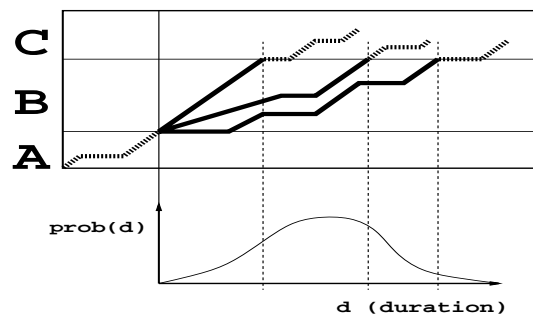
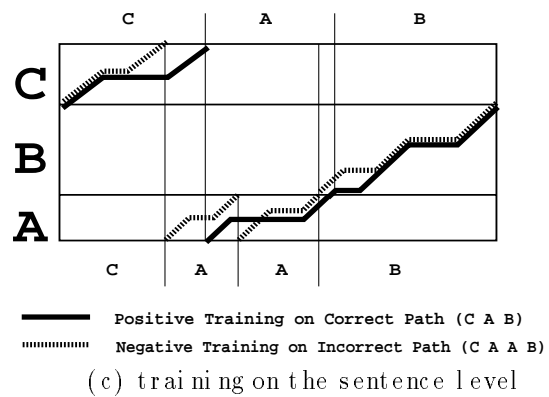


(a) alignment across word boundaries



(b) duration dependent word penalties



(c) training on the sentence level

Figure 2: Various techniques to improve sentence level recognition performance

Speaker Dependent (CMU Alph Data)			
500/2500 train, 100/500 crossvalidation, 400/2000 test sentences/words			
speaker	SPHINX [HF91]	MS-TDNN [HF91]	our MS-TDNN
njnt	96.0	97.5	98.5
mlbs	83.9	89.7	91.1
naem	-	-	94.6
fcaw	-	-	98.8
flgt	-	-	86.9
fee	-	-	91.0

Speaker Independent (Resource Management Spell-Mode)			
109 (ca. 11000) train, 11 (ca. 900) test speaker (words).			
	SPHINX [HF92]	our MS-TDNN	
	+ Senone		gender specific
	88.7	90.4	90.8
			92.0

Table 1: Word accuracy (in % on the test sets) on speaker dependent and speaker independent connected letter tasks.

In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, Toronto, Ontario, Canada, May 1991. IEEE

[FC90] M. Fanty and R. Cole. Spoken letter recognition. In *Proceedings of the Neural Information Processing Systems Conference NIPS*, Denver, November 1990.

[Haf92] P. Haffner. Connectionist Word-Level Classification in Speech Recognition. In *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 1992.

[HF91] P. Haffner, M. Franzini, and A. Waibel. Integrating Time Alignment and Neural Networks for High Performance Continuous Speech Recognition. In *Proc. Int. Conf. on Acoustics, Speech, and Signal Processing*. IEEE, 1991.

[HB92] M.-Y. Hwang and X. Huang. Subphonetic Modeling with Markov States - Senone. In *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 133 - 137. IEEE, 1992.

[HW90] J. Hampshire and A. Waibel. A Novel Objective Function for Improved Phoneme Recognition Using Time Delay Neural Networks. *IEEE Transactions on Neural Networks*, June 1990.

P. Haffner and A. Waibel. Multi-state Time Delay Neural Networks for Continuous Speech Recognition. In *NIPS(4)*. Morgan Kaufman, 1992.

Waibel. Multi-Speaker/Speaker-Independent Architectures for Time Delay Neural Network. In *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 1993.

and R. A. Cole. Speaker-independent Phonetic Recognition of English Letters. In *Proceedings of the IJCNN*

Dynamic Programming Algorithm for Continuous Speech Recognition. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 1984.

