

Phoneme Recognition Using Time-Delay Neural Networks

A. Waibel T. Hanazawa G. Hinton* K. Shikano
K. Lang[†]

ATR Interpreting Telephony Research Laboratories

October 30, 1987

*University of Toronto and Canadian Institute for Advanced Research
[†]Carnegie-Mellon University

TR-I-0006

Phoneme Recognition
Using Time-Delay Neural Networks

時間遅れ神経回路網による音韻認識

A. Waibel, T. Hanazawa, G. Hinton,
K. Shikano and K. Lang

ワイベル、花沢利行、ヒントン、鹿野清宏、ラング

1987.10

Abstract

In this paper we present a Time Delay Neural Network (TDNN) approach to phoneme recognition which is characterized by two important properties: 1.) Using a 3 layer arrangement of simple computing units, a hierarchy can be constructed that allows for the formation of arbitrary nonlinear decision surfaces. The TDNN learns these decision surfaces automatically using error back-propagation[1]. 2.) The time-delay arrangement enables the network to discover acoustic-phonetic features and the temporal relationships between them independent of position in time and hence not blurred by temporal shifts in the input.

As a recognition task, the speaker-dependent recognition of the phonemes "B", "D", and "G" in varying phonetic contexts was chosen. For comparison, several discrete Hidden Markov Models (HMM) were trained to perform the same task. Performance evaluation over 1946 testing tokens from three speakers showed that the TDNN achieves a recognition rate of 98.5 % correct while the rate obtained by the best of our HMMs was only 93.7 %. Closer inspection reveals that the network "invented" well-known acoustic-phonetic features (e.g., F2-rise, F2-fall, vowel-onset) as useful abstractions. It also developed alternate internal representations to link different acoustic realizations to the same concept.

ATR Interpreting Telephony Research Laboratories
ATR 自動翻訳電話研究所

Contents

1	Introduction	6
2	Time Delay Neural Networks (TDNN)	7
2.1	A TDNN Architecture for Phoneme Recognition	8
2.2	Learning in a TDNN	11
3	Hidden Markov Models (HMM)	14
3.1	An HMM for Phoneme Recognition	15
3.2	Learning in an HMM	15
4	Recognition Experiments	16
4.1	Experimental Conditions	16
4.2	Results	17
4.3	The Learned Internal Representations of a TDNN	18
5	Conclusion and Summary	20

List of Figures

1	A Time Delay Neural Network (TDNN) unit	9
2	The architecture of the TDNN	10
3	TDNN output error vs. number of learning iterations (increasing training set size)	12
4	Hidden Markov Model	15
5	Learning in a Hidden Markov Model	16
6	Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "B"s using a HMM	22
7	Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "D"s using a HMM	23
8	Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "G"s using a HMM	24
9	Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "B"s using a TDNN	25
10	Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "D"s using a TDNN	26
11	Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "G"s using a TDNN	27
12	TDNN Activation patterns for "DA"	28
13	TDNN Activation patterns for "DO"	29
14	Weights on connections from 16 coefficients over 3 time frames to each of the 8 hidden units in the first layer	30
15	TDNN Activation patterns for "GA" embedded in an utterance	31
16	TDNN Activation patterns for "GA" in utterance initial position	32
17	TDNN Activation patterns for "DO" misaligned by +30 msec	33
18	TDNN Activation patterns for "DO" misaligned by -30 msec	34

List of Tables

1	Recognition results for three speakers over test data using TDNN and HMM	18
---	---	----

Abstract

In this paper we present a Time Delay Neural Network (TDNN) approach to phoneme recognition which is characterized by two important properties: 1.) Using a 3 layer arrangement of simple computing units, a hierarchy can be constructed that allows for the formation of arbitrary nonlinear decision surfaces. The TDNN learns these decision surfaces automatically using error back-propagation[1]. 2.) The time-delay arrangement enables the network to discover acoustic-phonetic features and the temporal relationships between them independent of position in time and hence not blurred by temporal shifts in the input.

As a recognition task, the speaker-dependent recognition of the phonemes "B", "D", and "G" in varying phonetic contexts was chosen. For comparison, several discrete Hidden Markov Models (HMM) were trained to perform the same task. Performance evaluation over 1946 testing tokens from three speakers showed that the TDNN achieves a recognition rate of 98.5 % correct while the rate obtained by the best of our HMMs was only 93.7 %. Closer inspection reveals that the network "invented" well-known acoustic-phonetic features (e.g., F2-rise, F2-fall, vowel-onset) as useful abstractions. It also developed alternate internal representations to link different acoustic realizations to the same concept.

Acknowledgement

The authors would like to express their gratitude to Dr. Akira Kurematsu, president of ATR Interpreting Research Laboratories, for his enthusiastic encouragement and support, which made this research possible. We are also indebted to the members of the Speech Processing Department at ATR and Mr. Fukuda at Apollo Computer for assistance in the various stages of this research.

1 Introduction

In recent years, the advent of new learning procedures and the availability of high speed parallel supercomputers have given rise to a renewed interest in connectionist models of intelligence[1]. These models are particularly interesting for cognitive tasks that require massive constraint satisfaction, i.e., the parallel evaluation of many clues and facts and their interpretation in the light of numerous interrelated constraints. Cognitive tasks, such as vision, speech, language processing and motor control are also characterized by a high degree of uncertainty and variability and it has proved difficult to achieve good performance for these tasks using standard serial programming methods. Complex networks composed of simple computing units are attractive for these tasks not only because of their "brain-like" appeal, but because they offer ways for automatically designing systems that can make use of multiple interacting constraints. In general, such constraints are too complex to be easily programmed and require the use of automatic learning strategies. Such learning algorithms now exist (For an excellent review, see Lippman[2]) and have been demonstrated to discover interesting internal abstractions, in their attempts to solve a given problem[1,3,4,5]. Learning is most effective when used in an architecture that is appropriate for the task. Indeed, the experiments reported in this paper suggest that as much prior knowledge as possible should be built into the network.

Naturally, these techniques will have far-reaching implications for the design of automatic speech recognition systems, if proven successful in comparison to already existing techniques. Lippmann[6] has compared several kinds of neural networks with other classifiers and evaluated their ability to create complex decision surfaces. Other studies have investigated actual speech recognition tasks and compared them to psychological evidence in speech perception[7] or to existing speech recognition techniques[8,9]. Speech recognition experiments using neural nets have so far mostly been aimed at isolated word recognition (mostly the digit recognition task) [10,11,12,13] or phonetic recognition with predefined constant[14,15] or variable phonetic contexts[16,14,17].

A number of these studies report very encouraging recognition performance[16], but only few comparisons to existing recognition methods exist. Some of these comparisons found performance similar to existing methods[9,11], but others found that networks perform worse than other techniques[8]. One might argue that this state of affairs is encouraging considering the amount of fine-tuning that has gone into optimizing the more popular, established techniques. Nevertheless, better comparative performance figures are needed before neural networks can be considered as a viable alternative for speech recognition systems.

One possible explanation for the mixed performance results obtained so far may be limitations in computing resources leading to short-cuts that limit performance. Another more serious limitation, however, is the inability of most neural network architectures to deal properly with the dynamic nature of speech. Two important aspects of this are for a network to represent temporal relation-

ships between acoustic events, while at the same time providing for invariance under translation in time. The specific movement of a formant in time, for example, is an important cue to determining the identity of a voiced stop, but it is irrelevant whether the same set of events occurs a little sooner or later in the course of time. Without translation invariance a neural net requires precise segmentation, to align the input pattern properly. Since this is not always possible in practice, learned features tend to get blurred (in order to accommodate slight misalignments) and their performance deteriorates.

In the present paper, we describe a Time Delay Neural Network (TDNN), which addresses both of these aspects and demonstrate through extensive performance evaluation that superior recognition results can be achieved using this approach. In the following section, we begin by introducing the architecture and learning strategy of a TDNN aimed at phoneme recognition. Next, we compare the performance of our TDNNs with one of the more popular current recognition techniques. In section 3, we therefore describe several Hidden Markov Models (HMM), under development at ATR[18]. Both techniques are then evaluated over a testing database. We report the results in section 4 of this paper and show that substantially higher recognition performance is achieved by the TDNN than by the best of our HMMs. We also take a close look at the internal representation that the TDNN learns for this task. It discovers a number of interesting linguistic abstractions which we show by way of examples. The implications of these results are then discussed and summarized in the final section of this paper.

2 Time Delay Neural Networks (TDNN)

To be useful for speech recognition, a layered feed forward neural network must have a number of properties. First, it should have multiple layers and sufficient interconnections between units in each of these layers. This is to ensure that the network will have the ability to learn complex non-linear decision surfaces[2,6]. Second, the network should have the ability to represent relationships between events in time. These events could be spectral coefficients, but might also be the output of higher level feature detectors. Third, the actual features or abstractions learned by the network should be invariant under translation in time¹. Fourth, the learning procedure should not require precise temporal alignment of the labels that are to be learned. Fifth, the number of weights in the network should be small compared to the amount of training data so that the network is forced to encode the training data by extracting regularity. In the following, we describe a TDNN architecture that satisfies all of these criteria and is designed explicitly for the recognition of phonemes, in particular, the voiced stops "B", "D" and "G".

¹In vision, solutions to the similar problem of shift-invariance have been proposed by use of a "Neocognitron"[19].

2.1 A TDNN Architecture for Phoneme Recognition

The basic unit used in many neural networks computes the weighted sum of its inputs and then passes this sum through a non-linear function, most commonly a threshold or sigmoid function[2.1]. In our TDNN, this basic unit is modified by introducing delays D_1 through D_N as shown in Fig.1. The J inputs of such a unit now will be multiplied by several weights, one for each delay and one for the undelayed input. For $N = 2$, and $J = 16$, for example, 48 weights will be needed to compute the weighted sum of the 16 inputs, with each input now measured at three different points in time. In this way a TDNN unit has the ability to relate and compare current input with the past history of events. The sigmoid function was chosen as the non-linear output function F due to its convenient mathematical properties[20,5].

For the recognition of phonemes, a three layer net is constructed². Its overall architecture and a typical set of activities in the units are shown in Fig.2.

At the lowest level, 16 melscale spectral coefficients serve as input to the network. Input speech, sampled at 12 kHz, was hamming windowed and a 256-point FFT computed every 5 msec. Melscale coefficients were computed from the power spectrum as in[21] and adjacent coefficients in time collapsed resulting in an overall 10 msec frame rate. The coefficients of an input token (in this case 15 frames of speech centered around the hand-labeled vowel onset) were then normalized to lie between -1.0 and +1.0 with the average at 0.0. Fig.2 shows the resulting coefficients for the speech token "BA" as input to the network, where positive values are shown as black and negative values as grey squares.

