

# SPEEDING UP THE SCORE COMPUTATION OF HMM SPEECH RECOGNIZERS WITH THE BUCKET VORONOI INTERSECTION ALGORITHM

J. Fritsch, I. Rogina, T. Sloboda, A. Waibel  
{fritsch,rogina,sloboda,waibel}@ira.uka.de

Interactive Systems Laboratories  
University of Karlsruhe — Germany  
Carnegie Mellon University — USA

## ABSTRACT

With increasing sizes of speech databases, speech recognizers with huge parameter spaces have become trainable. However, the time and memory requirements for high accuracy realtime speaker-independent continuous speech recognition will probably not be met by the available hardware for a reasonable price for the next few years. This paper describes the application of the Bucket Voronoi Intersection algorithm to the JANUS-2 speech recognizer, which reduces the time for the computation of HMM emission probabilities with large Gaussian mixtures by 50% to 80%.

## 1. INTRODUCTION

Although the computation of Gaussians is only a part (for very large vocabularies, even a small part) of the overall run time, speeding it up does reduce the reaction time of the recognizer, and especially the time for training significantly. When computing the log probability of a Gaussian mixture, many speech recognizers do not use all Gaussians but only the top  $n$ . We have found that in our system using only the one Gaussian with the highest probability is almost as good as using the sum of more Gaussians. We have also found that using the Euclidean distance instead of the Mahalanobis distance for finding that most probable Gaussian does not decrease recognition accuracy too much. This reduces the computation of an HMM emission probability to a two part process: First, find the centroid that has the smallest Euclidean distance to the current speech sample, and second, compute the value of the Gaussian (multiplied with its mixture weight) for that centroid. So instead of computing  $n$  Gaussians, where  $n$  is the size of the mixture, we only have to compute one Gaussian plus we have to run an algorithm for finding the closest centroid. For this we use the Bucket Voronoi Intersection (BVI) algorithm [1]. It was introduced for high speed vector quantization of low-dimensional vectors. However, we have found that it is still good enough for 16-dimensional speech vectors. In this paper we describe experiments in which we have investigated the effect of the BVI-algorithm on the run-time behavior and the recognition accuracy of the JANUS-2 speech recognizer [2, 3].

## 2. THE BUCKET VORONOI INTERSECTION ALGORITHM

For a detailed discourse on the Bucket Voronoi Intersection (BVI) algorithm see [1].

All points in the feature space having the same nearest-neighbor codebook vector define a Voronoi region. The set of all Voronoi regions constitutes a disjoint partitioning of the feature space. The aim of the algorithm is to approximate this partitioning with a top-down tree search.

The principle behind it is a binary tree. Each node of the tree represents a hyperplane in the feature space. When classifying a sample vector, the tree is descended from the root down to a leaf. At every node, a decision is made to descend into the left or the right successor node, depending on the sample vector being on the left or on the right side of the current hyperplane. So every step down the tree reduces the size of the search space.

When the tree descending algorithm has finally reached a leaf node, there will be only a few codebook vectors left whose Voronoi region is intersecting with the remaining search space, which is called a bucket. The set of all buckets constitutes a disjoint partitioning of the feature space. Depending on how deep we descend the tree, we get different buckets and a different partitioning of the feature space. In higher levels of the tree we get larger buckets, which contain more Voronoi regions.

Although the bucket sizes decrease monotonically with increasing tree depth, there is no guarantee to reach the optimal case of having only one codebook vector per bucket. For that reason, there is a tradeoff between speed up and memory requirements of the tree.

The time for traversing the tree is not the critical factor. Let  $d$  be the depth of the tree,  $b$  the average bucket size and  $n$  the codebook size, then we will have to compute  $d$  hyperplane comparisons, plus  $b$  Gaussians instead of  $n$  Gaussians. Since the BVI-algorithm only uses hyperplanes of the form  $x_i = t$ , deciding on which side of the hyperplane a vector  $\mathbf{y}$  is located takes only one simple floating point comparison  $y_i < t$ . A full binary tree of depth  $d$  has  $2^d$  leaves (buckets), so the memory requirements for storing the tree grow exponentially. Since we are usually using feature spaces that are at least 16-dimensional, the limit for the depth of the trees will be determined rather by the amount of available memory, than by the run-time requirements.

