

## Speech Recognition for English E-Set

- E-Set: {b, c, d, e, g, p, t, v, z}, total nine English letters with the same ending vowel “E”.
- Speech Recognition: conducted in a multi-speaker, isolated-word mode.
- Training Set: 900 tokens from 100 speakers
- Testing Set: another 900 tokens from the same 100 speakers
- Speech Signal:
  - ◆ Sampled at 6.67 kHz,
  - ◆ Processed by using an eighth order LPC analysis, with a 45 ms window and a 15 ms shift.
  - ◆ A 24-component feature vector, consisting of 12 LPC-derived liftered cepstral coefficients and 12 delta cepstral coefficients.



## Speech Recognition Baseline Model:

- Each letter is modeled by a 5-state left-to-right HMM
- Five Gaussian mixture components per state are used for characterizing each state observation density.
- The HMM parameters were estimated using the segmental k-means algorithm
- Recognition Rates: 61.7% for the testing set, and 80.2% for the training set.

## Discrimination-based Approaches [Su and Lee 94]

■ Input observation vector  $X$ :  $X^t = [S_1^1, S_2^1, S_3^1, \dots, S_3^9, S_4^9, S_5^9]$

- ◆ where  $S_j^i$  is either the averaged or the accumulated state log-likelihood (score) for state  $j$  or word  $i$ . There are 45 elements in this *score vector*.

■ Subspace Projection:

- ◆ Maximin Algorithm is proposed for feature selection:

- ◆ Divergence  $D_{i,j}(k)$  between the class  $i$  and the class  $j$ , for  $1 \leq i, j \leq 9$ ;  $i < j$ ,  $1 \leq k \leq 45$ .

$$D_{ij}(k) = \int_{Y_k} (P_i(Y_k) - P_j(Y_k)) \cdot \ln \left( \frac{P_i(Y_k)}{P_j(Y_k)} \right) dY_k$$

- ◆ Find  $D_{\min}(k) = \min_{i,j} D_{i,j}(k)$  for  $1 \leq i, j \leq 9$ ;  $i < j$ ,  $1 \leq k \leq 45$ .
- ◆ Sort  $D_{\min}(k)$  in descending order, then the sequence of subspaces is obtained

## Discrimination-based Approaches (cont.):

### ■ Subspace Projection (cont.):

#### ◆ Result:



◆ Table I: The 45 Features Indices Listed in the Order of Descending Divergence Values (Left to Right First, then Top to Bottom)



◆ Table II: Corresponding Divergence Values Listed in Descending Order



◆ Figure 1: Recognition rate as a function of dimensionality for a subspace-based recognizer



◆ Figure 2: Recognition rate as a function of iteration for a robust subspace-based recognizer

## Discrimination-based Approaches (cont.):

### ■ Weighted HMM Algorithm:

- ◆ A linear discriminator is used:

$$g(X, W_i) = W_i^t \cdot X$$

$$\hat{C} = \arg \max_i (W_i^t \cdot X)$$

- ◆ Result:



- ◆ Figure 3: Recognition rate as a function of iteration number for a WHMM-based recognizer



- ◆ Figure 4: Recognition rate as a function of iteration number for a robust WHMM recognizer