This input layer is then fully interconnected to a layer of 8 time delay hidden units, where $J = 16$ and $N = 2$ (i.e., 16 coefficients over three frames with time delay 0, 1 and 2). An alternative way of seeing this is depicted in Fig.2. It shows the inputs to these time delay units expanded out spatially into a 3 frame window, which is passed over the input spectrogram. Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window. The particular choice of 3 frames (30 msec) was motivated by earlier studies[22], that suggested that a 30 msec window might be sufficient to represent low level acoustic-phonetic events for phoneme recognition. It was also the optimal choice among a number of alternative designs evaluated by Lang[23] on a similar task.

In the second hidden layer, each of 3 TDNN units looks at a 5 frame window of activity levels in hidden layer 1 (i.e., $J = 8$, $N = 4$). The choice of a larger 5 frame window in this layer was motivated by the intuition that higher level units should learn to make decisions over a wider range in time based on more local abstractions at lower levels.

Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent

²Lippmann[2.6] demonstrated recently that three layers can encode arbitrary pattern recognition decision surfaces

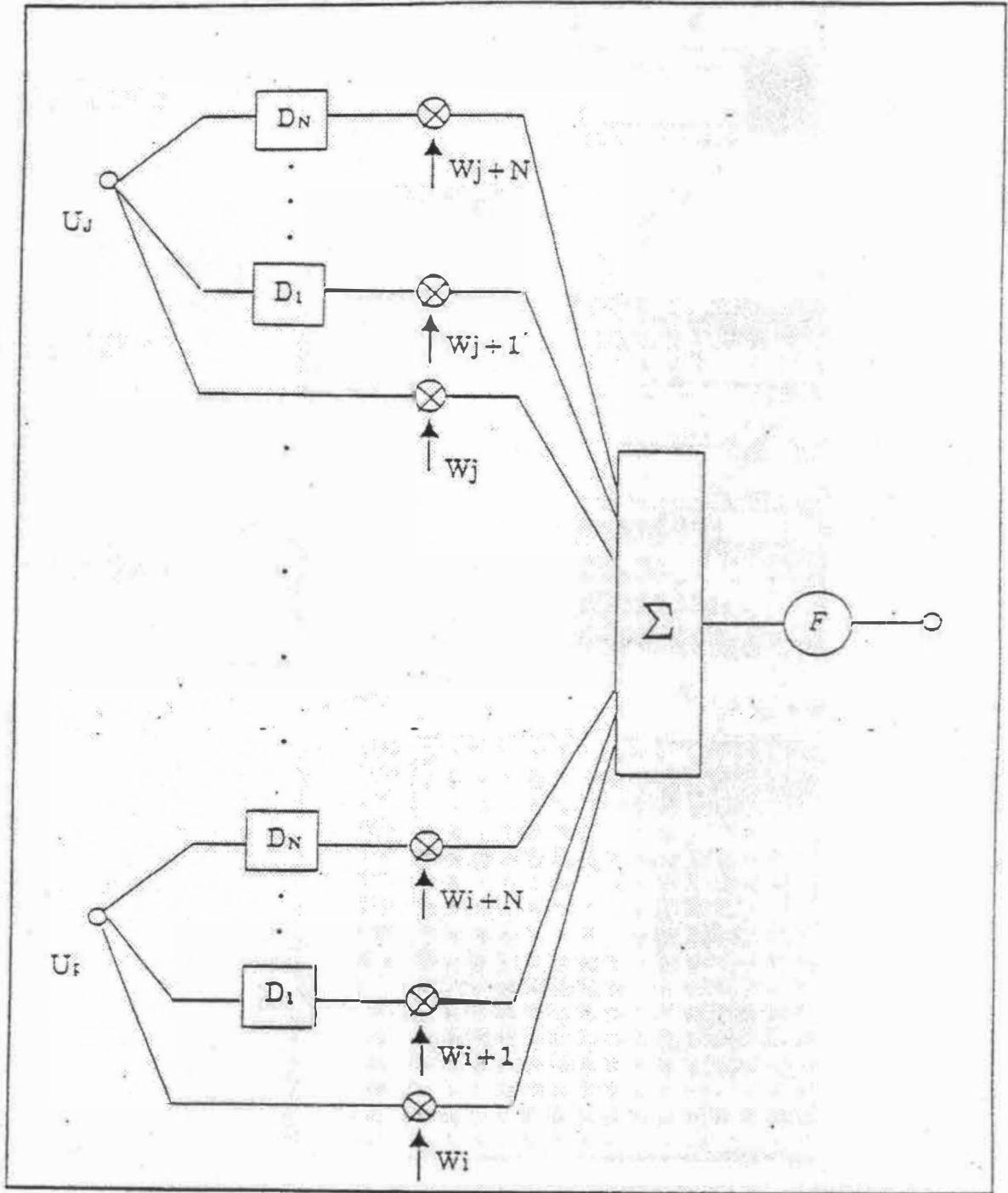


Figure 1: A Time Delay Neural Network (TDNN) unit

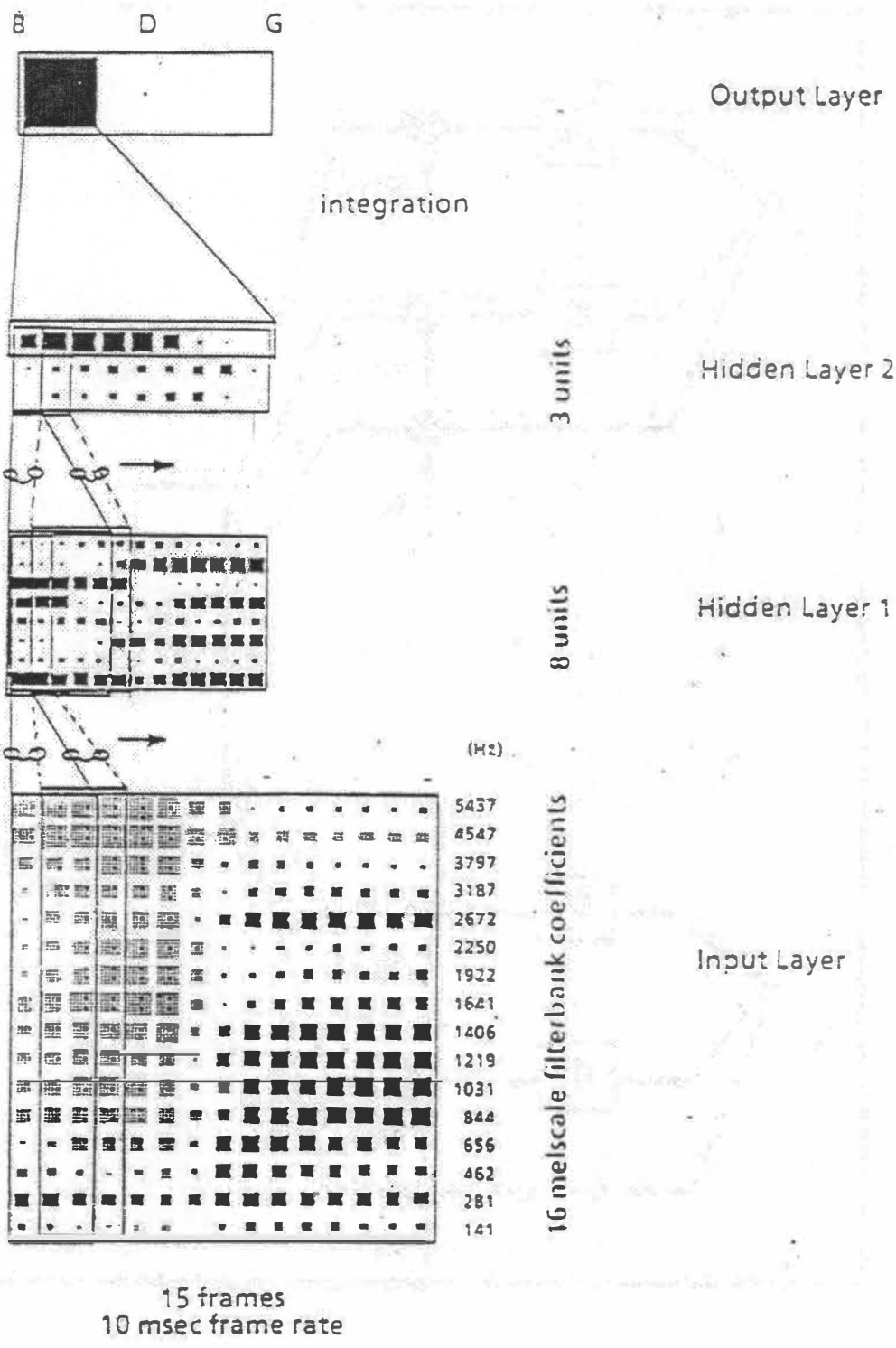


Figure 2: The architecture of the TDNN

output unit (shown in Fig.2 over 9 frames for the "B" output unit). In practice, this summation is implemented simply as another TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2³.

When the TDNN has learned its internal representation, it performs recognition by passing input speech over the TDNN units. In terms of the illustration of Fig.2 this is equivalent to passing the time delay windows over the lower level units' firing patterns. At the lowest level, these firing patterns simply consist of the sensory input, i.e., the spectral coefficients.

Each TDNN unit outlined in this section has the ability to encode temporal relationships within the range of the N delays. Higher layers can attend to larger time spans, so local short duration features will be formed at the lower layer and more complex longer duration features at the higher layer. The learning procedure ensures that each of the units in each layer has its weights adjusted in a way that improves the network's overall performance.

2.2 Learning in a TDNN

Several learning techniques exist for optimization of neural networks[1,2,24]. For the present network we adopt the Back-propagation Learning Procedure[20,5]. This procedure performs two passes through the network. During the forward pass, an input pattern is applied to the network with its current connection strengths (initially small random weights). The outputs of all the units at each level are computed starting at the input layer and working forward to the output layer. The output is then compared with the desired output and its error calculated. During the backward pass, the derivative of this error is then propagated back through the network, and all the weights are adjusted so as to decrease the error[20,5]. This is repeated many times for all the training tokens until the network converges to producing the desired output.

In the previous section we described a method of expressing temporal structure in a TDNN and contrasted this method to training a network on a static input pattern (spectrogram), which results in shift sensitive networks (i.e., poor performance for slightly misaligned input patterns) as well as less crisp decision making in the units of the network (caused by misaligned tokens during training).