## 3. COMPUTING THE BVI-TREE

Since it is extremely expensive to compute the real boundaries of a high dimensional Voronoi region, we approximate the Voronoi region with a cuboid whose edges are parallel to the coordinates. These approximate regions generally overlap each other. The boundaries of the cuboids are determined by encoding a sufficiently large set of training vectors. (see Fig. 1).

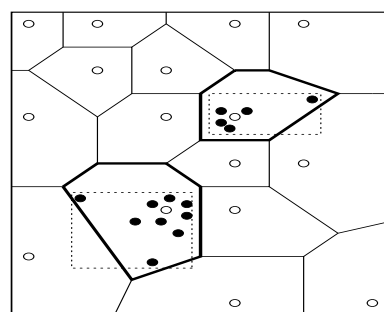


Fig. 1: approximated Voronoi regions

A cuboid-approximated Voronoi region is defined entirely on one side of a hyperplane if all the training vectors that fall into the region are on the same side of the plane. With this approach we get a very simple decision rule, but we introduce a possible classification error (Fig. 2).

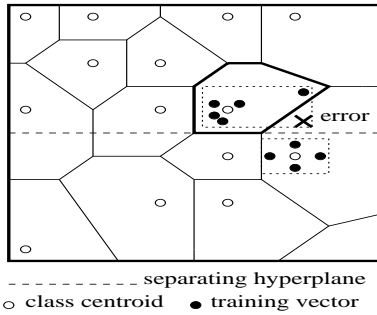


Fig. 2: classification error in approximated Voronoi regions

The error rate can be reduced by increasing the number of training vectors. The more vectors we use for training the more it is likely that the approximate cuboid of a Voronoi region will contain the entire region. Fig. 3 shows the average classification error rate, depending on the number of training vectors and the depth of the BVI-trees.

In our experiments, we have found that a low classification error rate for the nearest neighbor is not important for a good speech recognition accuracy (see Fig. 4)

The objective of a good BVI-tree is to have as few Voronoi regions in every bucket as possible. The average size of a bucket decreases with the depth of the tree, while the memory requirements and error rate grow exponentially, limiting the tree size. We have conducted experiments with trees up to a depth of 12. Fig 5 shows the average bucket size depending on the depth of the BVI-tree.

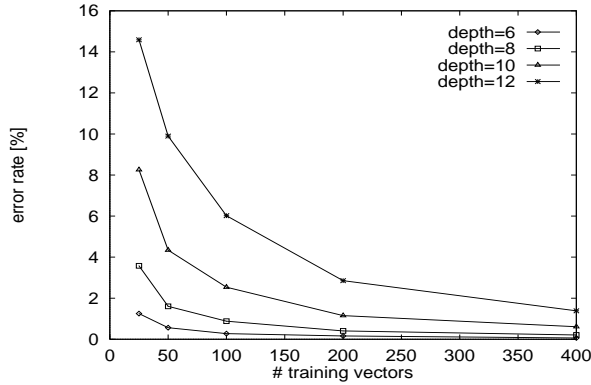


Fig. 3: average classification error rate

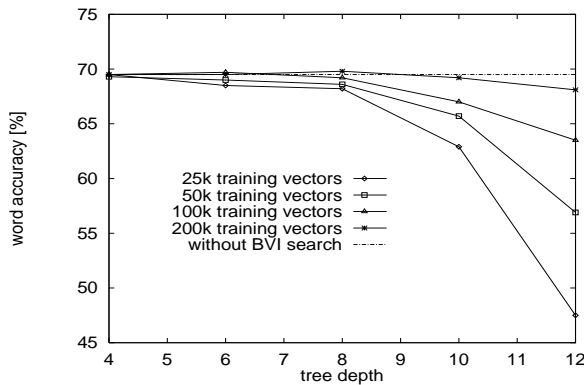


Fig. 4: recognition accuracy using BVI-search

#### 4. RUNTIME BEHAVIOR

The speedup in the HMM-emission probability computation can be approximated by the average mixture size divided by the average bucket size. Of course, the relative speedup for

the entire system is smaller. Fig. 6 shows the speedups for training and testing sessions with the JANUS-2 speech recognizer.