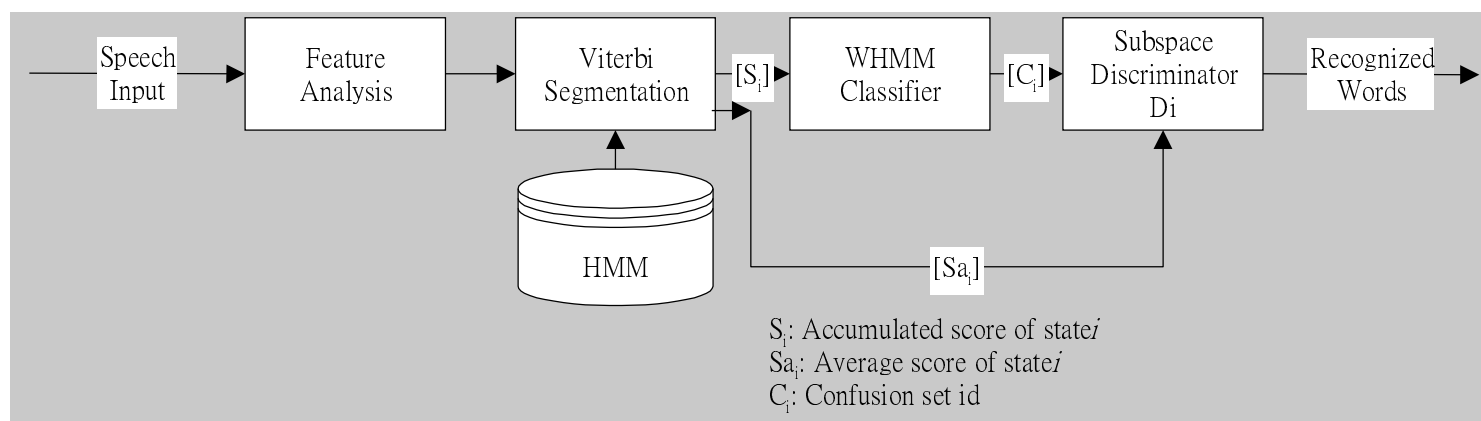


- ◆ Figure 5: Recognition rate as a function of margin ratio for a robust WHMM recognizer

## Discrimination-based Approaches (cont.):

### ■ A Two-Stage Discriminator:

- ◆ Block Diagram: A block diagram of the proposed two-stage recognition system



- ◆ Result:



- ◆ Table III: Performance of Each Confusion Set in Each Stage

## Discrimination-based Approaches (cont.):



**TABLE I**

THE 45 FEATURE INDICES LISTED IN THE ORDER OF DESCENDING  
DIVERGENCE VALUES (LEFT TO RIGHT FIRST, THEN TOP TO BOTTOM)

1 (B <sub>1</sub> )	6 (C <sub>1</sub> )	31 (T <sub>1</sub> )	41 (Z <sub>1</sub> )	26 (P <sub>1</sub> )
21 (G <sub>1</sub> )	17 (E <sub>2</sub> )	37 (V <sub>2</sub> )	36 (V <sub>1</sub> )	16 (E <sub>1</sub> )
22 (G <sub>2</sub> )	18 (E <sub>3</sub> )	11 (D <sub>1</sub> )	38 (V <sub>3</sub> )	2 (B <sub>2</sub> )
12 (D <sub>2</sub> )	10 (C <sub>5</sub> )	4 (B <sub>4</sub> )	28 (P <sub>3</sub> )	43 (Z <sub>3</sub> )
30 (P <sub>5</sub> )	20 (E <sub>5</sub> )	34 (T <sub>4</sub> )	27 (P <sub>2</sub> )	7 (C <sub>2</sub> )
32 (T <sub>2</sub> )	19 (E <sub>4</sub> )	14 (D <sub>4</sub> )	33 (T <sub>3</sub> )	45 (Z <sub>5</sub> )
13 (D <sub>3</sub> )	9 (C <sub>4</sub> )	40 (V <sub>5</sub> )	39 (V <sub>4</sub> )	42 (Z <sub>2</sub> )
25 (G <sub>5</sub> )	15 (D <sub>5</sub> )	24 (G <sub>4</sub> )	8 (C <sub>3</sub> )	29 (P <sub>4</sub> )
5 (B <sub>5</sub> )	3 (P <sub>3</sub> )	23 (G <sub>3</sub> )	35 (T <sub>5</sub> )	44 (Z <sub>4</sub> )

## Discrimination-based Approaches (cont.):



**TABLE II**

THE 45 CORRESPONDING DIVERGENCE VALUES LISTED IN DESCENDING ORDER

---

0.099410	0.066085	0.042937	0.038249	0.037496
0.028690	0.020272	0.019982	0.017487	0.011872
0.009741	0.005572	0.005305	0.004210	0.004037
0.004005	0.003803	0.003728	0.003656	0.003621
0.003214	0.003200	0.003085	0.002916	0.002690
0.002147	0.002064	0.001937	0.001900	0.001661
0.001651	0.001637	0.001496	0.001221	0.000969
0.000969	0.000596	0.000562	0.000384	0.000381
0.000229	0.000194	0.000187	0.000076	0.000042