To achieve the desired learning behavior, we need to ensure that the network is exposed to *sequences* of patterns and that it is allowed (or encouraged) to learn about the most powerful cues and sequences of cues among them. Conceptually, the back-propagation procedure is applied to speech patterns that are stepped through in time. An equivalent way of achieving this result is to use a spatially expanded input pattern, i.e., a spectrogram plus some constraints on

³Note, however, that as for all units in this network (except the input units), the output units are also connected to a permanently active threshold unit. In this way, the dc-bias of each output unit can still be adjusted for optimal classification.

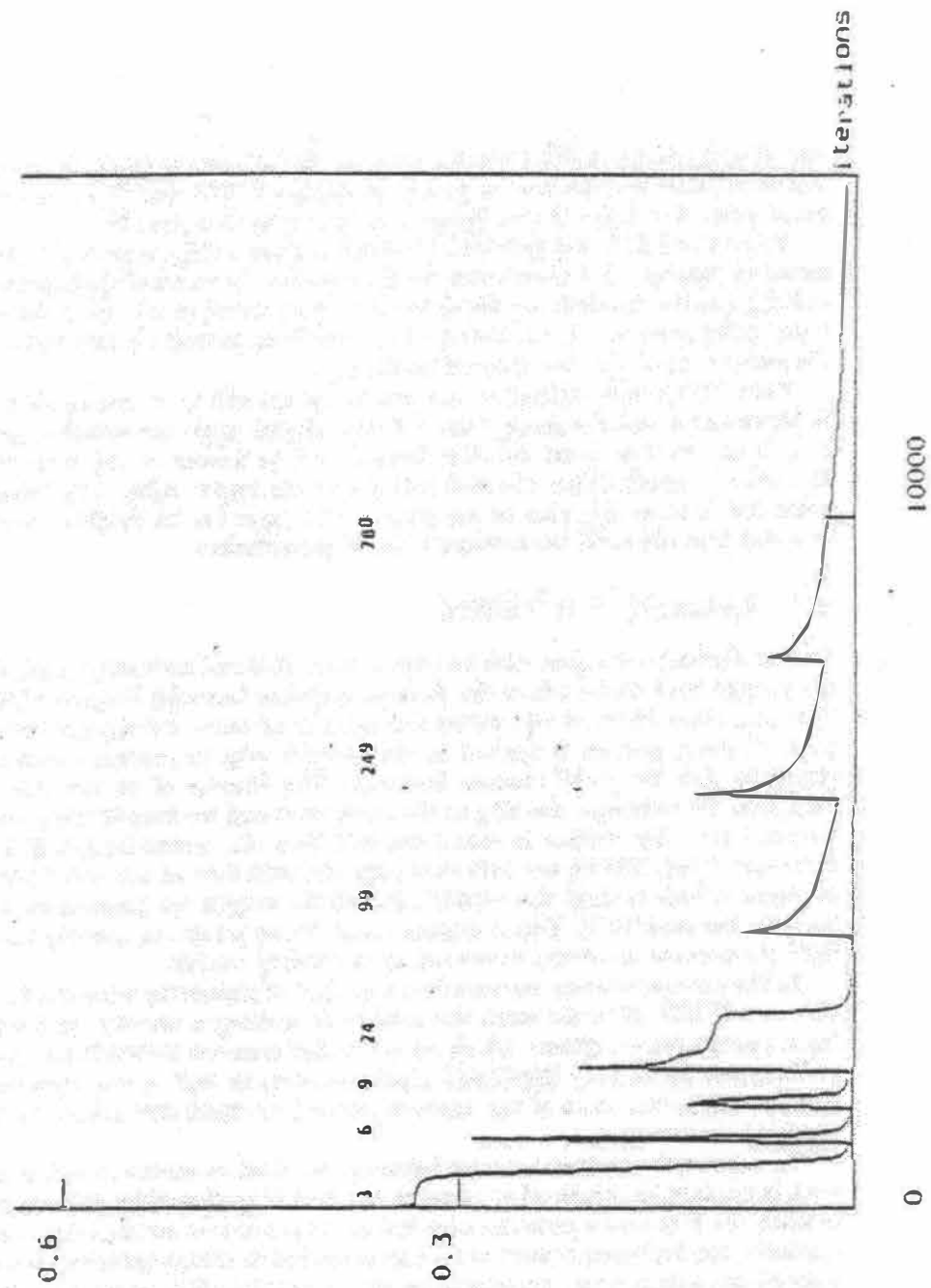


Figure 3: TDNN output error vs. number of learning iterations (increasing training set size)

the weights. Each collection of TDNN-units described above is duplicated for each one frame shift in time. In this way the whole history of activities is available at once. Since the shifted copies of the TDNN-units are mere duplicates and are to look for the same acoustic event, the weights of the corresponding connections in the time shifted copies must be constrained to be the same. To realize this, we first apply the regular back-propagation forward and backward pass to all time shifted copies as if they were separate events. This yields different error derivatives for corresponding (time shifted) connections. Rather than changing the weights on time-shifted connections separately, however, we actually update each weight on corresponding connections by the same value, namely by the *average* of all corresponding time-delayed weight changes⁴. Fig.2 illustrates this by showing in each layer only two connections that are linked to (constrained to have the same value as) their time shifted neighbors. Of course, this applies to all connections and all time shifts. In this way, the network is forced to discover useful acoustic-phonetic features in the input, regardless of when in time they actually occurred. This is an important property, as it makes the network independent of errorprone preprocessing algorithms, that otherwise would be needed for time alignment and/or segmentation. In section 4.3, we will show examples of grossly misaligned patterns that are properly recognized, due to this property.

The procedure described here is computationally rather expensive, due to the many iterations necessary for learning a complex multidimensional weight space and the number of learning samples. In our case, about 800 learning samples were used and between 20,000 and 50,000 iterations of the back-propagation loop were run over all training samples. Two steps were taken, to perform learning within reasonable time. First, we have implemented our learning procedure in C and FORTRAN on a 4 processor Alliant supercomputer. The speed of learning can be improved considerably by computing the forward and backward sweeps for several different training samples in parallel on different processors. Further improvements can be gained by vectorizing operations and possibly assembly coding the innermost loop. Our present implementation achieves about a factor of 9 speedup over a VAX 8600, but still leaves room for further improvements (Lang[23] for example reports a speedup of a factor of 120 over a VAX11/780 for an implementation running on a Convex supercomputer). The second step taken towards improving learning time is given by a staged learning strategy. In this approach we start optimizing the network based on 3 prototypical training tokens only⁵. In this case convergence is achieved rapidly, but the network will have learned a representation that generalizes poorly to new and different patterns. Once convergence is achieved, the network is presented with approximately twice the number of tokens and learning continues until convergence.

⁴Note that in the experiments reported below these weight changes were actually carried out each time the error derivatives from all training samples had been computed[5].

⁵Note that for optimal learning, the training data is presented by always alternating tokens for each class. Hence we start the network off by presenting 3 tokens, one for each class.

Fig.3 shows the progress during a typical learning run. The measured error is $1/2$ the squared error of all the output units, normalized for the number of training tokens. In this run the number of training tokens used were 3,692,499 and 780. As can be seen from Fig.3, the error briefly jumps up every time more variability is introduced by way of more training data. The network is then forced to improve its representation to discover clues that generalize better and to deemphasize those that turn out to be merely irrelevant idiosyncracies of a limited sample set. Using the full training set of 780 tokens this particular run was continued until iteration 35,000 (Fig.3 shows the learning curve only up to 15,000 iterations). With this full training set small learning steps have to be taken and learning progresses slowly. In this case a step size of 0.002 and a momentum[5] of 0.1 was used. The staged learning approach was found to be useful to move the weights of the network rapidly into the neighborhood of a reasonable solution, before the rather slow fine tuning over all training tokens begins.

Despite these speedups, learning runs still take in the order of several days. A number of programming tricks[23] as well as modifications to the learning procedure[25] are not implemented yet and could yield another factor of 10 or more in learning time reduction. It is important to note, however, that the amount of computation considered here is necessary *only for learning* of a TDNN and *not for recognition*. Recognition can easily be performed in better than real time on a workstation or personal computer. The simple structure makes TDNNs also well suited for standardized VLSI-implementation. The detailed knowledge could be learned "off-line" using substantial computing power and then downloaded in the form of weights onto a real-time production network.

3 Hidden Markov Models (HMM)

As an alternative recognition approach we have implemented several Hidden Markov Models (HMM) aimed at phoneme recognition. HMMs are currently the most successful and promising approach [26,27,28] in speech recognition as they have been successfully applied to the whole spectrum of recognition tasks. Excellent performance was achieved at all levels from the phonemic level[29,30,31,32] to word recognition[33,28] and to continuous speech recognition[34]. HMMs' success is partially due to their ability to cope with the variability in speech by means of stochastic modeling. In the following sections, we describe the HMMs developed in our laboratory. They were aimed at phoneme recognition, more specifically the voiced stops "B", "D" and "G". Several experiments with variations on these models are described elsewhere[18].

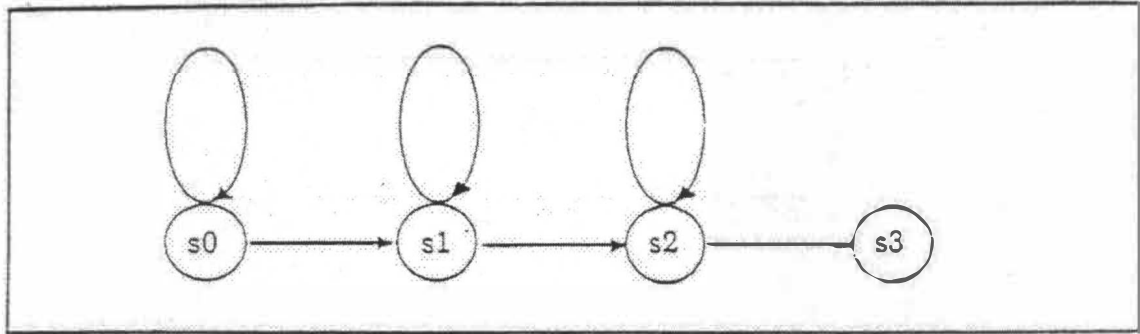


Figure 4: Hidden Markov Model

3.1 An HMM for Phoneme Recognition

The acoustic front end for Hidden Markov Modeling is typically a vector quantizer that classifies sequences of short-time spectra. Such a representation was chosen as it is highly effective for HMM-based recognizers[34].

