

# SUBSPACE HIGH-DENSITY DISCRETE HIDDEN MARKOV MODEL FOR AUTOMATIC SPEECH RECOGNITION

Guoli Ye, Brian Mak

Department of Computer Science and Engineering  
The Hong Kong University of Science and Technology

{yeguoli, mak}@cse.ust.hk

## ABSTRACT

Compared with continuous density hidden Markov model (CDHMM), discrete HMM (DHMM) has inherent attractive properties: it takes only  $O(1)$  time to get the state output probability, and the discrete features, compared with cepstral coefficients, could be encoded in fewer bits, lowering the bandwidth requirement in distributed speech recognition architecture. Unfortunately, the recognition performance of conventional DHMM is significantly worse than that of CDHMM due to the large quantization error and the use of multiple independent streams. One way to reduce the quantization error and to improve the recognition accuracy, is to use a very large codebook. However, the training data is usually not sufficient to robustly train such high density DHMM (HDDHMM). In this paper, we investigate a subspace approach together with subvector quantization to solve the training problem of HDDHMM. The resulting model is called subspace HDDHMM (SHDDHMM). On both Resource Management and Wall Street Journal 5K-vocabulary task, when compared with its CDHMM counterpart, SHDDHMM shows comparable performance in recognition accuracy, with faster decoding speed and lower bandwidth requirement.

**Index Terms**— subspace modeling, subvector quantization, high-density discrete HMM, semi-continuous HMM

## 1. INTRODUCTION

For many current ASR systems, CDHMM with Gaussian mixture model (GMM) as state output distribution, is commonly adopted, due to its simple parametric form and good generalization ability. On the other hand, DHMM with discrete density as state output distribution, which was used in the early stage of speech recognition history, is almost abandoned nowadays, mainly because its recognition accuracy is not as good as that of CDHMM.

The unsatisfactory performance of traditional DHMM mainly arises from two reasons:

- Large quantization error: unlike GMM which has a smooth probability density function, discrete density imposes hard partitions in the feature space, with each partition represented by a codeword. All the different feature vectors falling into the same partition are quantized into the same codeword, which inevitably causes quantization error. Due to the limit of memory, and the amount of training data, traditional DHMM usually uses a codebook size of 256 to 1024. Such a small codebook size usually results in large quantization error, which badly hurts the recognition performance.

This work was supported by the Research Grants Council of the Hong Kong SAR under the grant number HKUST617008.

- Incorrect stream independent assumption: due to the use of large feature dimension (e.g., the common 39-dimensional MFCC vector) in ASR, multiple-stream DHMM is employed to allow manageable codebook sizes. The inability to model correlation among streams also limits the performance of DHMM.

With the advance of computer technology, we would like to revisit DHMM to see if some of its limitations could be overcome nowadays. Compared with CDHMM, DHMM is attractive to us for the following reasons:

- Fast decoding: in CDHMM, it is found that the Gaussian likelihood computation may take about 30-70% of the total recognition time [1]. What is worse, the computation time grows linearly with the number of Gaussian components. Large and complex systems usually use a large mixture of Gaussians per state, thus resulting in a heavy computation burden in decoding. In DHMM with each state storing a discrete density table, to get the state output probability only requires a  $O(1)$  time table lookup, which is independent of the table size<sup>1</sup>.
- Bandwidth saving in distributed speech recognition (DSR) systems: DSR systems become more common with the emergence of cloud computing technology. In these DSR systems, ASR components are distributed between clients and servers. Speech features may be extracted locally at a client, and then transmitted to a remote server, where recognition is carried out. The discrete features used by DHMM could be encoded in much fewer bits than MFCC features used in CDHMM, lowering the bandwidth requirement.

In recent years, there is the resurrection of using discrete systems to model speech [2, 3]. Unlike the unsupervised vector quantization methods (e.g., k-means clustering) commonly used in traditional DHMM, Droppo [2] proposed partitioning the feature space by an acoustic decision tree using MMI-based purity function at each tree node. By making use of the class label knowledge, which is not used in k-means clustering, the resulting discrete system showed comparable performance as that of GMM, but only on a simple frame classification task.