---



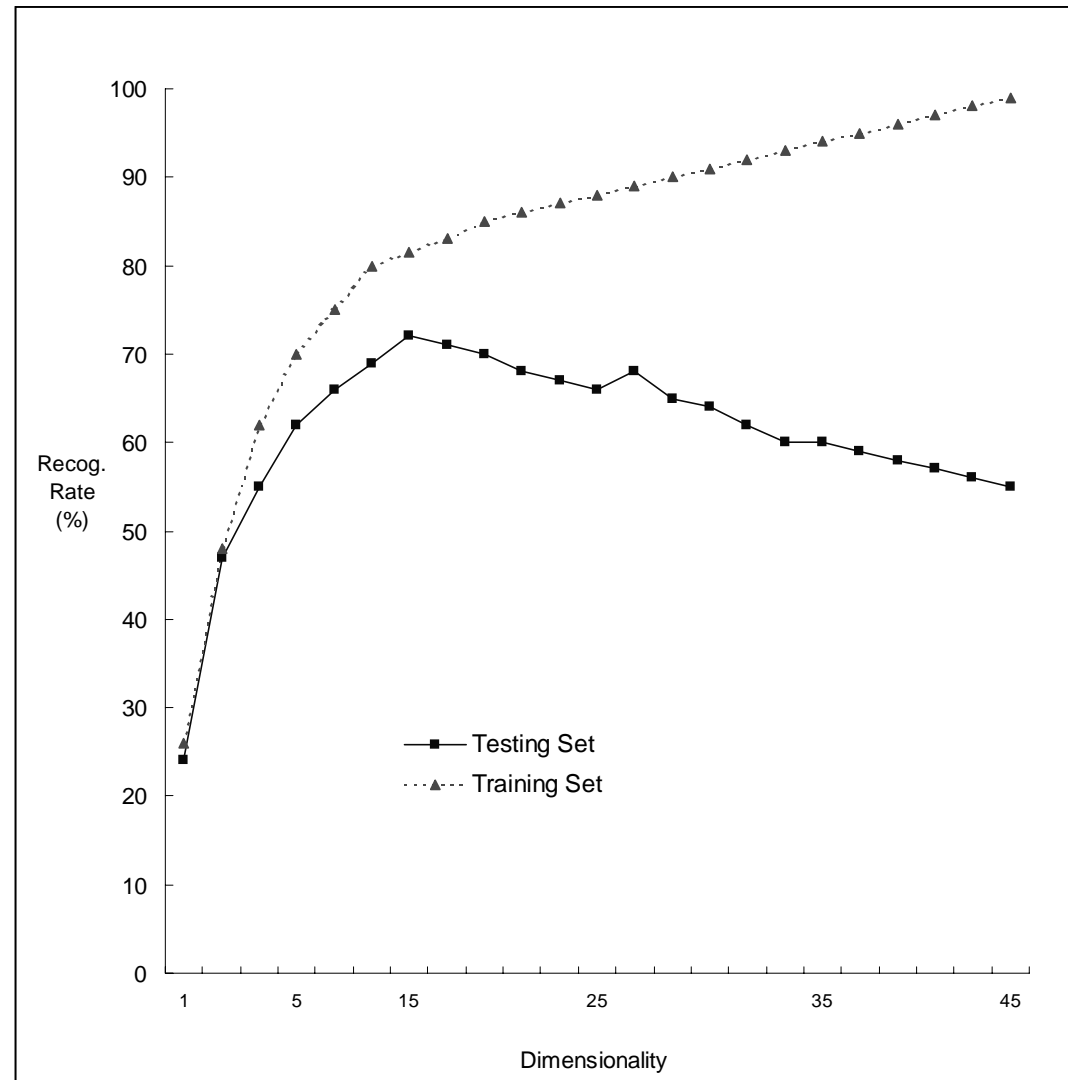
## Discrimination-based Approaches (cont.):