Input speech was sampled at 12kHz, preemphasized by $(1 - 0.97z^{-1})$ and windowed using a 256-point Hamming window every 3 msec. Then a 12-order LPC analysis was carried out. A codebook of 256 LPC spectrum envelopes was generated from 216 phonetically balanced words. The Weighted Likelihood Ratio[35,36] augmented with power values (PWLR)[37,36] was used as LPC distance measure for vector quantization.

A typical HMM was adopted in this paper as shown in Fig.4. It has four states and six transitions.

3.2 Learning in an HMM

The HMM probability values were trained using vector sequences of phonemes according to the forward-backward algorithm[26]. The vector sequences for "B", "D" and "G" include a consonant part and five frames of the following vowel. This is to model important transient informations, such as formant movement and has lead to improvements over context insensitive models [18].

The HMM was trained using about 250 phoneme tokens of vector sequences per speaker and phoneme (see details of the training database below). Fig.5 shows for a typical training run the average log probability normalized by the

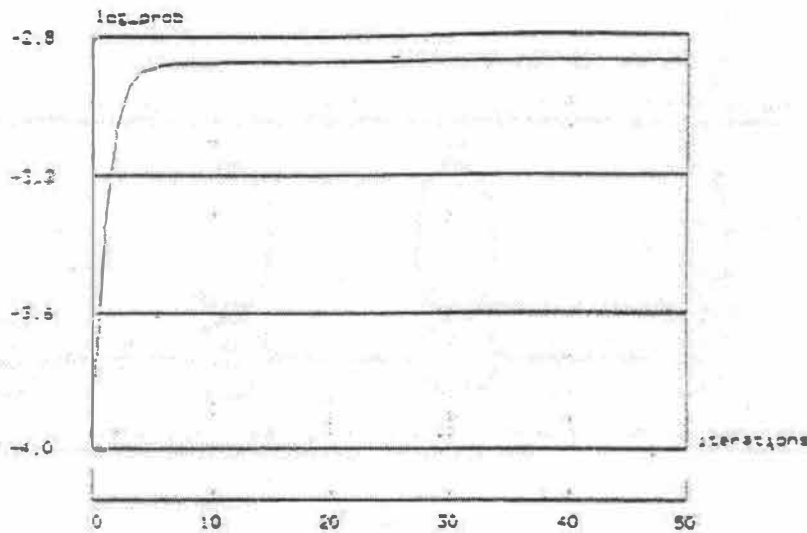


Figure 5: Learning in a Hidden Markov Model

number of frames. Training was continued until the increase of the average log probability between iterations became less than 2×10^{-3} .

Typically, about 10 to 20 learning iterations are required for 256 tokens. A training run takes about one hour on a VAX 8700. Floor values were set on the output probabilities to avoid errors caused by zero-probabilities. We have experimented with composite models, which were trained using a combination of context-independent and context-dependent probability values as suggested by Schwartz et al.[29,30]. In our case, no significant improvements were attained.

4 Recognition Experiments

We now turn to an experimental evaluation of the two techniques described in the previous sections. To provide a good framework for comparison, the same experimental conditions were given to both methods. For both, the same training data was used and both were tested on the same testing database as described below.

4.1 Experimental Conditions

For performance evaluation, we have used a large vocabulary database of 5240 common Japanese words[38]. These words were uttered in isolation by three male native Japanese speakers (MAU, MHT and MNM, all professional announcers). All utterances were recorded in a sound proof booth and digitized

at a 12 kHz sampling rate. The database was then split into a training set (the even numbered files) and a testing set (the odd numbered files). Both the training and the testing data, therefore, consisted of 2620 utterances each, from which the actual phonetic tokens were extracted.

The phoneme recognition task chosen for this experiment was the recognition of the voiced stops, i.e., the phonemes "B", "D" and "G". The actual tokens were extracted from the utterances using manually selected acoustic-phonetic labels provided with the database[38]. For speaker MAU, for example, a total of 219 "B"s, 203 "D"s and 260 "G"s were extracted from the training and 227 "B"s, 179 "D"s and 252 "G"s from the testing data. Both recognition schemes, the TDNNs and the HMMs, were trained and tested speaker-dependently. Thus in both cases, separate networks were trained for each speaker.

In our database, no preselection of tokens was performed. All tokens labeled as one of the three voiced stops were included. It is important to note, that since the consonant tokens were extracted from entire utterances and *not* read in isolation, a significant amount of phonetic variability exists. Foremost, there is the variability introduced by the phonetic context out of which a token is extracted. The actual signal of a "BA" will therefore look significantly different from a "BI" and so on. Second, the position of a phonemic token within the utterance introduces additional variability. In Japanese, for example, a "G" is nasalized, when it occurs embedded in an utterance, but not in utterance initial position. Both of our recognition algorithms are only given the phonemic identity of a token and must find their own ways of representing the fine variations of speech. Since recognition results based on the training data are not meaningful⁶, we report in the following only the results from open testing, i.e., from performance evaluation over the separate testing data set.

4.2 Results

Table 1 shows the results from the recognition experiments described above. As can be seen, for all three speakers, the TDNN yields considerable performance improvements over our HMM. Averaged over all three speakers, the error rate is reduced from 6.3% to 1.5%, a more than four fold reduction in error.

Fig.6 through Fig.11 show scatter plots of the recognition outcome for the test data for speaker MAU, using the HMM and the TDNN. For the HMM (see Fig.6 through Fig.8), the log probability of the next best matching *incorrect* token is plotted against the log probability⁷ of the correct token, e.g., "B", "D" and "G". In Fig.9 through Fig.11, the activation levels from the TDNN's output units are plotted in the same fashion. We should caution the reader that these plots are not easily comparable, as the two recognition methods have

⁶Particularly, for neural networks such results would be grossly misleading since good performance could in principle be achieved by memorization of the training patterns, rather than by generalization.

⁷normalized by number of frames

speaker	number of tokens	number of errors	recognition rate	TDNN	number of errors	recognition rate	HMM
MAU	b(227)	4	98.2	98.8	18	92.1	92.9
	d(179)	3	98.3		6	96.7	
	g(252)	1	99.6		23	90.9	
MHT	b(208)	2	99.0	99.1	8	96.2	97.2
	d(170)	0	100		3	98.2	
	g(254)	4	98.4		7	97.2	
MNM	b(216)	11	94.9	97.5	27	87.5	90.9
	d(178)	1	99.4		13	92.7	
	g(256)	4	98.4		19	92.6	

Table 1: Recognition results for three speakers over test data using TDNN and HMM

been trained in quite different ways. We present this result here to show some interesting properties of the two techniques. The most striking observation that can be made from these plots is that the output units of a TDNN have a tendency to fire with high confidence as can be seen from the cluster of dots in the lower right hand corner of the scatter plots. Most output units tend to fire strongly for the correct phonemic class and not at all for any other, a property that is encouraged by the learning procedure. One possible consequence of this is that rejection thresholds could be introduced to improve recognition performance. If one were to eliminate among speaker MAU's tokens all those whose highest activation level is less than 0.5 and those which result in two or more closely competing activations (i.e., are near the diagonal in the scatter plots), 2.6% of all tokens would be rejected, while the remaining substitution error rate would be less than 0.46%.

4.3 The Learned Internal Representations of a TDNN

Given the encouraging performance of our TDNNs, a closer look at the learned internal representation of the network is warranted. What are the properties or

abstractions that the network has learned that appear to yield a very powerful description of voiced stops? Fig.13 and Fig.12 show two typical instances of a "D" out of two different phonetic contexts ("DA" and "DO", respectively). In both cases, only the correct unit, the "D-output unit" fires strongly, despite the fact that the two input spectrograms differ considerably from each other. If we study the internal firings in these two cases we can see that the network has learned to use alternate internal representations to link variations in the sensory input to the same higher level concepts. A good example is given by the firings of the third and fourth hidden unit in the first layer above the input layer. As can be seen from Fig.13, the fourth hidden unit fires particularly strongly after vowel onset in the case of "DO", while the third unit shows stronger activation after vowel onset in the case of "DA".

Fig.14 shows the significance of these different firing patterns. Here the connection strengths for the eight moving TDNN units are shown, where white and black blobs represent positive and negative weights, respectively, and the magnitude of a weight is indicated by the size of the blob. In this figure, the time delays are displayed spatially as a 3 frame window of 16 spectral coefficients. Conceptually, the weights in this window form a moving acoustic-phonetic feature detector, that fires when the pattern for which it is specialized is encountered in the input speech. In our example, we can see that hidden unit number 4 (which was activated for "DO") has learned to fire when a falling (or rising) second formant starting at around 1600 Hz is found in the input (see filled arrow in Fig.14). As can be seen in Fig.13, this is the case for "DO" and hence the firing of hidden unit 4 after voicing onset (see row pointed to by the filled arrow in Fig.13). In the case of "DA" (see Fig.12) in turn, the second formant does not fall significantly, and hidden unit 3 (pointed to by the filled arrow in Fig.14) fires instead. From Fig.14 we can verify that TDNN-unit 3 has learned to look for a steady (or only slightly falling) second formant starting at about 1800 Hz. The connections in the second and third layer then link the different firing patterns observed in the first hidden layer into one and the same decision.

Another interesting feature can be seen in the bottom hidden unit in hidden layer number 1 (see Fig.12, Fig.13 and compare with the weights of hidden unit 1 displayed in Fig.14). This unit has learned to take on the role of finding the segment boundary of the voiced stop. It does so in reverse polarity, i.e., it is always on *except* when the vowel onset of the voiced stop is encountered (see unfilled arrow in Fig.13 and Fig.12). Indeed, the higher layer TDNN-units subsequently use this "segmenter" to base the final decision on the occurrence of the right lower features at the right point in time.

In the previous example, we have seen that the TDNN can account for variations in phonetic context. Fig.15 and Fig.16 show examples of variability caused by the relative position of a phoneme within a word. In Japanese, a "G" embedded in a word tends to be nasalized as seen in the spectrum of a "GA" in Fig.15. Fig.16 shows a word initial "GA". Despite the striking differences between these two input spectrograms, the network's internal alternate repre-

sentations manage to produce in both cases crisp output firings for the right category.

Fig. 17 and Fig. 18, finally, demonstrate the shift-invariance of the network. They show the same token "DO" of Fig. 13, misaligned by +30 msec and -30 msec, respectively. Despite the gross misalignment, (note that significant transitional information is lost by the misalignment in Fig. 18) the correct result was obtained reliably. A close look at the internal activation patterns reveals that the hidden units' feature detectors do indeed fire according to the events in the input speech, and are not negatively affected by the relative shift with respect to the input units.