Teunen and Akamine [3] introduced mixture models and soft decisions to alleviate the quantization error of discrete density. Mixture

<sup>1</sup>The time for finding the codeword for an input vector in DHMM depends on the vector quantization (VQ) method. e.g., if the simplest LBG method is used, finding the codeword takes  $O(N)$  time for a codebook size of  $N$ ; on the other hand, if tree VQ is used, it takes  $O(\log N)$  time. Furthermore, for an input vector, the codeword only needs to be found once, and then used by all the active states.

models make use of the smoothing property of ensemble method, and soft decisions make soft partitions rather than hard partitions. However, these techniques incur additional computational cost. Furthermore, compared with CDHMM, there was still a 43% relative increase in word error rate (WER) of the proposed model on WSJ 5K task [3].

In [4], we proposed the *high-density discrete hidden Markov model* (HDDHMM) to reduce the quantization error of discrete density by simply increasing the codebook size. In theory, if the codebook size is infinitely large, then the quantization error will become infinitely small. In practice, limited by computational resources, we need to address the following three problems for HDDHMM:

- requirement of large memory to store the model
- how to find the codeword for an acoustic vector fast?
- lack of training data to train such a high-density model

The first two problems are not as serious as one may think. With the advance in semiconductor memory and its falling price, large memory requirement could be easily satisfied nowadays. To find the codeword for an acoustic vector fast, tree-structured, or subvector-quantized (SVQ) codebooks could be used instead of linear codebook. Unlike linear codebook in which the search time for a codeword grows linearly with the codebook size  $N$ , the tree-structured or SVQ codebooks only take  $O(\log N)$  time. In this paper, SVQ codebooks are used.

The lack of training data turns out to be a real challenge for HDDHMM with a codebook consisting of *tens of thousands* of codewords. For such a large codebook, even with today's large inventory of speech corpora, it is still hard to get sufficient training data for it. In [4], HDDHMM was directly converted from CDHMM without any retraining. In this paper, we investigate a subspace approach in which each state output discrete density is treated as a vector, and it is required to lie in the subspace spanned by a global pool of bases. The bases are further constrained with the use of SVQ in order to reduce the number of model parameters. Consequently, the number of model parameters in the resulting HDDHMM is greatly reduced and re-estimation of the model parameters becomes possible. We call our new HDDHMM "*subspace high-density discrete hidden Markov model*" (SHDDHMM).

## 2. SUBSPACE HIGH-DENSITY DISCRETE HMM

The overview structure of SHDDHMM is shown in Fig. 1. The details are explained in the following sections.

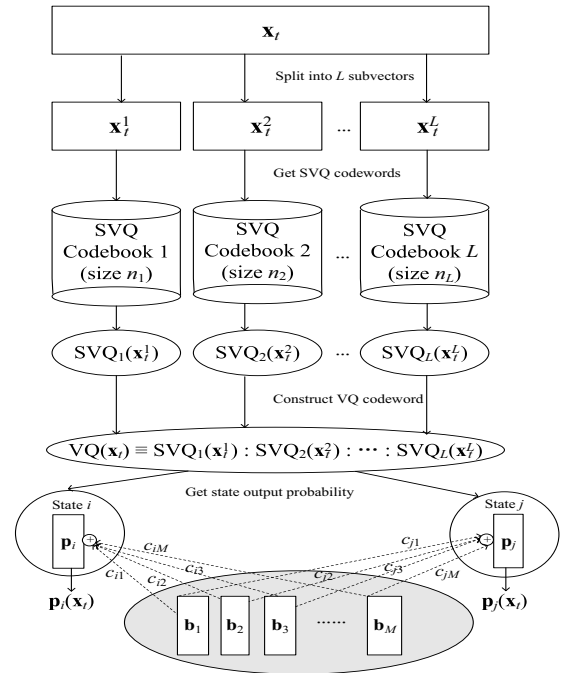
### 2.1. Subvector-quantized (SVQ) Codebooks