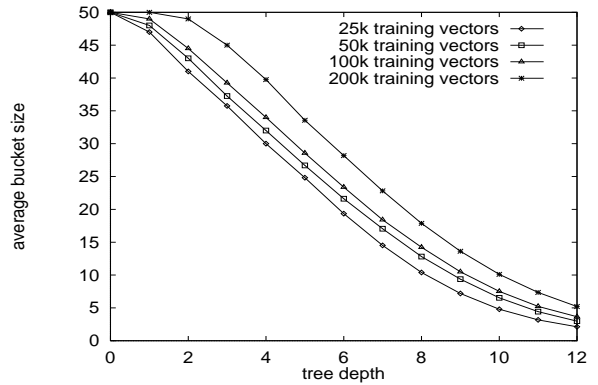


Fig. 5: average bucket size depending on tree depth

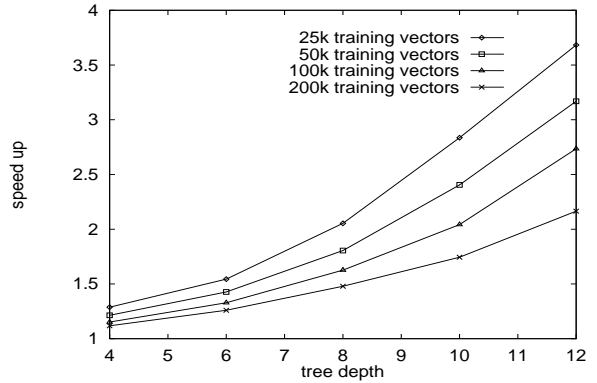


Fig. 6: speedup of BVI-score computation

#### 5. RECOGNITION ACCURACY

We have found that the recognition accuracy of the speech recognizer does not suffer from the possible classification errors of the BVI-algorithm. Fig. 4 shows the word accuracy on the German Spontaneous Scheduling Task (GSST) [3, 2] for different amounts of training data for the BVI-algorithm.

#### 6. CONCLUSION

We present first results of our ongoing research on speeding up the score computation with the BVI-algorithm. Although the algorithm was developed for data compression applications, we successfully integrated this fast vector quantization method into an HMM speech recognizer.

#### REFERENCES

- [1] Ramasubramanian, V.; Paliwal, K. K.: *Fast K-dimensional Tree Algorithms for Nearest Neighbor Search with Application to Vector Quantization Encoding*, IEEE Transactions on Signal Processing, Vol. 40, No. 3, March 1992.
- [2] M. Woszczyna, N. Aoki-Waibel, F.D. Buø, N. Coccaro, K. Horiguchi, T. Kemp, A. Lavie, A. McNair, T. Polzin, I. Rogina, C.P. Rose, T. Schultz, B. Suhm, M. Tomita, A. Waibel: *JANUS 93: Towards Spontaneous Speech Translation*, Proceedings of the ICASSP 1994, Adelaide, volume 1, pp 345-348.
- [3] M. Woszczyna, N. Coccaro, A. Eisele, A. Lavie, A. McNair, T. Polzin, I. Rogina, C.P. Rose, T. Sloboda, M. Tomita, J. Tsutsumi, N. Aoki-Waibel, A. Waibel, W. Ward: *Recent Advances in Janus, a Speech to Speech Translation System*, Proceedings of the EUROSPEECH, Berlin, 1993.