Three important properties of the TDNNs have therefore been observed. First, our TDNN was able to learn without human interference meaningful linguistic abstractions such as formant tracking and segmentation. Second, we have demonstrated that it has learned to form alternate representations linking different acoustic events with the same higher level concept. In this fashion it can implement trading relations between lower level acoustic events leading to robust recognition performance. Third, we have seen that the network is shift-invariant and does not rely on precise alignment or segmentation of the input.

5 Conclusion and Summary

In this paper we have presented a Time Delay Neural Network (TDNN) approach to phoneme recognition. We have shown that this TDNN has two desirable properties related to the dynamic structure of speech. First, it can learn the temporal structure of acoustic events and the temporal relationships between such events. Second, it is translation invariant, that is, the features learned by the network are insensitive to shifts in time. Examples demonstrate that the network was indeed able to learn acoustic phonetic features, such as formant movements and segmentation, and use them effectively as internal abstractions of speech.

The TDNN presented here has two hidden layers and has the ability to learn complex non-linear decision surfaces. This could be seen from the network's ability to use alternate internal representations and trading relations among lower level acoustic-phonetic features, in order to arrive robustly at the correct final decision. Such alternate representations have been particularly useful for representing tokens that vary considerably from each other due to their different phonetic environment or their position within the original speech utterance.

Finally, we have evaluated the TDNN on the recognition of three acoustically similar phonemes, the voiced stops "B", "D" and "G". In extensive performance evaluation over testing data from three speakers, the TDNN achieved an average recognition score of 98.5 %. For comparison, we have applied various Hidden Markov Models to the same task and only been able to reach recognize 93.7 %

of the tokens correctly. We would like to note, that many variations of HMMs have been attempted and many more variations of both HMMs and TDNs are conceivable. Some of these variations could potentially lead to significant improvements over the results reported in this study. Our goal here is to present TDNs as a new and successful approach for speech recognition. Their power lies in their ability to develop shift-invariant internal abstractions of speech and use them in trading relations for making optimal decisions. This holds significant promise for speech recognition in general, as it could overcome the representational weaknesses of existing techniques when faced with uncertainty and variability in real life signals.

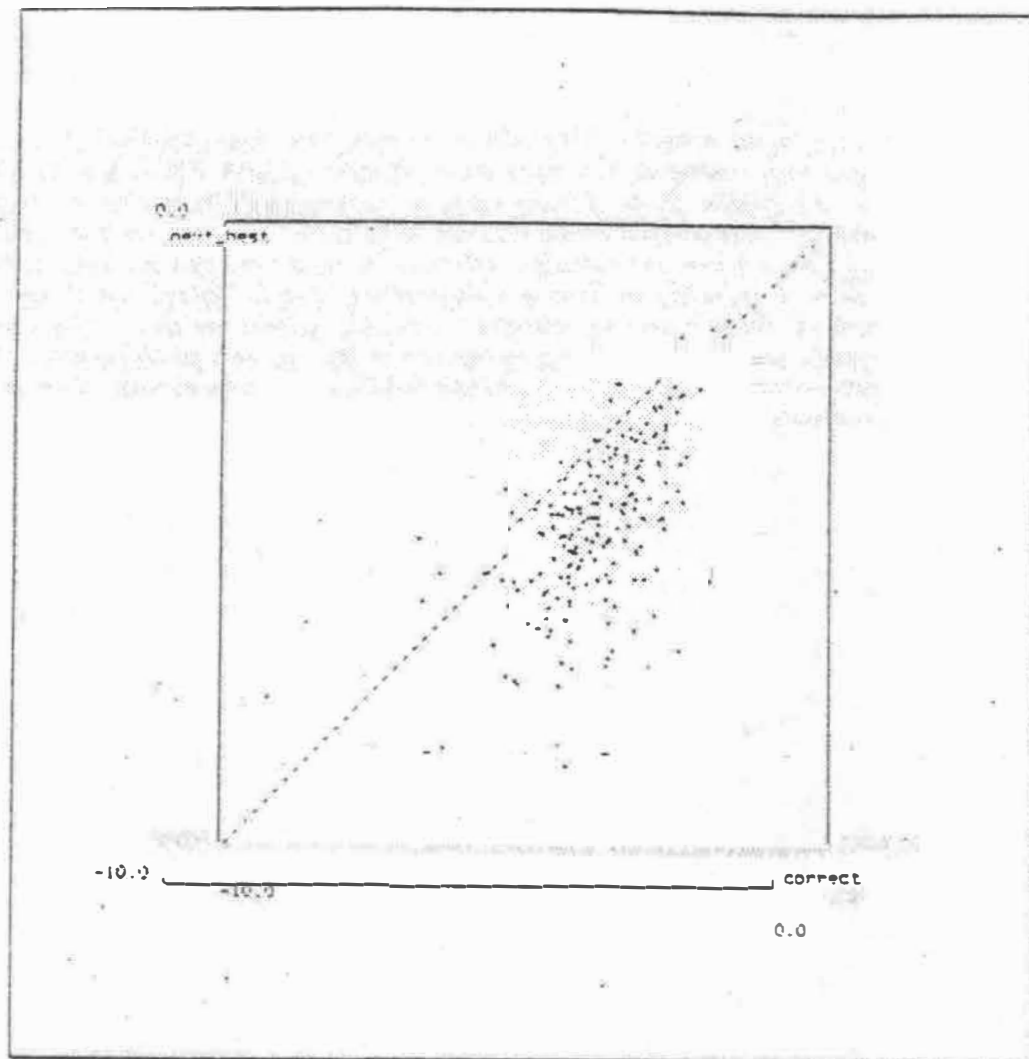


Figure 6: Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "B"s using a HMM

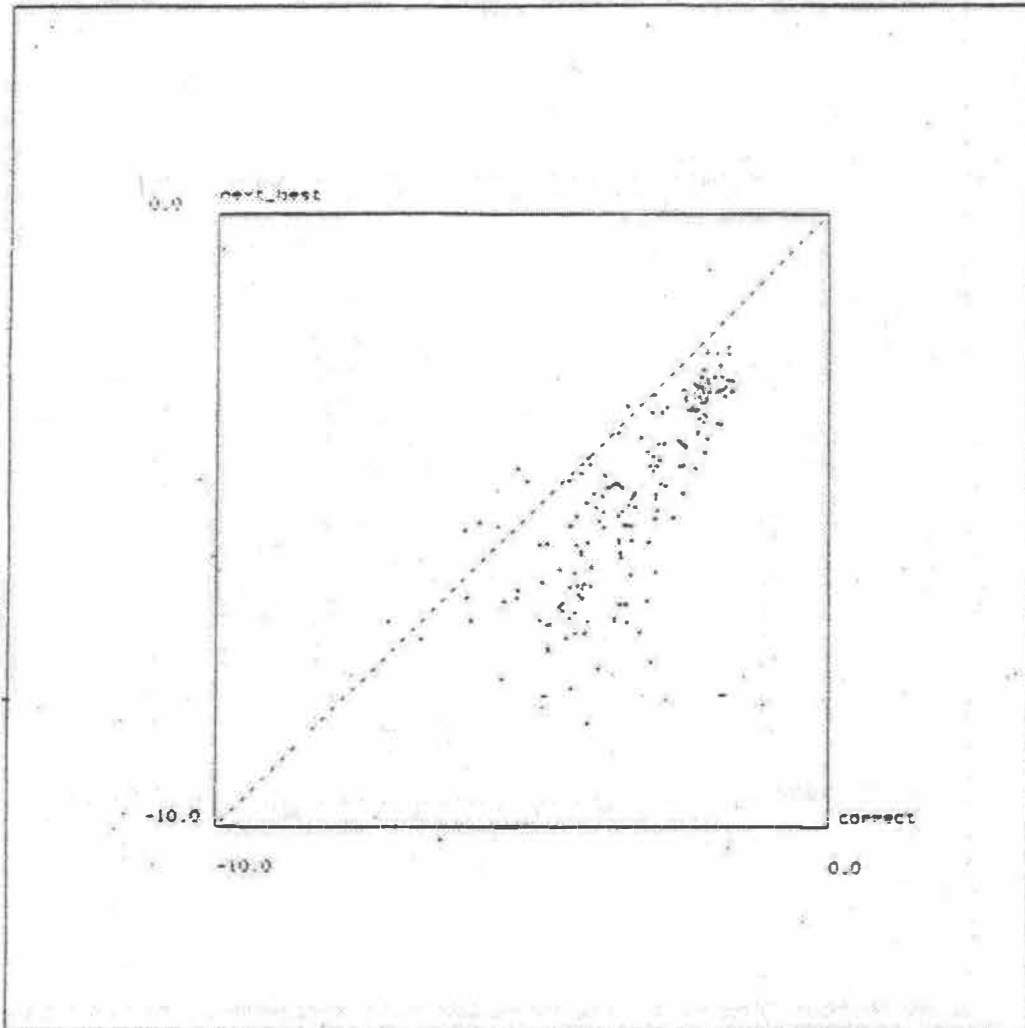


Figure 7: Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "D"s using a HMM.