This section gives an introduction of SVQ codebooks, which are used by SHDDHMM. To construct the SVQ codebooks, each  $D$ -dimensional acoustic vector  $\mathbf{x}_t$  is split into  $L$  subvectors<sup>2</sup>,  $\mathbf{x}_t = [\mathbf{x}_t^1, \mathbf{x}_t^2, \dots, \mathbf{x}_t^L]$ . Subvectors of each split are vector-quantized (by k-means) to create an SVQ codebook for the split. Let  $n_i$  denote the SVQ codebook size of the  $i$ th split.

Given a new acoustic vector, it will first be split into  $L$  subvectors. The subvector in the  $i$ th split is then compared with the corresponding SVQ codebook to find its SVQ codeword in  $O(n_i)$  time. After getting all the  $L$  SVQ codewords of an input vector, its full-space VQ codeword is constructed as the products of SVQ codewords, one from each of the  $L$  splits as:

$$\text{VQ}(\mathbf{x}_t) \equiv \text{SVQ}_1(\mathbf{x}_t^1) : \text{SVQ}_2(\mathbf{x}_t^2) : \dots : \text{SVQ}_L(\mathbf{x}_t^L),$$

<sup>2</sup>When  $L = D$ , the SVQ codebooks become SQ codebooks used in [4].



**Fig. 1.** SHDDHMM architecture overview. The shaded ellipse represents the global pool of bases, spanning the subspace. Each state output discrete density table (i.e.,  $\mathbf{p}_i$ , and  $\mathbf{p}_j$ ), lies in the subspace. The state-dependent weights (the dotted lines) and the global pool of bases (the shaded ellipse) are temporary parameters, which only exist during model training; the final model only stores  $\mathbf{p}_i$  and  $\mathbf{p}_j$ .

where  $\text{VQ}(\cdot)$  and  $\text{SVQ}_i(\cdot)$  represent the VQ codeword given the full-space acoustic vector and the SVQ codeword in the  $i$ th split given the acoustic subvector in that split respectively. Notice that although SVQ is employed, it is only used to efficiently index a VQ codeword in the original full-space through the combinatorial effect of per-split SVQ codewords.

Since construction of full-space codeword from SVQ codewords would take no time, the time to get the full-space codeword for a given input is determined by the cost to find the  $L$  SVQ codewords, which is  $(\sum_{i=1}^L n_i)$ . Thus, by exploiting SVQ codebooks, finding codeword for an input vector is significantly quicker than that of the linear VQ codebook with the same codebook size, which takes  $O(\prod_{i=1}^L n_i)$  time. Finally, the constructed full-space codeword is used as an index to get the state output probability by a table lookup, which takes only  $O(1)$  time.

### 2.2. SHDDHMM with Constrained Bases

Like DHMM, SHDDHMM usually has multiple streams. For simplicity, we only discuss SHDDHMM with one stream in this section. It should be easy to extend it to multiple streams.

#### 2.2.1. Basic Structure

The basic structure of SHDDHMM could be simply written as:

$$\mathbf{p}_j(\text{VQ}(\mathbf{x}_t)) = \sum_{m=1}^M c_{jm} \mathbf{b}_m(\text{VQ}(\mathbf{x}_t)) \quad (1)$$

where  $j$  is an HMM state;  $\mathbf{p}_j$  is the output discrete density table of state  $j$  with  $\mathbf{p}_j(\cdot)$  being the discrete density function;  $\mathbf{b}_1, \dots, \mathbf{b}_M$  are a pool of  $M$  discrete density tables shared globally by all the states;  $c_{j1}, \dots, c_{jM}$  are the state-dependent weights, satisfying the constraint that  $\sum_{m=1}^M c_{jm} = 1$ . By viewing a discrete density table as a vector, the vector  $\mathbf{p}_j \in \mathbb{R}^N$  thus lies in the subspace spanned by  $M$  bases  $\mathbf{b}_1, \dots, \mathbf{b}_M$ <sup>3</sup>. We call it “subspace” as the number of bases  $M$  in the model is usually significantly smaller than the dimension  $N$  of vector  $\mathbf{p}_j$ , which is the discrete density table size or codebook size.

