

ENVIRONMENTAL SNIFFING: ROBUST DIGIT RECOGNITION FOR AN IN-VEHICLE ENVIRONMENT

Murat Akbacak and John H.L. Hansen

Robust Speech Processing Group, Center for Spoken Language Research
University of Colorado at Boulder, Boulder, CO, 80302, U.S.A

{murat, jhlh}@cslr.colorado.edu

Web: <http://cslr.colorado.edu>

ABSTRACT

In this paper, we propose to integrate an *Environmental Sniffing* [1] framework, into an in-vehicle hands-free digit recognition task. The framework of Environmental Sniffing is focused on detection, classification and tracking changing acoustic environments. Here, we extend the framework to detect and track acoustic environmental conditions in a noisy-speech audio stream. Knowledge extracted about the acoustic environmental conditions is used to determine which environment dependent acoustic model to use. Critical Performance Rate (CPR), previously considered in [1], is formulated and calculated for this task. The sniffing framework is compared to a ROVER solution for automatic speech recognition (ASR) using different noise conditioned recognizers in terms of Word Error Rate (WER) and CPU usage. Results show that the model matching scheme using the knowledge extracted from the audio stream by Environmental Sniffing does a better job than a ROVER solution both in accuracy and computation. A relative 11.1% WER improvement is achieved with a relative 75% reduction in CPU resources.

1. INTRODUCTION

Significant advances in ASR technology have been achieved in applications where the environmental noise condition is constant. Most recently, ASR research focus has shifted to real-world environments where changing environmental noise conditions represent significant challenges in maintaining ASR performance.

All efforts in the field of noisy speech recognition have been directed at reducing the mismatch between training and operating conditions such as speech enhancement, noise resistant features, re-training and multi-style training, and model adaptation. Each solution has both advantages and disadvantages [2].

Today, state of the art ASR systems use a parallel bank of recognizers in a ROVER paradigm [3] to take advantage of the methods mentioned above. This framework seeks to reduce word error rates for ASR by exploiting differences in the nature of the errors made by multiple speech recognizers which use different features in the feature extraction step, different noise compensation schemes in the enhancement step, or different model adaptation schemes. The disadvantage of this framework is the high computational power it requires. There are also many open questions: for example how important is the combination order of the system hypotheses?, which recognizers should be used?, how many systems

should we combine?, is it advantageous to preprocess or normalize the systems' outputs prior to combination? Most researchers take the approach of using more recognizers, and try to make each ASR engine different in some meaningful way (i.e., different features, trained in different noise, etc.) to leverage the potential differences in recognition errors.

In [1], we addressed the problem of changing acoustic environmental conditions in speech tasks by proposing a new framework entitled *Environmental Sniffing* to detect, classify and, track changing acoustic environmental conditions and extract knowledge about the environmental noise. The goal is to do smart tracking of environmental conditions and direct the ASR engine to use a solution specific to each environmental condition.

The organization of our paper is as follows. In Section 2, we specialize the general framework of sniffing environmental noise for an in-vehicle hands-free digit recognition task. In Section 3, algorithm formulation of environmental sniffing in a noisy-speech scenario is presented. Section 4 includes the formulation of the critical performance rate (CPR) of Environmental Sniffing for the digit car task. In Section 5, evaluations of the framework integrated into an in-vehicle ASR engine is presented. Section 6 discusses some further research issues for sniffing with conclusions given in Section 7.

2. SYSTEM ARCHITECTURE

Fig.1 shows a proposed robust ASR system for in-vehicle route information originally presented in [4].

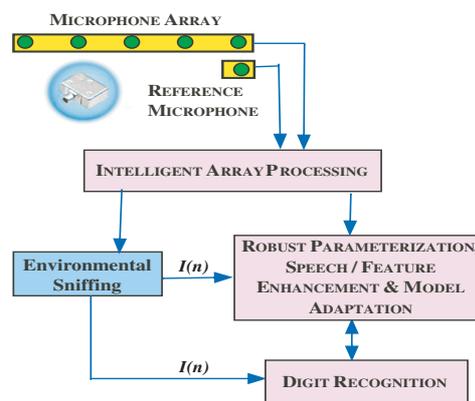


Fig. 1. An in-vehicle digit recognition system.