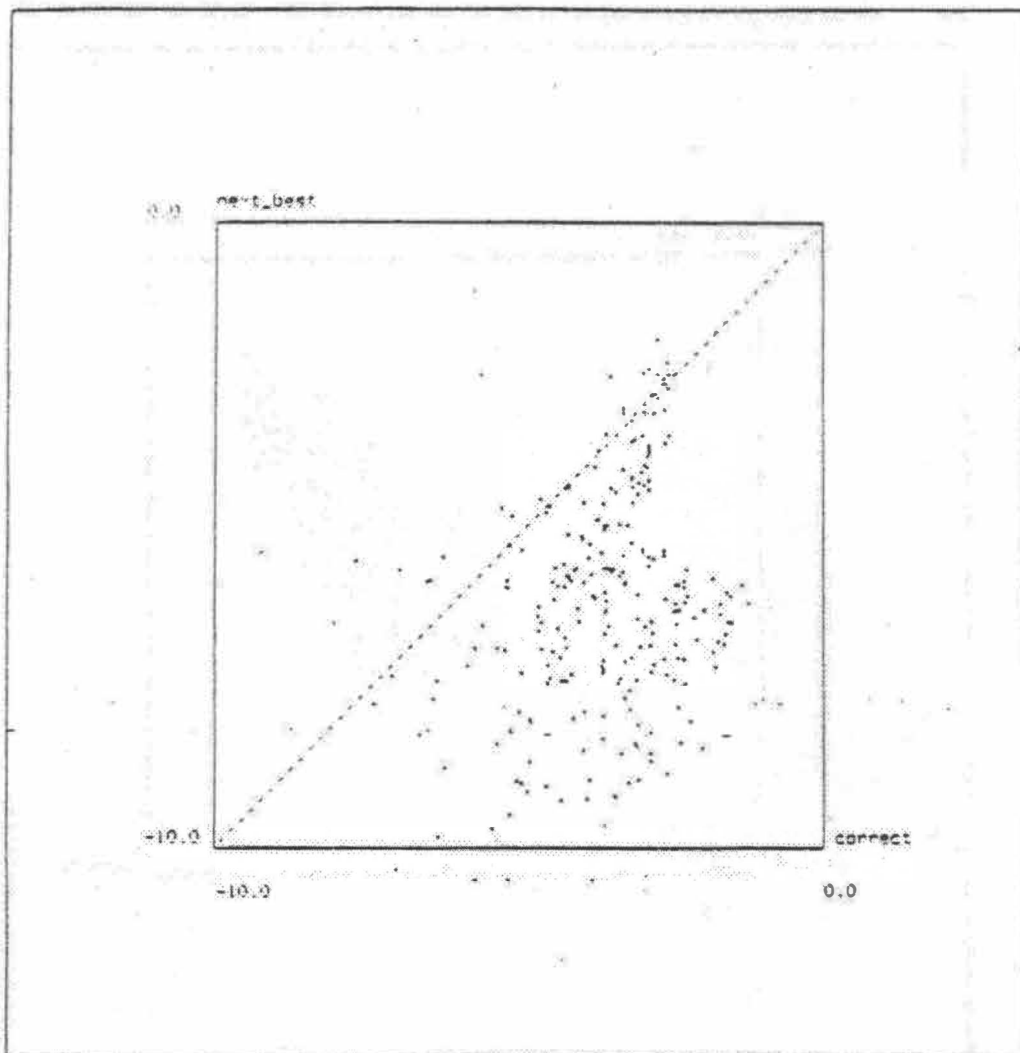


Figure 8: Scatter plot showing log probabilities for the best matching incorrect case vs. the correctly recognized "G"s using a HMM