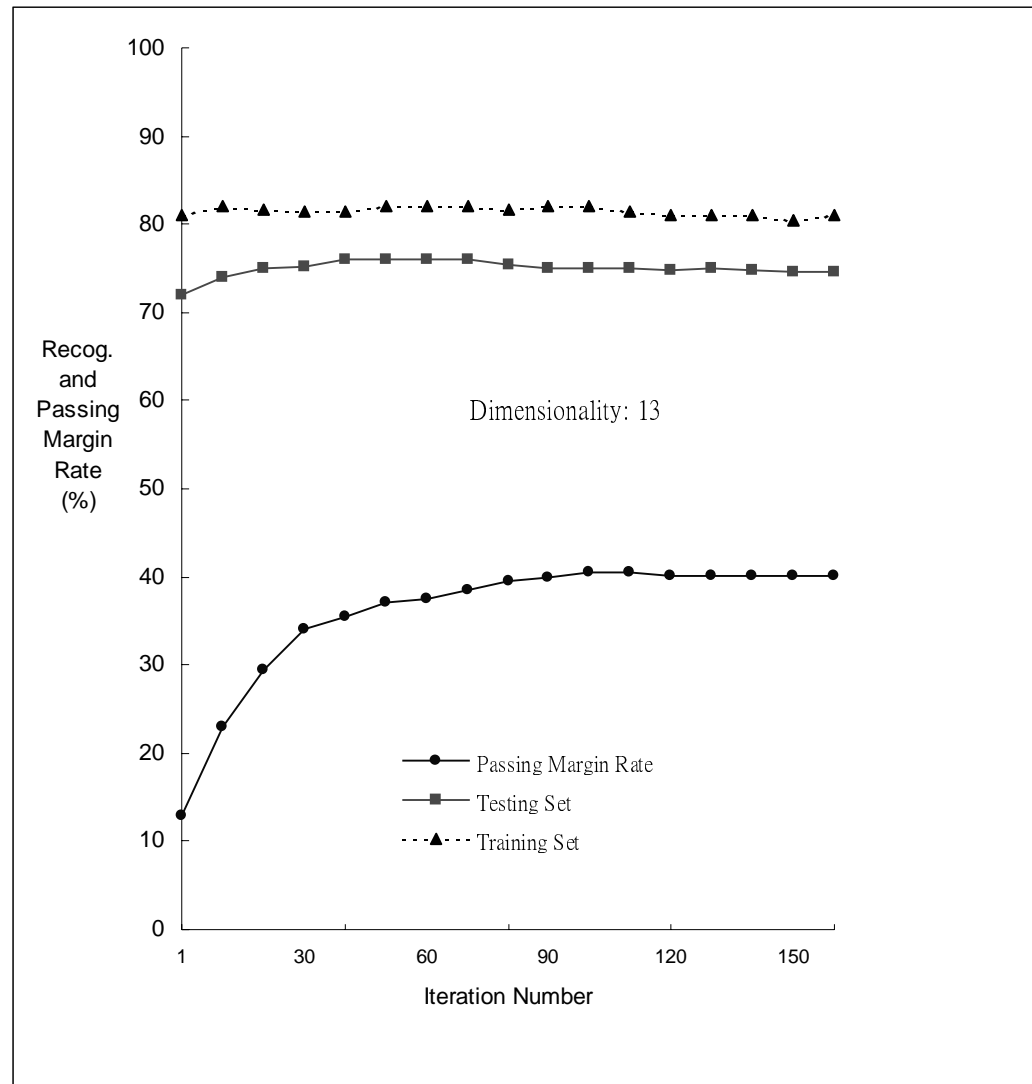
**TABLE III**  
PERFORMANCE OF EACH CONFUSION SET IN EACH STAGE

	First Stage (%)	Second Stage (%)
B	68.9	74.7
C	88.9	89.6
D	64.3	73.9
E	87.6	89.7
G	79.6	84.5
P	72.0	73.1
T	68.0	75.7
V	73.2	75.6
Z	73.3	78.1
Average	74.9	79.4