The motivation for selecting this environment is the huge diversity of acoustic environmental conditions and the need to maintain near real-time performance for route navigation dialogs.

In Fig. 1, we see that environmental sniffing plays a central role in determining the environment information which could be used to direct front-end array processing, parameterization, speech enhancement, model adaptation, or ASR model selection for effective speech recognition. Therefore, the environmental sniffer could be a passive system and simply provide information $I(n)$ to any prior or subsequent speech processing tasks. In contrast, the sniffer could instead take control and direct appropriate microphone array processing, feature selection/processing, and/or adjust model adaptation depending on the environmental knowledge and confidence. For the purpose of this paper, the Environmental Sniffing framework will be employed for ASR model selection.

3. ENVIRONMENTAL SNIFFING FOR NOISY SPEECH

In [1], we focused on extracting knowledge concerning the acoustic environmental noise using a noise-only audio database containing 8 noise conditions in a car environment. Here, we present a broad class monophone recognition based system for sniffing noisy-speech data, as shown in Fig. 2. In addition to 8 noise conditions, the acoustic condition set contains also the clean condition-CL as shown in Table 1.

Acoustic Condition Set		
1	N1	Idle noise consisting of the engine running with no movement and windows closed
2	N2	City driving without traffic and windows closed
3	N3	City driving with traffic and windows closed
4	N4	Highway driving with windows closed
5	N5	Highway driving with windows 2 inches open
6	N6	Highway driving with windows half-way down
7	N7	Windows 2 inches open in city traffic
8	NX	Others
9	CL	Clean (i.e., noise-free)

Table 1. In-vehicle acoustic conditions considered.

After defining a set of broad phone classes (e.g., STP- stop, FRC- fricative, NSL- nasal, VWL- vowel, SIL- silence, etc.), an HMM is trained for each (*broad phone class, acoustic condition*) pair. As an example, an HMM for the pair (FRC,N1) is trained from a clean database of fricatives degraded by acoustic condition N1. These acoustic models are used during the broad class monophone recognition.

As Fig. 2 shows, the incoming audio stream is first segmented into acoustically homogeneous speech blocks using our T^2 -BIC [5] segmentation scheme with a low false alarm penalty (i.e. false alarms are tolerable to ensure we capture all potential marks, both true and false). For each segment, a lattice is generated in an FST (Finite State Transducer) format via phoneme recognition. During decision smoothing, the resulting phone-lattice of each segment is combined with an FST representing the noise language model. The costs of noise transitions in the FST representing the noise language model is inversely proportional with the transition probabilities presented in [1].

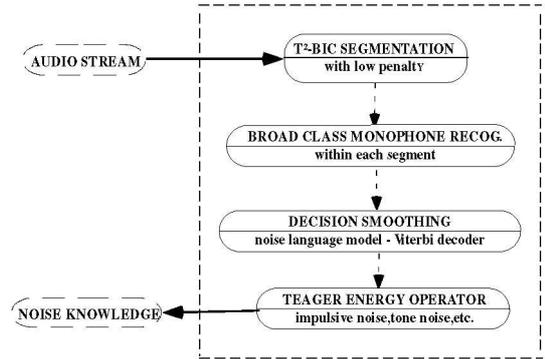


Fig. 2. Environmental Noise Sniffing.

4. CRITICAL PERFORMANCE RATE

In [1], we defined a critical performance rate (CPR) in a general sense. We now specialize the formulation of CPR to a specific case where Environmental Sniffing framework is used for model selection within an ASR system. The Environmental Sniffing framework determines the initial acoustic model to be used according to the environmental knowledge it extracts. The knowledge in this context, will consist of the acoustic condition types with time tags.

Let us denote the error matrix for noise classification as ϵ :

$$\epsilon = \begin{bmatrix} e_{11} & e_{12} & \dots & e_{1N} \\ e_{21} & e_{22} & \dots & e_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ e_{N1} & e_{N2} & \dots & e_{NN} \end{bmatrix}. \quad (1)$$

For $i = j$, $1 \leq i, j, \leq N$, e_{ij} is zero, and for $i \neq j$, e_{ij} is the classification error rate (in a range 0-1) for the error type where the i^{th} noise class is classified as the j^{th} noise class.