Since the codebook size  $N$  (i.e., the length of each basis vector) in SHDDHMM is usual very large (tens of thousands), even if we use the subspace representation, there are still too many parameters in the system. Thus, similar as in [5], we made the assumption that: the SVQ codewords of different subvectors are conditionally independent given the basis. This assumption imposes a constraint on each basis as:

$$\mathbf{b}_m(\text{VQ}(\mathbf{x}_t)) = \prod_{i=1}^L \mathbf{b}_m^i(\text{SVQ}_i(\mathbf{x}_t^i)) \quad (2)$$

where  $\mathbf{b}_m^1 \in \mathbb{R}^{n_1}, \dots, \mathbf{b}_m^L \in \mathbb{R}^{n_L}$  are the SVQ discrete density tables in split 1 to  $L$  respectively, which are used to construct the original basis  $\mathbf{b}_m \in \mathbb{R}^N$ . Thus, the “real” parameters in SHDDHMM with constrained bases are  $c_{jm}$ ’s and  $\mathbf{b}_m^i$ ’s, which have significantly fewer number of parameters compared with the original system parameters  $\mathbf{p}_j$ ’s.

Let us use an example to illustrate the reduction of system parameters in the proposed model. Let  $I = 3132$  be the number of states in our tied-state system,  $M = 1024$  be the number of bases,  $L = 2$  so that a feature vector is split into 2 subvectors, and the SVQ codebook sizes for them are  $n_1 = 256 = n_2 = 256$  respectively. Thus, the equivalent full-space codebook size,  $N \equiv n_1 \times n_2 = 65536$ . In the original system, there are  $I$  states, each with a discrete table of size  $N$ , giving a total of  $NI = 205,258,752$  parameters! On the other hand, in the proposed model, there will be only  $MI = 3,207,168$  parameters for the weights  $c_{jm}$ ’s, and  $M(n_1 + n_2) = 524,288$  parameters for the constrained bases  $\mathbf{b}_m^i$ ’s, resulting in a total of 3,731,456 model parameters, which is 55-fold reduction compared with the original system.

Notice that the weight parameters  $c_{jm}$ ’s and the basis parameters  $\mathbf{b}_m^i$ ’s only exist during the training stage, which make training with limited amount of data feasible. After training, we use these parameters to compute the state output discrete densities  $\mathbf{p}_j$ ’s by Eqn.(1) and (2). Finally, we discard the weights and basis parameters, and only store the state output discrete densities. Thus, the final SHDDHMM for decoding has the same architecture of traditional DHMM, but just with a much larger codebook. In other words, the attractive properties of DHMM are naturally inherited by SHDDHMM (e.g., it takes only  $O(1)$  table lookup to get the state output probability).

### 2.2.2. Connection with Other Models

SHDDHMM shares some resemblance with the semi-continuous or tied-mixture HMM (SCHMM). As a result, many techniques developed for SCHMM, such as adapting mixture weights for speaker adaptation [6], are also applicable to SHDDHMM. Despite the similarity, there are two key differences between the two models. Firstly,

<sup>3</sup>Strictly speaking, due to the constraint on  $c_{jm}$ ’s,  $\mathbf{p}_j$  only lies in part of the subspace spanned by bases  $\mathbf{b}_1, \dots, \mathbf{b}_M$ .

SHDDHMM is a discrete model, while SCHMM is a continuous model. Secondly, unlike SCHMM, the density pool and state-dependent weights of SHDDHMM only exist during the training stage. The final SHDDHMM for decoding only stores the state output discrete densities, and discards the pool and mixture weights.