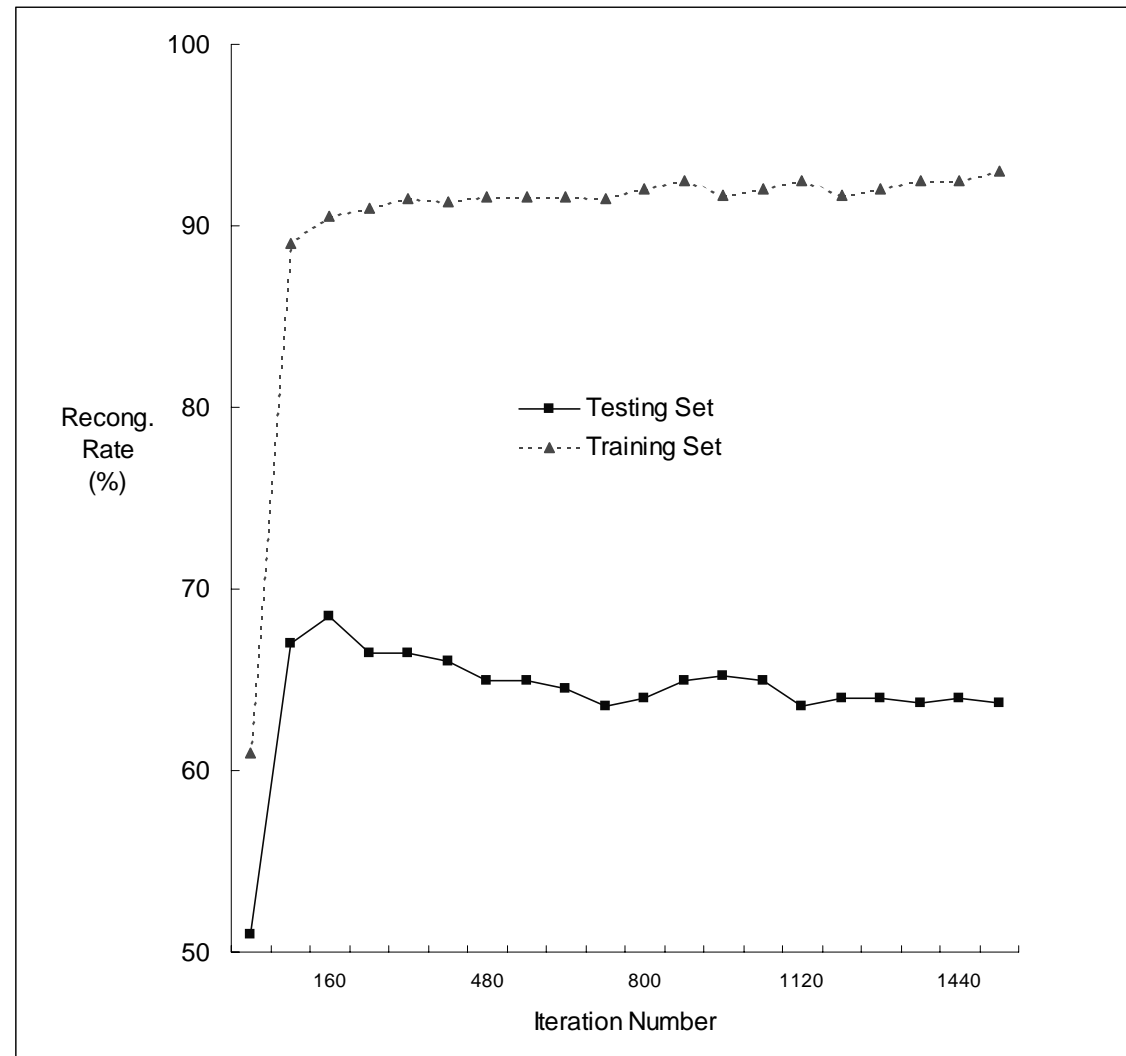
# Subspace-based Recognizer: Recognition Rate vs. Dimensionality