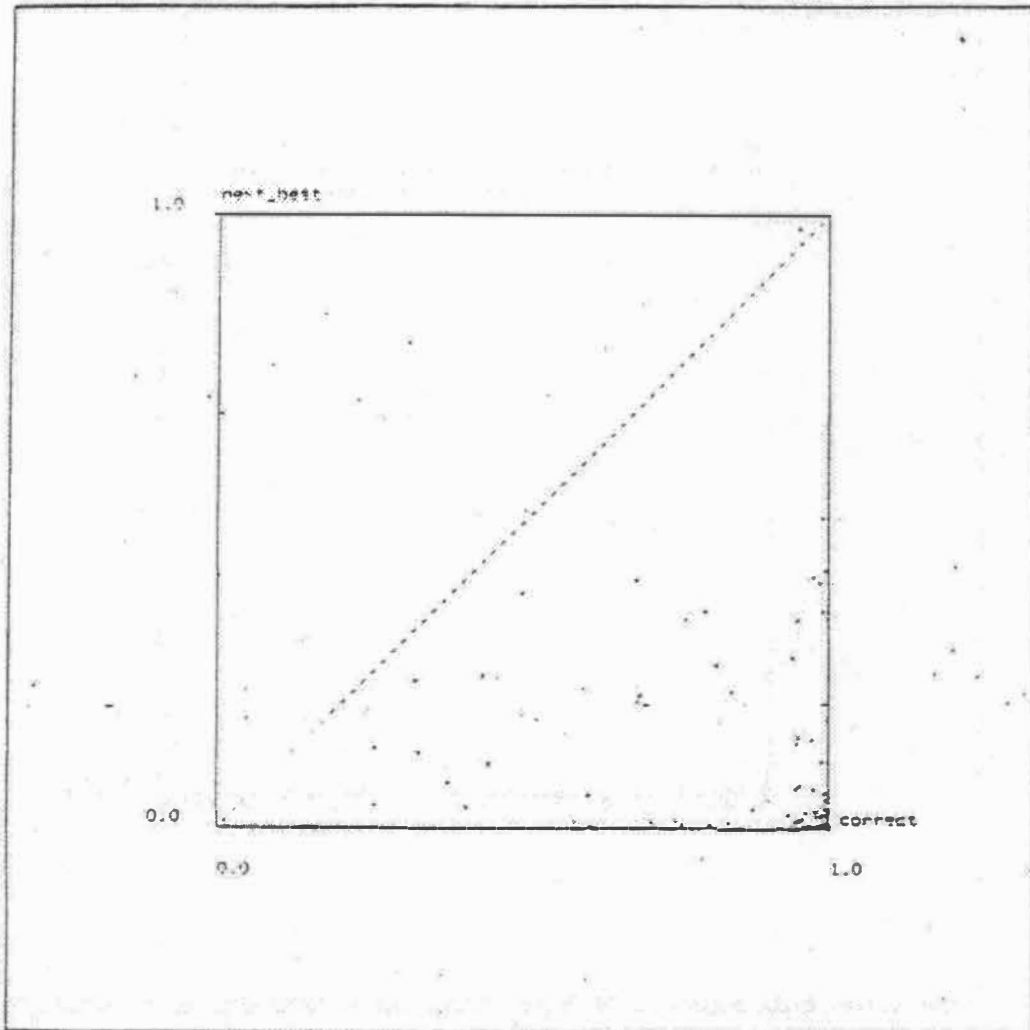


Figure 9: Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "B"s using a TDNN

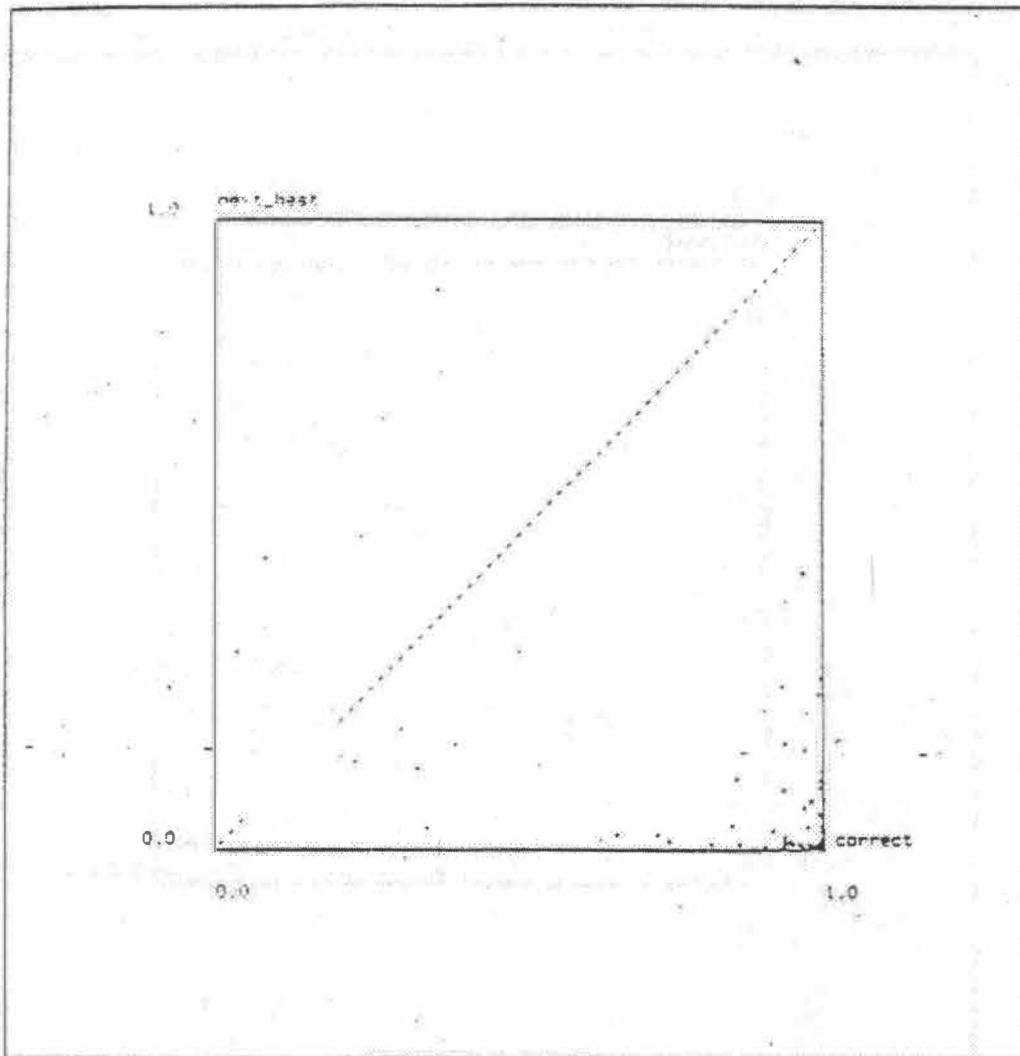


Figure 10: Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "D"s using a TDNN

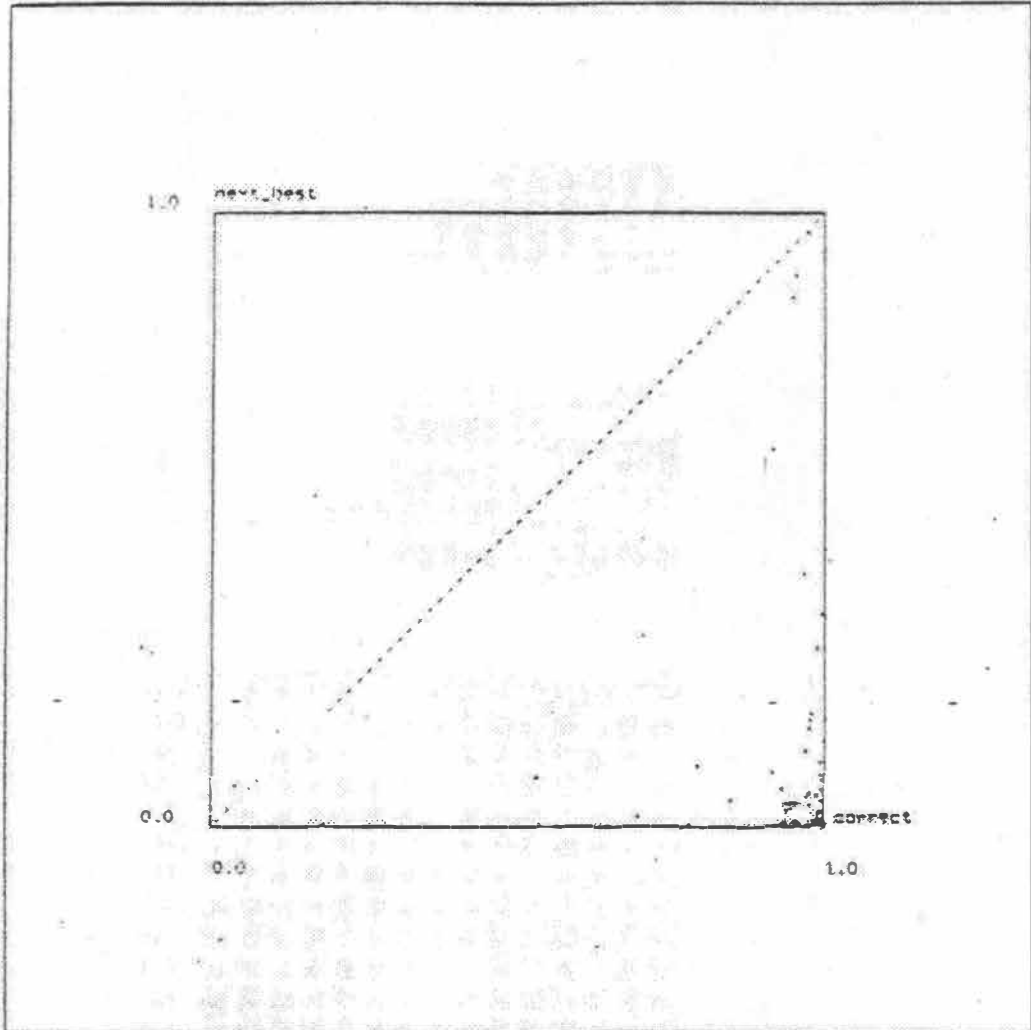


Figure 11: Scatter plot showing activation levels for the best matching incorrect case vs. the correctly recognized "G"s using a TDNN

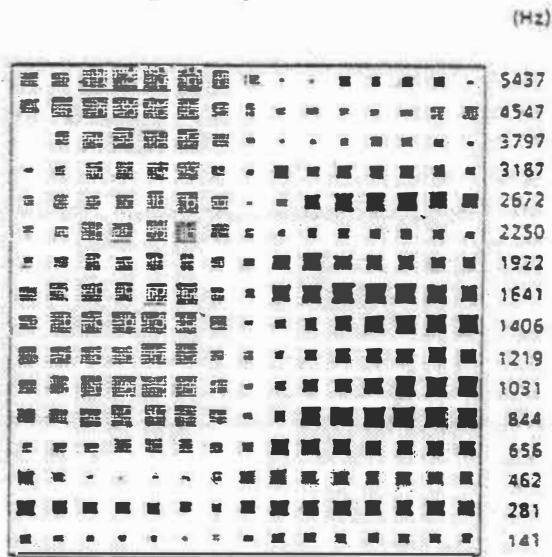
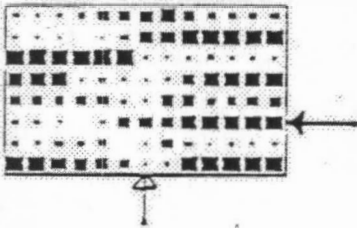
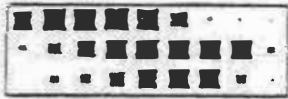


Figure 12: TDNN Activation patterns for "DA"

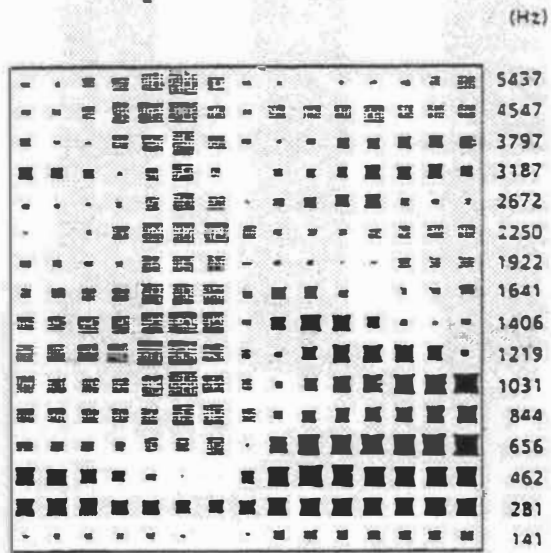
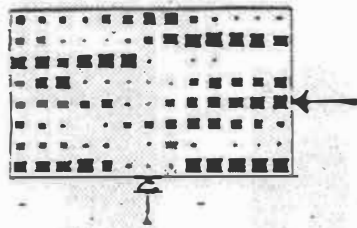
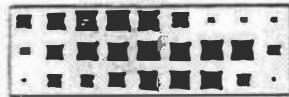
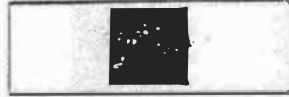


Figure 13: TDNN Activation patterns for "DO"

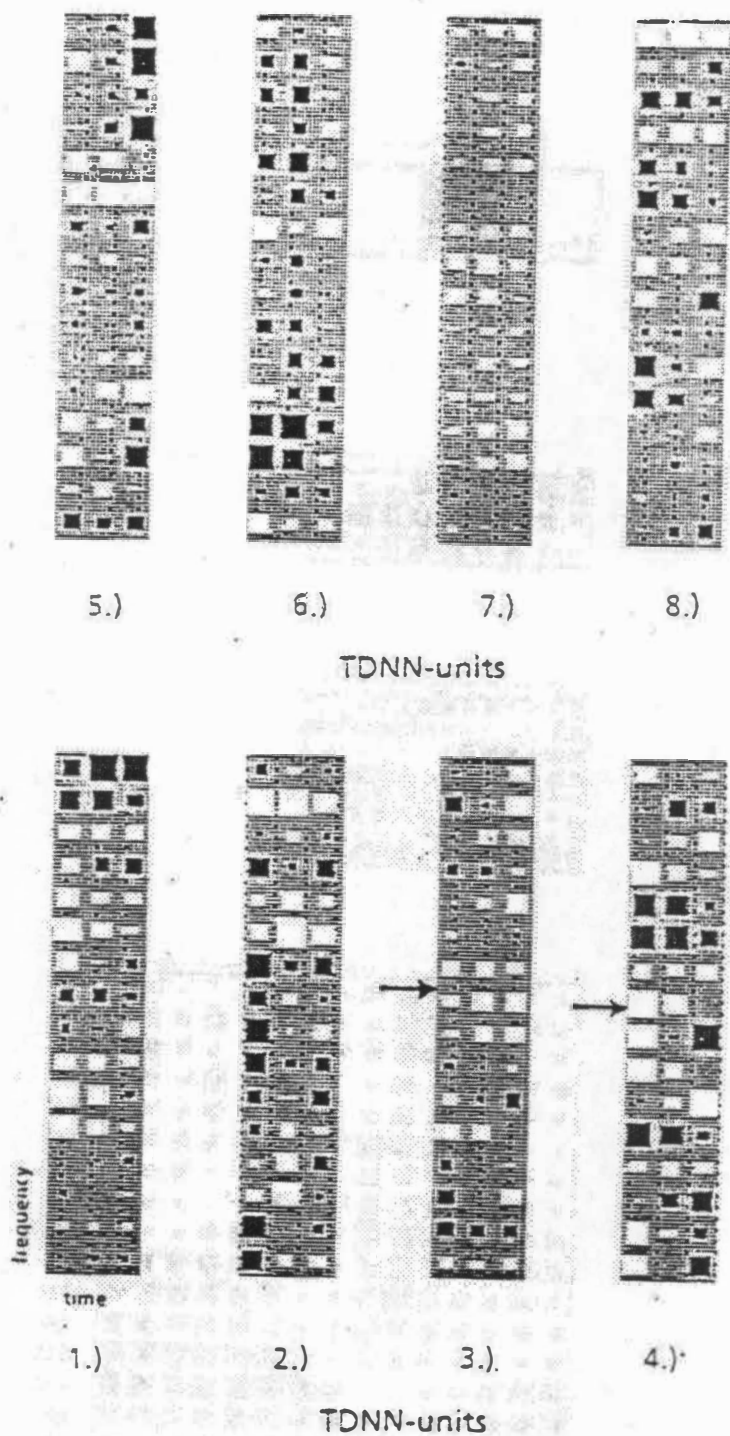
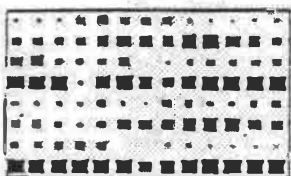


Figure 14: Weights on connections from 16 coefficients over 3 time frames to each of the 8 hidden units in the first layer



(Hz)

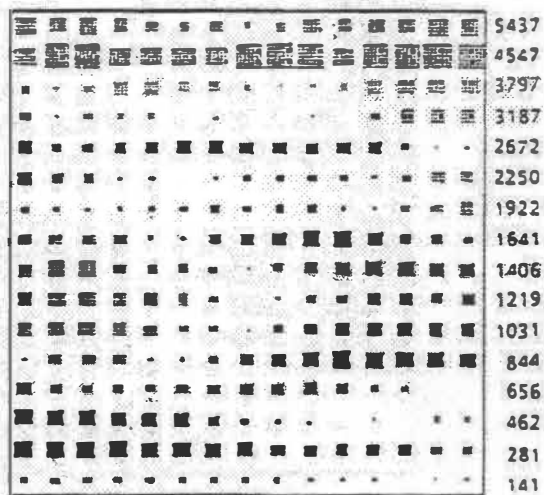


Figure 15: TDNN Activation patterns for "GA" embedded in an utterance

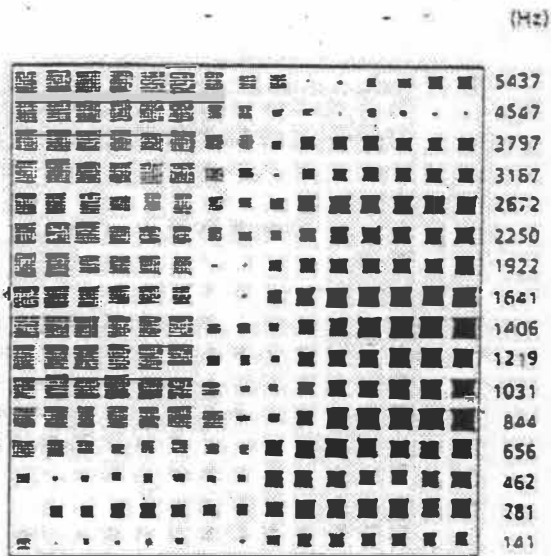
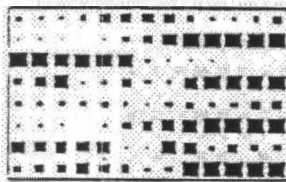
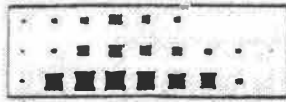
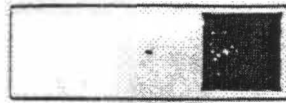
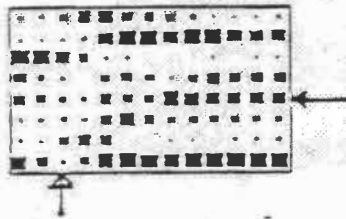