The subspace concept of SHDDHMM also has some connection to the recently proposed subspace Gaussian mixture model (SGMM) [7], except that SGMM uses subspace technique to derive means of the state, while SHDDHMM uses it to generate the state output discrete table.

SHDDHMM is also related to subvector-quantized discrete mixture HMM (DMHMM) [5], which also uses subvector partition. However, in the final model, for each state, DMHMM still keeps a mixture of discrete tables, while SHDDHMM only stores a single discrete table. As a result, to get the state output probability during decoding, DMHMM needs to evaluate each discrete mixture and sum up their probabilities, while SHDDHMM only takes one table lookup in  $O(1)$  time, irrespective of the number of mixtures. Furthermore, the subspace technique to reduce parameters in SHDDHMM is not utilized in DMHMM.

### 2.3. Training of SHDDHMM

Instead of flat-start training, a 1-stream SHDDHMM with  $M$  bases is first initialized from SCHMM, and is then reestimated. The basic steps of generating SHDDHMM are listed:

STEP 1: A 1-stream SCHMM with a global pool of  $M$  diagonal-covariance Gaussians is first trained.

STEP 2: Initialization of SHDDHMM: The mixture weights  $c_{jm}$ ’s of SHDDHMM are copied from SCHMM. The SVQ discrete density tables  $\mathbf{b}_m^i$ ’s of SHDDHMM are initialized using SCHMM as:

$$\mathbf{b}_m^i(l) = \frac{\sum_{t: \text{SVQ}_i(\mathbf{x}_t^i)=l} \gamma_t(m)}{\sum_{t=1}^T \gamma_t(m)} \quad (3)$$

where  $l$  is a SVQ codeword,  $\mathbf{b}_m^i(l)$  is the probability value from the entry indexed by  $l$  of SVQ discrete table  $\mathbf{b}_m^i$ ,  $\gamma_t(m)$  is the posterior probability of observation  $\mathbf{x}_t$  belonging to mixture  $m$  at time  $t$  given the SCHMM and the training observation sequence  $\mathbf{x}_1, \dots, \mathbf{x}_T$ .

STEP 3: Retraining of SHDDHMM: The updating of SVQ discrete density tables in SHDDHMM uses exactly the same Eqn.(3), except that the posterior probability  $\gamma_t(m)$  is based on SHDDHMM of the previous iteration, rather than SCHMM. SHDDHMM weights  $c_{jm}$ ’s are also updated. This step is usually run for a few iterations.

STEP 4: Generation of final SHDDHMM: After the reestimation is complete, the state output discrete densities  $\mathbf{p}_j$ ’s are computed using the weights  $c_{jm}$ ’s and SVQ discrete tables  $\mathbf{b}_m^i$ ’s by Eqn.(1) and (2). The final SHDDHMM only stores the state output discrete densities, and discards the weights and SVQ discrete tables.

## 3. EXPERIMENTAL EVALUATION

The proposed SHDDHMM was evaluated on Resource Management (RM) and Wall Street Journal (WSJ) databases. For both databases, the conventional 39-dimensional MFCC vectors (with energy, delta, and delta-delta) were extracted.

**Table 1.** Baseline model performance (WER) on WSJ SI-84

System (Task)	CDHMM	DHMM	SCHMM
SI-84 (Nov'92 5K closed)	4.46%	7.90%	4.82%

The RM systems used the standard SI-109 training data, and were tested on Feb'89 and Feb'91 test set, with the standard word-pair grammar. The WSJ systems were trained using the WSJ0 SI-84 training data with 7,138 utterances, and evaluated on the Nov'92 5K closed vocabulary task, together with the standard 5K-word trigram language model. For each task, a development data set was set out to tune system parameters. All systems used cross-word triphones.

The HTK software was modified for SHDDHMM training and decoding. All experiments were run on a Linux machine with a 3GHz CPU and 4GB RAM. All recognition results are shown in word error rate (WER). The discussion will focus mainly on the WSJ task; for RM experiments, only the final WERs are reported.