Assume that there are N initial acoustic models¹ to be used during recognition, each corresponding to an environmental condition. These models can be trained by simply re-training HMMs for N different acoustic conditions. Assume that there is enough diversity among noise conditions so that for a noise type during decoding, using the matched acoustic model as an initial model during model adaptation yields the lowest WER. Let us define a matrix W as follows:

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1N} \\ w_{21} & w_{22} & \dots & w_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ w_{N1} & w_{N2} & \dots & w_{NN} \end{bmatrix} \quad (2)$$

where w_{ij} represents the WER value for the case where test tokens are from the i^{th} noise class, but the j^{th} acoustic model which is trained from the j^{th} acoustic condition is used as an initial model. Using the matrix W , we can assign a cost value for each error type so that each error rate e_{ij} can be weighted by the normalized cost values to calculate the Critical Performance Rate (CPR) of the Environmental Sniffing framework. For the error type where

¹If there are M ($M \leq N$) initial models, the \mathbf{W} matrix will still be $N \times N$, since some noise classes will use the same acoustic model, and the cost of errors among these noise classes will be zero.

the i^{th} noise class is classified as the j^{th} noise class, the cost is $\Delta w_{ij-ii} = w_{ij} - w_{ii}$, which is the performance deviation of the ASR engine by using the j^{th} acoustic model during decoding instead of using the correct i^{th} acoustic model.

Since some noise conditions occur more frequently than others, each noise condition will have an *a priori* probability denoted as follows:

$$\vec{\mathbf{a}} = [a_1 \quad a_2 \quad \dots \quad a_N] \quad (3)$$

Now, we can formulate the Critical Performance Rate as:

$$\begin{aligned} CPR &= 1 - \sum_{i=1}^N a_i \sum_{j=1}^N \frac{\Delta w_{ij-ii}}{n_{ij}} e_{ij} \\ &= 1 - \sum_{i=1}^N a_i \sum_{j=1}^N \frac{w_{ij} - w_{ii}}{\frac{\sum_{k=1, k \neq i}^N w_{ik} - (N-1)w_{ii}}{N-1}} e_{ij} \\ &= 1 - \sum_{i=1}^N a_i \sum_{j=1}^N \frac{w_{ij} - w_{ii}}{\frac{\sum_{k=1}^N w_{ik} - Nw_{ii}}{N-1}} e_{ij} \\ &= 1 - \sum_{i=1}^N a_i \sum_{j=1}^N C_{ij} e_{ij} \end{aligned} \quad (4)$$

where n_{ij} is the normalization term and C_{ij} is the normalized cost value for the error type where the i^{th} noise class is classified as the j^{th} noise class.

In matrix form, Eq. 4 becomes:

$$CPR = 1 - \text{diag}\{\mathbf{C} \cdot \epsilon^T\} \cdot \vec{\mathbf{a}}^T \quad (5)$$

where \mathbf{C} is the normalized cost matrix having entries C_{ij} .

If all noise conditions have equal *a priori* probabilities $1/N$, and all error types have equal costs (e.g., each error type has the same impact on the subsequent system's performance) then we obtain

$$CPR = 1 - \sum_{i=1}^N \frac{1}{N} \sum_{j=1}^N e_{ij} \quad (6)$$

The goal, in terms of performance, is to optimize the critical performance rate rather than optimizing the environmental noise classification performance rate, since it is more important to detect and classify noise conditions that have a more significant impact on ASR performance.

We can use Eq. 4 by setting $n_{ij} = 1$ (without normalizing the costs) to calculate $1 - CPR$, which will actually be the expected performance deviation from the highest achievable performance when we use Environmental Sniffing. In our example, using the matched acoustic model for an environmental condition will achieve the lowest WER. Using the expected performance deviation, we can estimate the performance of model matching employing the Environmental Sniffer.

5. EVALUATIONS

In our evaluations, we degraded the TI-DIGIT database at random SNR values ranging from -5 dB to +5 dB (i.e., -5,-3,-1,+1,+3,+5 dB SNR) with 8 different in-vehicle noise conditions using the noise database from [1]. A 2.5-hour noise data set was used to degrade the training set of 4000 utterances, and the 0.5 hour set was used to degrade the test set of 500 utterances (i.e., open noise degrading condition). Each digit utterance was degraded with only

one acoustic noise condition. Next, an HMM was trained for each (*broad phone class, acoustic condition*) pair. The phoneme classes for digits were mapped to 7 broad classes (including SIL- silence). Since we have 9 acoustic conditions (including CL- clean) in the acoustic condition set, at the end we had $7 \times 9 = 63$ HMMs. Each (*broad phone class, acoustic condition*) was listed in the lexicon during decoding. A total of 7 silence models were also used as filler models.