Figure 16: TDNN Activation patterns for "GA" in utterance initial position



(M2)

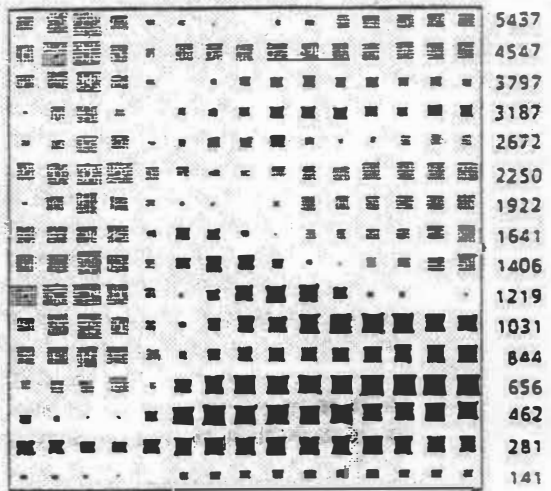


Figure 17: TDNN Activation patterns for "DO" misaligned by +30 msec

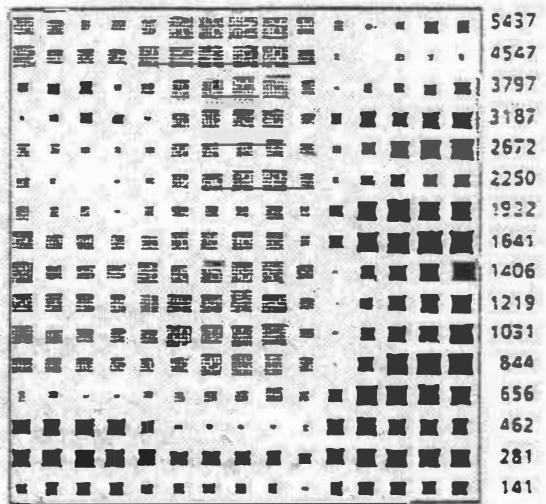
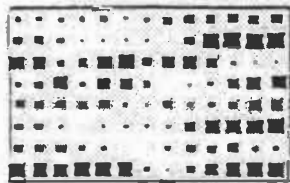
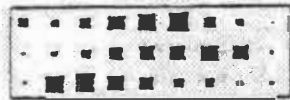
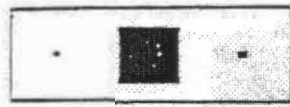


Figure 1.3: TDNN Activation patterns for "DO" misaligned by -30 msec

References

- [1] D.E. Rumelhart and J.L. McClelland. *Parallel Distributed Processing; Explorations in the Microstructure of Cognition*. Volume I and II, MIT Press, Cambridge, MA, 1986.
- [2] R.P. Lippmann. An introduction to computing with neural nets. *IEEE ASSP Magazine*, 4-22, April 1987.
- [3] D.C. Plaut, S.J. Nowlan, and G.E. Hinton. *Experiments on Learning by Back Propagation*. Technical Report CMU-CS-86-126, Carnegie-Mellon University, June 1986.
- [4] T.J. Sejnowski and C.R. Rosenberg. *NETalk: A Parallel Network that Learns to Read Aloud*. Technical Report JHU/EECS-86/01, Johns Hopkins University, June 1986.
- [5] D.E. Rumelhart, G.E. Hinton, and R.J. Williams. Learning representations by back-propagating errors. *Nature*, 323:533-536, October 1986.
- [6] Huang W.Y. and Lippmann R.P. Comparison between neural net and conventional classifiers. In *IEEE International Conference on Neural Networks*, June 1987.
- [7] J.L. McClelland and J.L. Elman. *Interactive Processes in Speech Perception: The TRACE Model*, chapter 15, pages 58-121. MIT Press, Cambridge, MA, 1986.
- [8] S.M. Peeling, R.K. Moore, and M.J. Tomlinson. The multi-layer perceptron as a tool for speech pattern processing research. In *Proceedings ICA Autumn Conference on Speech and Hearing*, 1986.
- [9] H. Bourlard and C.J. Wellekens. Multilayer perceptrons and automatic speech recognition. In *IEEE International Conference on Neural Networks*, June 1987.
- [10] B. Gold, Lippmann R.P., and M.L. Malpass. Some neural net recognition results on isolated words. In *IEEE International Conference on Neural Networks*, June 1987.
- [11] Lippmann R.P. and B. Gold. Neural-net classifiers useful for speech recognition. In *IEEE International Conference on Neural Networks*, June 1987.
- [12] D.J. Burr. A neural network digit recognizer. In *IEEE International Conference on Systems, Man, and Cybernetics*, October 1986.
- [13] D. Lubensky. Learning spectral-temporal dependencies using connectionist networks. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, April 1988. to be published.

- [14] R.L. Watrous and L. Shastri. Learning phonetic features using connectionist networks: an experiment in speech recognition. In *IEEE International Conference on Neural Networks*, June 1987.
- [15] R.W. Prager, T.D. Harrison, and F. Fallside. Boltzmann machines for speech recognition. *Computer, Speech and Language*, 3-27, March 1986.
- [16] J.L. Elman and D. Zipser. *Learning the Hidden Structure of Speech*. Technical Report, University of California, San Diego, February 1987.
- [17] R.L. Watrous, L. Shastri, and A.H. Waibel. Learned phonetic discrimination using connectionist networks. In *European Conference on Speech Technology*, pages 377-380, Edinburgh, September 1987.
- [18] T. Hanazawa, T. Kawabata, and K. Shikano. Discrimination of Japanese voiced stops using hidden markov model. In *Conference of the Acoustical Society of Japan*, pages 19-20, October 1987. (in Japanese).
- [19] K. Fukushima, S. Miyake, and T. Ito. Neocognitron: a neural network model for a mechanism of visual pattern recognition. *IEEE Transactions on System, Man, and Cybernetics*, SMC-13(5):826-834, September/October 1983.
- [20] D.E. Rumelhart, G.E. Hinton, and R.J. Williams. *Learning Internal Representations by Error Propagation*, chapter 8, pages 318-362. MIT Press, Cambridge, MA, 1986.
- [21] A. Waibel and B. Yegnanarayana. *Comparative Study of Nonlinear Time Warping Techniques in Isolated Word Speech Recognition Systems*. Technical Report, Carnegie-Mellon University, June 1981.
- [22] S. Makino and K. Kido. Phoneme recognition using time spectrum pattern. *Speech Communication*, 225-237, June 1986.
- [23] K. Lang. Connectionist speech recognition. July 1987. PhD thesis proposal, Carnegie-Mellon University.
- [24] G.E. Hinton. Connectionist learning procedures. *Artificial Intelligence*, 1987. (in press).
- [25] M.A. Franzini. Speech recognition with back propagation. In *Ninth Annual Conference of the IEEE/Engineering in Medicine and Biology Society*, November 1987.
- [26] F. Jelinek. Continuous speech recognition by statistical methods. *Proceedings of the IEEE*, 64(4):532-556, April 1976.
- [27] J. K. Baker. *Stochastic Modeling as a Means of Automatic Speech Recognition*. PhD thesis, Carnegie-Mellon University, April 1975.

- [28] L. R. Bahi, S. K. Das, P. V. de Souza, F. Jelinek, S. Katz, R. L. Mercer, and M. A. Picheny. Some experiments with large-vocabulary isolated-word sentence recognition. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, April 1984.
- [29] R. Schwartz, Y. Chow, O. Kimball, S. Roucos, M. Krasner, and J. Makhoul. Context-dependent modeling for acoustic-phonetic recognition of continuous speech. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, April 1985.
- [30] A.-M. Derouault. Context-dependent phonetic markov models for large vocabulary speech recognition. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 360-3, April 1987.
- [31] K.F. Lee and H. Hsiao-Wuen. Speaker-independent phoneme recognition using hidden markov models. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, April 1987.
- [32] P. Brown. *The Acoustic-Modeling Problem in Automatic Speech Recognition*. PhD thesis, Carnegie-Mellon University, May 1985.
- [33] L. R. Rabiner, B. H. Juang, S. E. Levinson, and M. M. Sondhi. Recognition of isolated digits using Hidden Markov Models with continuous mixture densities. *AT&T Technical Journal*, 64(6):1211-33, July-August 1985.
- [34] Y.L. Chow, M.O. Dunham, O.A. Kimball, M.A. Krasner, G.F. Kubala, J. Makhoul, S. Roucos, and R.M. Schwartz. BYBLOS: the BBN continuous speech recognition system. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 89-92, April 1987.
- [35] M. Sugiyama and K. Shikano. LPC peak weighted spectral matching measures. *Institute of Electrical and Communication Engineers of Japan*, 64-A(5):409-416, 1981. in Japanese.
- [36] K. Shikano. *Evaluation of LPC Spectral Matching Measures for Phonetic Unit Recognition*. Technical Report, Carnegie-Mellon University, May 1985.
- [37] K. Aikawa and K. Shikano. Spoken word recognition using vector quantization in power-spectrum vector space. *Institute of Electrical and Communication Engineers of Japan*, 68-D(3), March 1985. in Japanese.
- [38] Y. Sagisaka, K. Takeda, S. Katagiri, and H. Kuwabara. *Japanese Speech Database with Fine Acoustic-Phonetic Transcriptions*. Technical Report, ATR Interpreting Telephony Research Laboratories, May 1987.

TR-I-0006

Phoneme Recognition
Using Time-Delay Neural Networks

時間遅延神経回路網による音韻認識

A. Waibel, T. Hanazawa, G. Hinton,
K. Shikano and K. Lang

ワイバル、花沢利行、ヒントン、鹿野清宏、ラング

Advanced Telecommunications Research Institute International

ATR

ATR TECHNICAL REPORT

(株)エイ・ティ・アール自動翻訳電話研究所

ATR Interpreting Telephony