### 3.1. Baseline Model Performance on WSJ SI-84

For SI-84, SHDDHMM was compared with 1-stream CDHMM, traditional 4-stream DHMM, and 4-stream SCHMM. The standard 4 streams were adopted in our 4-stream models: the 4 streams consist of the 12 static MFCCs, 12  $\Delta$ MFCCs, 12  $\Delta\Delta$ MFCCs, and the 3 energies respectively. All models had 3,132 tied states. The CDHMM had 16 mixture components per state. The 4-stream DHMM had 512 codewords for each MFCC stream, and 256 codewords for the energy stream. Finally, the 4-stream SCHMM, which was used to initialize the SHDDHMM, had 1,024 tied mixtures per stream; only the top 128 Gaussians were used for its decoding (we found in the experiments that top 128 Gaussians for decoding in SCHMM is a good trade-off between accuracy and speed). For each model, we tried out best effort to figure out its optimal setting.

All the baseline model results are shown in Table 1. These results are comparable to others in the literature for the same task. As expected, compared with CDHMM, DHMM has a significant degradation in the recognition performance. SCHMM is close to CDHMM in WER, with only 0.36% absolute degradation.

### 3.2. Subvector-partition and Bit-allocation Schemes

4-stream SHDDHMMs of different subvector-partition and bit-allocation schemes were initialized from SCHMM. For fair comparison, all the SHDDHMMs had the same total number of bits for each stream, which were 16, 16, 16 and 15 respectively (i.e., the codebook size for stream 1 to 4 was  $2^{16}$ ,  $2^{16}$ ,  $2^{16}$  and  $2^{15}$  respectively, which resulted in a model size of roughly 1G bytes). We came up with these codebook size by considering both the memory and training data limit of our system. It can be seen that the codebook size of SHDDHMM is significantly larger than traditional DHMM, which was only  $2^9$ ,  $2^9$ ,  $2^9$  and  $2^8$  for each of the four streams.

Three subvector-partition and bit-allocation schemes used in the experiment are depicted in Table 2, with each of the last 3 columns representing a scheme. Each scheme describes a different subvector partition and its corresponding bit allocation gives the best WER of that partition<sup>4</sup>. The notation of the scheme in Table 2 is explained

<sup>4</sup>We did not exhaust all possible bit allocations. Instead, we use the heuristic of putting more bits to the lower-order MFCCs, making use of the fact that the lower-order MFCCs are more important for recognition task.

**Table 2.** Different subvector-partition and bit-allocation schemes for 4-stream SHDDHMM

Streams	Scheme s1	Scheme s2	Scheme s3
Stream 1	$[1(2)]^4[1(1)]^8$	$3(6)3(4)[3(3)]^2$	$6(9)6(7)$
Stream 2	$[1(2)]^4[1(1)]^8$	$3(6)3(4)[3(3)]^2$	$6(9)6(7)$
Stream 3	$[1(2)]^4[1(1)]^8$	$3(6)3(4)[3(3)]^2$	$6(9)6(7)$
Stream 4	$[1(5)]^3$	$[1(5)]^3$	$[1(5)]^3$

**Table 3.** Performance of SHDDHMM with different subvector-partition and bit-allocation schemes on WSJ SI-84

Scheme	WER (Reestimation?)	
	no	yes
s1	6.82%	6.11%
s2	5.62%	5.32%
s3	5.53%	<b>4.89%</b>

as follows. The basic building block is  $w(b)$ , where  $w$  is the subvector dimension, and  $b$  is the number of bits allocated to the subvector. And  $[w(b)]^C$  is equivalent to  $\underbrace{w(b) \cdots w(b)}_{C \text{ in total}}$ . For example, in

scheme s2, stream 1 (consisting of 12 static MFCCs) was split into four 3-dimensional subvectors, and 6 bits were allocated to the first subvector, 4 bits were allocated to the second subvector, and 3 bits were allocated to each of the last two subvectors.

