: an architecture for Dialogue Systems". In Proc. ICSLP, Beijing, China, 2000.
- [4] L. Chase, "Error-Responsive Feedback Mechanisms for Speech Recognizers", Ph.D. thesis, Carnegie Mellon University, Tech. Report CMU-RI-TR-97-18, April 1997
- [5] L. Chase, "Word and acoustic confidence annotation for large vocabulary speech recognition". In Proc. Eurospeech, Rhodes, Greece, 1997. pp 815-818.
- [6] S.O. Kamppari, T.J. Hazen, "Word and phone level acoustic confidence scoring". Proc. ICASSP, Istanbul, Turkey, 2000. pp III-1799, III-1802.
- [7] R. San-Segundo, B. Pellom, W. Ward., and JM. Pardo., "Confidence measures for dialogue management in the CU communicator system". Proc. ICASSP, Istanbul, Turkey, 2000. pp III-1237, III-1240.
- [8] T. Hazen, T. Burianek, J. Polifroni, S. Seneff, "Recognition Confidence Scoring for Speech Understanding Systems," Proc. of the ISCA ITRW ASR2000 Workshop on Automatic Speech Recognition: Challenges for the new Millenium, Paris France, September 2000, pp. 213-220.