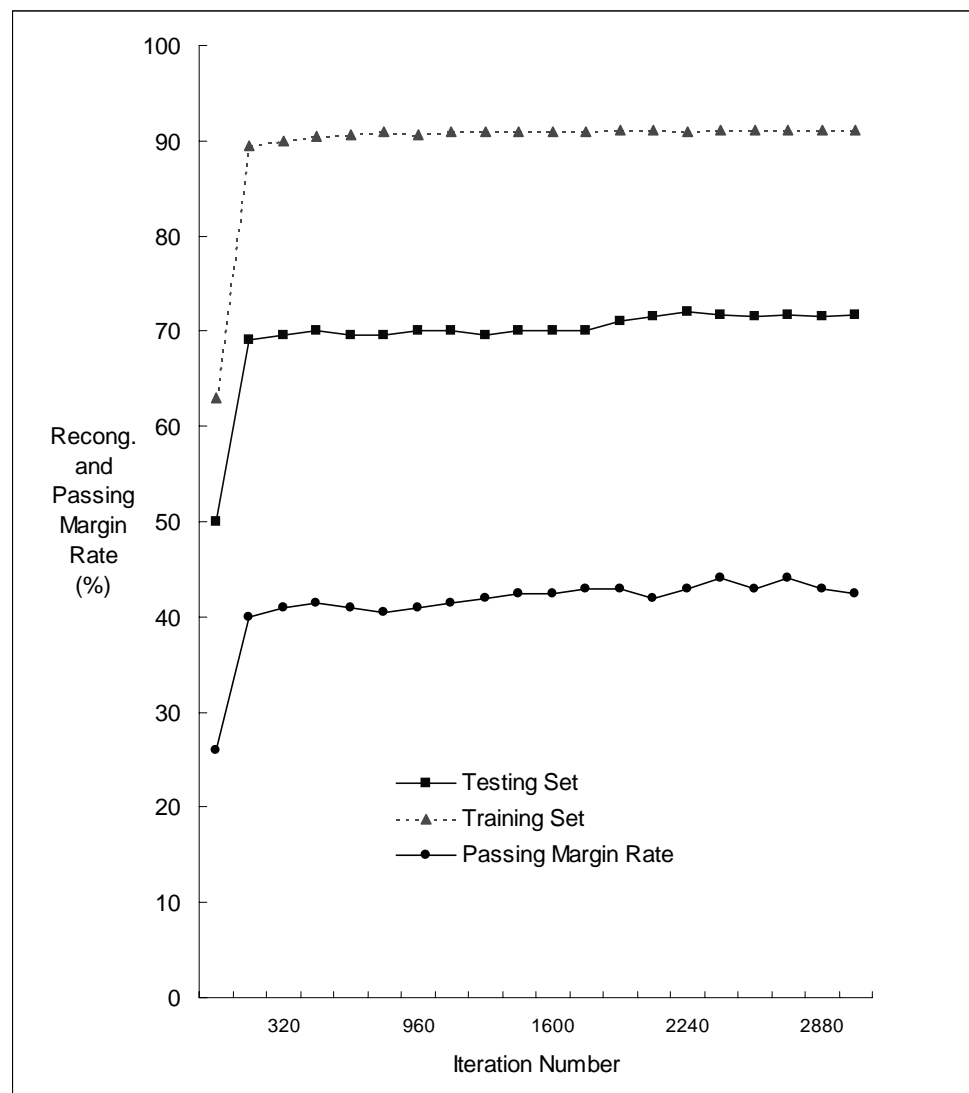
# Robust Subspace-based Recognizer: Recognition Rate vs. Iteration



# WHMM-based Recognizer: Recognition Rate vs. Iteration Number



# Robust WHMM Recognizer: Recognition Rate vs. Iteration Number



# Robust WHMM Recognizer: Recognition Rate vs. Margin Ratio