K. Lang. Phoneme Recognition. In *IEEE Transactions on*

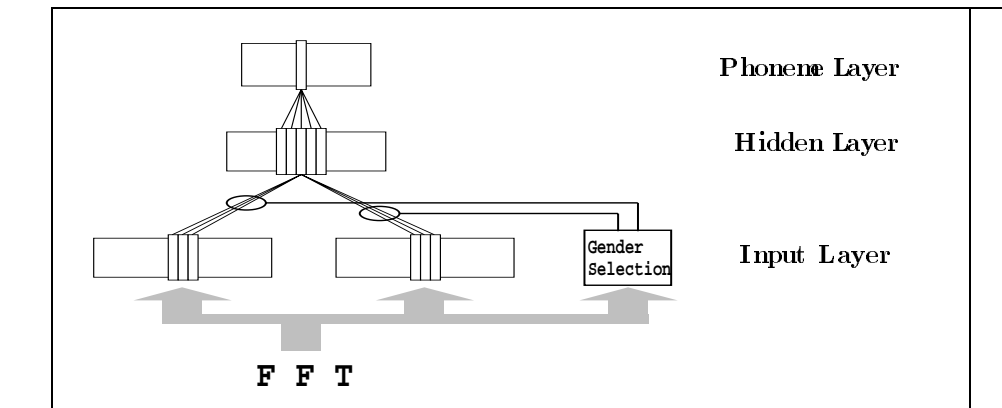


Figure 4: A network architecture with gender-specific and shared connections. Only the front-end TDNN is shown.

ing, cross-validation and test set, respectively. The DARPA Resource Management Spell-Mode Data were used for speaker independent testing. This data base contains about 1700 sentences, spelled by 85 male and 35 female speakers. The speech of 7 male and 4 female speakers was set aside for the test set, one sentence from all 109 and all sentences from 6 training speakers were used for crossvalidation. Table 1 summarizes our results. With the help of the training technique above we were able to outperform previously reported [HFW1] speaker independent results as well as the HMM based SPHINX System

5 SUMMARY AND FUTURE WORK

We have presented a connectionist speech recognition system for connected letter recognition. New training techniques for the 15-level recognition enabled our MS-TDNN to outperform other systems of this kind as well as a state-of-the-art HMM based system. In the future, gender-specific subnets, we are experimenting with “internal speaker models” for a more robust system. In the future we will also experiment with other techniques.

Acknowledgements

The authors gratefully acknowledge the support of DARPA. We wish to thank M. Nair for keeping the system running and the ideas presented here.

References

[CFG91]
In

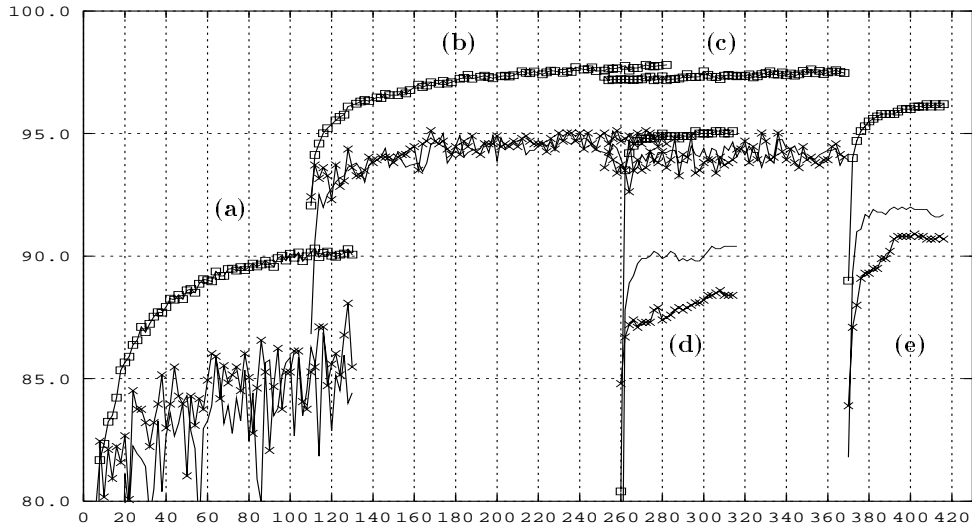


Figure 3: Learning curves (a = bootstrapping, b, c = word level (excerpted words), d, e = sentence level training (continuous speech)) on the training (\square), cross validation (-) and test set (x) for the speaker-independent RMSpell-Møde data.