In general, the subvector dimension increases from scheme s1 to s3. We did not try subvector dimensions larger than the one in scheme s3, because the time of finding codewords increases significantly after that. (In the extreme case, SVQ becomes full-space VQ, and the time of finding codewords will be  $O(N)$  where  $N$  is the size of the full-space VQ codebook.)

### 3.3. Effects of Schemes and Re-estimation on WSJ SI-84

Table 3 lists the SHDDHMM results with the three different subvector-partition and bit-allocation schemes. For reestimation, the column labelled with "no" gives results of the model that was initialized from an SCHMM without any reestimation, whereas the column labelled with "yes" gives the results of the final reestimated SHDDHMM. From the table, we observe that:

- From scheme s1 to s3, the WER keeps dropping. As expected, the performance is usually better when the subvector dimension is larger, as larger dimension results in smaller quantization error in full space.
- Initialization of SHDDHMM from SCHMM is effective (i.e., even without reestimation, the model is reasonably good), which provides a good starting point for reestimation.
- For each scheme, model reestimation always reduces the WER significantly.

### 3.4. Recognition Accuracy of Various Systems

Besides the WSJ task, we also tested the various modeling methods on the RM task. The RM models consisted of 1,589 tied states. The CDHMM had 6 mixtures per state, and the SCHMM had 512 tied-mixtures per stream. Their SHDDHMM counterpart was built in a

**Table 4.** Model performance (WER) on RM and WSJ

System (Task)	CDHMM	SCHMM	SHDDHMM
RM (Feb'89)	2.85%	2.77%	2.66%
RM (Feb'91)	2.94%	2.94%	2.98%
WSJ SI84 (Nov'92 5k)	4.46%	4.82%	4.89%

procedure similar to the one used in the WSJ task. Table 4 compares their performances on both RM and WSJ.

It is observed that, for each task, the performance of SHDDHMM is similar to that of its SCHMM counterpart. Since SCHMM could be viewed as a continuous version of SHDDHMM, it shows that the quantization error inherent in a discrete model is largely alleviated in SHDDHMM. On the WSJ task, the performance of both SCHMM and SHDDHMM is slightly worse than that of CDHMM. The small performance degradation is probably due to the incorrect stream independence assumption underlying SCHMM and SHDDHMM. However, all the degradations are not statistically significant (according to the NIST matched pair sentence segment word error test with confidence level 0.95).

It is also worth pointing out that the performance difference between our SHDDHMM and CDHMM is smaller than some recently proposed discrete models (e.g., [3]) on a similar task. Since many techniques were tried in [3] to reduce the quantization error in various discrete models using traditional codebook sizes, and yet the performance of the resulting discrete models were still incomparable with their CDHMM counterparts, we believe that it is important to build DHMM with a large codebook in order for it to perform comparably with CDHMM.

### 3.5. Computation Time on WSJ SI-84

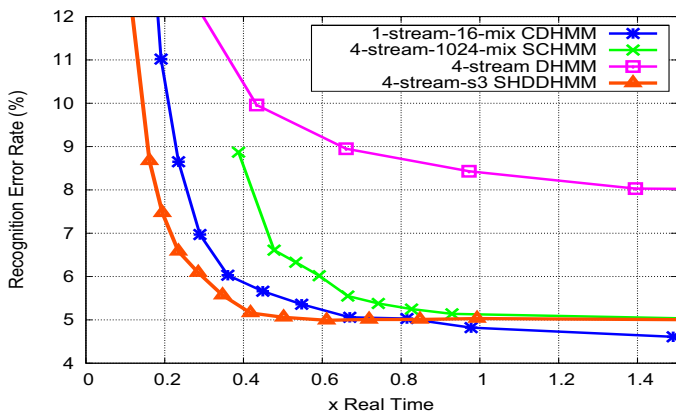
By fixing the number of tied states, we investigated the likelihood computation time for both CDHMM and SHDDHMM. The computation time ratio between the two models is found to be 40:1; it indicates that 98% of GMM computation was cut by SHDDHMM, as the latter only needs  $O(1)$  table lookup time to get likelihood. This figure is larger than most other fast GMM computation techniques. E.g., in [8], Chan implemented a set of fast GMM computation techniques and found that 80% cut-off was an empirical upper bound.