Each digit utterance was decoded into a sequence of (*broad phone class, acoustic condition*) pairs. At each leaf of the lattice, the (*broad phone class, acoustic condition*) pair was mapped to the corresponding acoustic condition (e.g., STP-N1 was mapped to N1). The resulting lattice in FST format is combined with the lattice representing the noise language model to find the most likely noise sequence.

For acoustic model training and decoding, we used CSLR's Large Vocabulary Continuous Speech Recognizer SONIC [6]. AT&T's FSM Toolkit [7] was used to combine the phone-lattice and the noise language model.

5.1. Sniffing Results:

Using the sniffing framework presented in Sec. 3, each utterance was assigned to an acoustic condition. Using the fact that there was only one acoustic condition within each utterance, the Environmental Sniffing framework did not allow noise transitions within an utterance. A noise classification rate of 82% was obtained. From this, a 9×9 error matrix ϵ was generated for use in digit recognition.

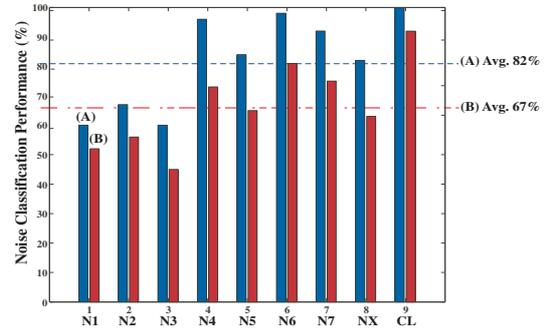


Fig. 3. Sniffing performance for each noise type (A) with(left bar in each pair)/(B) without(right bar in each pair) the prior knowledge that there was one environmental condition present in each utterance.

5.2. Digit Recognition Results:

5.2.1. Development Phase:

Environmental condition specific acoustic models were trained and used during recognition tests. Matrix \mathbf{W} was generated by testing different acoustic conditions using different acoustic models. By using the \mathbf{W} matrix, we calculated the normalized cost matrix \mathbf{C} using Eq. 4. Using Eq. 5, with the *a priori* noise probabilities,

$$\vec{\mathbf{a}} = [0.05 \quad 0.15 \quad 0.15 \quad 0.15 \quad 0.15 \quad 0.15 \quad 0.15 \quad 0.05],$$

CPR was calculated as 92.1%, and the expected performance deviation was found to be 0.22% when the Environmental Sniffer uses the knowledge that only one acoustic condition is present within each utterance.

5.2.2. Test Phase:

Having established the environmental sniffer, and normalized cost matrix C for directing ASR model selection, we now turn to ASR system evaluation. We tested and compared the following 3 system configurations:

- S1:** Model matching was done using *a priori* knowledge of the acoustic noise condition.
- S2:** Model matching was done based on the environmental acoustic knowledge extracted from Environmental Sniffing.
- S3:** All acoustic condition dependent models were used in a parallel multi-recognizer structure (e.g. ROVER) without using any noise knowledge and the recognizer hypothesis with the highest path score was selected.

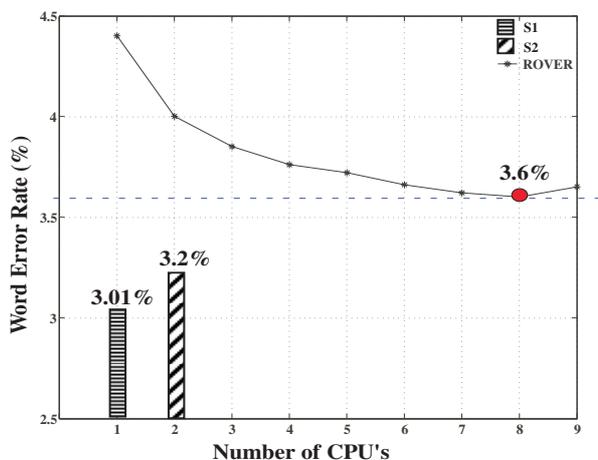


Fig. 4. Comparison of system configurations S1, S2, S3.