Operating curves of various SI-84 models were obtained by changing the beam width during Viterbi decoding, and they are compared in Fig. 2. The results show that SHDDHMM performs better than SCHMM and conventional DHMM. Compared with CDHMM with the same WER, SHDDHMM cuts the decoding time by 25% in most part of the operating curve. The 25% speedup is quite substantial because in our CDHMM system, the Gaussian computation takes at most 27%<sup>5</sup> of the total recognition time.

### 3.6. Bandwidth Saving

The conventional acoustic vectors in CDHMM consist of 39 MFCC coefficients, each of which is represented by a 4-byte float. Thus,  $4 \times 39 = 156$  bytes are used to encode a frame of MFCC coefficients. In contrast, the codeword in each stream of our 4-stream

<sup>5</sup>In our experiments, HDecode (in HTK) was used for decoding with block caching mechanism, which caches the state output likelihoods and saves Gaussian computation time. As a result, in our WSJ CDHMM system, when computing the time ratio between GMM computation and the total decoding for different beam widths, the maximum ratio is found to be 27%.



**Fig. 2.** Operating characteristics of various SI-84 models (finding codeword time is included for discrete models)

SHDDHMM is represented by a 2-byte unsigned short integer. As a consequence, only  $2 \times 4 = 8$  bytes are used to encode the discrete features in a frame, resulting in a 20-fold saving of bandwidth in DSR systems. Even for systems which only transfer the 13 static MFCCs and construct the full 39-dimensional MFCC vectors at the server side, the SHDDHMM feature could still result in a 6.5-fold saving of bandwidth.

## 4. CONCLUSION & FUTURE WORK

In this paper, we proposed to reduce the quantization error of DHMM with a very large codebook. Subspace modeling technique with constrained bases is applied to reduce the number of model parameters so that re-estimation of such high-density model becomes feasible. On both RM and WSJ 5K task, the proposed SHDDHMM gets similar recognition performance as that of its CDHMM counterpart, with faster decoding speed and lower bandwidth requirement.

Further directions involve generalizing the model to larger tasks, where the size of SHDDHMM is expected to be even larger.

## 5. REFERENCES

- [1] M. J. F. Gales *et al.*, "State-based Gaussian selection in large vocabulary continuous speech recognition using HMMs," in *IEEE Trans. on SAP*, vol. 7, pp. 152–161, 1999.
- [2] Jasha Droppo *et al.*, "Towards a non-parametric acoustic model: an acoustic decision tree for observation probability calculation," in *Proc. of Interspeech*, Sept, 2007.
- [3] J. Ajmera and M. Akamine, "Decision tree acoustic models for ASR," in *Proc. of Interspeech*, Sept, 2009.
- [4] Brian Mak *et al.*, "High-density discrete HMM with the use of scalar quantization indexing," in *Proc. of Eurospeech*, Sept 2005.
- [5] V. Digalakis *et al.*, "Efficient speech recognition using subvector quantization and discrete-mixture HMMs," in *Computer Speech and Language*, vol. 14, pp. 33–46, 2000.
- [6] X. Huang and K. F. Lee, "On speaker-independent, speaker-dependent, and speaker-adaptive speech recognition," in *IEEE Trans. on SAP*, vol. 1, pp. 150–157, 1993.
- [7] Daniel Povey *et al.*, "Subspace Gaussian mixture models for speech recognition," in *Proc. of ICASSP*, Mar 2010.
- [8] Chan, A. *et al.*, "Four-layer categorization scheme of fast GMM computation techniques in large vocabulary continuous speech recognition systems," in *Proc. of ICSLP*, 2004.