As shown in Fig. 4, system S1 achieved the lowest WER (i.e., 3.01%) since the models were matched perfectly to the acoustic condition during decoding. From the development phase, we know that the expected performance deviation was 0.22 for a model matching scheme employing Environmental Sniffing, which means that we can expect a WER value of $3.01 + 0.22 = 3.23\%$ for S2. Experimentally, the WER for S2 was 3.2% using 2 CPU's (1 CPU for digit recognition, 1 CPU for sniffing acoustic conditions), which was close to the expected value of 3.23% (Note: in Fig. 4, we plot system S2 2 CPU even though only 1 ASR engine was used). S3 achieved a WER of 3.6% by using 8 CPU's. When we compare S2 and S3, we see that a relative 11.1% WER improvement was achieved, while requiring a relative 75% **reduction** in CPU resources. These results confirm the advantage of using Environmental Sniffing over a ROVER paradigm.

6. DISCUSSION

There are two critical points to consider when integrating Environmental Sniffing into a speech task. First, and the most important, is to set up a configuration such as S1 where prior noise knowledge can be fully used to yield the lowest WER. This will require understanding of the sources of errors and finding specific solutions assuming that there is prior acoustic knowledge. For example, knowing which speech enhancement scheme or model adaptation scheme is best for a specific acoustic condition is required.

Secondly, a reliable cost matrix should be provided to the Environmental Sniffing so the subsequent speech task can calculate the expected performance in making an informed adjustment in the trade-off between performance and computation. For our experiments, we considered evaluation results for Environmental Sniffing where it is employed to find the *highest* possible acoustic condition so that correct acoustic condition dependent model could be used. This is most appropriate for the goal of determining a single solution for the speech task problem at hand. If the expected performance for the system employing Environmental Sniffing is lower than the performance of a ROVER system, it may be useful to find the n most probable acoustic condition types among N acoustic conditions. In the worst case, the acoustic condition knowledge extracted from Environmental Sniffing could be ignored and the system will reduce to the traditional ROVER solution.

Finally, in an ASR task, in addition to determining the number of systems that should be combined, even a ROVER paradigm could take advantage of Environmental Sniffing to address open questions such as the combination order of the ASR system hypotheses, engage or disable preprocessing or normalization of system outputs prior to combination, etc. The goal therefore has been to emphasize that direct estimation of environmental conditions should provide important information to tailor a more effective solution to robust speech recognition systems.

7. CONCLUSION

In this study, we have extended our previous proposed *Environmental Sniffing* [1] framework and integrated it into an in-vehicle hands-free digit recognition system. The critical performance rate (CPR) was formulated for this task. The sniffing framework was compared to a ROVER solution in terms of WER and CPU usage in a model matching task where environmental condition dependent models were used during decoding. In our experiments, the presented framework consistently outperformed the original ROVER solution by 11.1% in WER, while requiring 75% less CPU resources.

8. REFERENCES

- [1] M. Akbacak, J.H.L. Hansen, "Environmental Sniffing: Noise Knowledge Estimation for Robust Speech Systems," *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, April 2003.
- [2] Y. Gong, "Speech Recognition in Noisy Environments: A Survey," *Speech Communication*, v16-3, 261-91, 1995.
- [3] J. G. Fiscus, "A Post Processing System to yield reduced error rates: Recognizer Output Voting Error Reduction (ROVER)," *IEEE Workshop on Automatic Speech Recognition and Understanding*, 347-54, 1997.
- [4] J.H.L. Hansen, et al., "CU-Move: Analysis & Corpus Development for Interactive In-Vehicle Speech Systems," *Eurospeech*, v3, 2023-6, 2001.
- [5] B. Zhou, J.H.L. Hansen, "Unsupervised Audio Stream Segmentation and Clustering via the Bayesian Information Criterion," *Proc. ICSLP*, v3, 714-7, 2000.
- [6] B. L. Pellom, "Sonic: The University of Colorado Continuous Speech Recognizer," *Technical Report #TR-CSLR-2001-01*, 2001
- [7] AT&T FSM Library, <http://www.research.att.com/sw/tools/fsm/description.html>