3 G E N D E R S P E C I F I C S U B N E T S

A straightforward approach to building a more specialized system is to build two entirely individual networks for male and female speakers. If the gender of a speaker is known, during testing it is possible to use the “gender identification network”, which is simply a neural network with input units representing male and female speakers. This network classifies the speaker’s gender. The generalized network improved the accuracy (see table 1) to 91.3%. However, the network worked even better when the connections at the input layer were worked even better. The classification

in the same way as the phoneme boundaries within a word. Figure 2(a) shows an example in which the word to recognize is surrounded by a silence and a 'B', thus the left and right context (for all words to be recognized) is the phoneme 'sil' and 'b', respectively. The gray shaded area indicates the extension necessary to the DIW alignment. The diagram shows how a new boundary for the beginning of the word 'A' is found. As indicated in figure 3, this technique improves continuous recognition significantly, but it doesn't help for excerpted words.

2.2 WORD DURATION DEPENDENT PENALIZING OF INSERTION AND DELETION ERRORS

In "continuous testing mode", instead of looking at word units the well-known "One Stage DIW algorithm" [Ney84] is used to find an optimal path through an unspecified sequence of words. The short and confusable English letters cause many word insertion and deletion errors, such as "T E" vs. "T" or "O" vs. "O O", therefore proper duration modeling is essential.

As suggested in [HW92], minimum phoneme duration can be enforced by "duration duplication". In addition, we are modeling a duration and word dependent penalty

$$Pen_w(d) = \log(k + prob_w(d)),$$

where the pdf $prob_w(d)$ is approximated by the training data and k is a small constant to avoid zero probabilities.

This penalty is added to the accumulated score AS of the search path whenever it crosses the boundary of a word w in the alignment algorithm, as indicated in figure 2(b). The ratio λ

influence of the duration penalty, is another parameter.

There is no straightforward mathematical derivation of the "weight" λ_w to the insertion and deletion

errors. It is determined by gradient descent, which converges to a local minimum, i.e. we are trying to minimize the error.

2.3 ERROR ANALYSIS

Usually the MS-TDNN is trained to recognize continuously spoken sentences.

The training on the sentence "C A B", in which the alignment is

shown in figure 4, is done by

aligning the word boundaries

with the phoneme boundaries.

copied from the Phoneme Layer into the word models of the DIW Layer, where an optimal alignment path is found for each word. The activations along these paths are then collected in the word output units. All units in the DIW and Word Layer are linear and have no biases. 15 (25 to 100) hidden units per frame were used for speaker-dependent (-independent) experiments, the entire 26 letter network has approximately 5200 (8600 to 34500) parameters.

Training starts with "bootstrapping", during which only the front-end TDNN is used with fixed phoneme boundaries as targets. In a second phase, training is formed with word level targets. Phoneme boundaries are freely aligned with word boundaries in the DIW Layer. The error derivatives are backpropagated through the word units through the alignment path and the front-end TDNN. The choice of sensible objective functions is of great importance. For example, given (y_1, \dots, y_n) the output and $T = (t_1, \dots, t_n)$ the target at the phoneme level (bootstrapping), then the error function is $\sum_{i=1}^n |y_i - t_i|$, representing the correct phoneme. For word level training, see why the standard *Mean Squared Error* is not suitable for "1-out-of- n " coding. The error function for a target $(1, 0, \dots, 0)$ is $\sum_{i=1}^n |y_i - t_i|^2$, which is not suitable for this task.

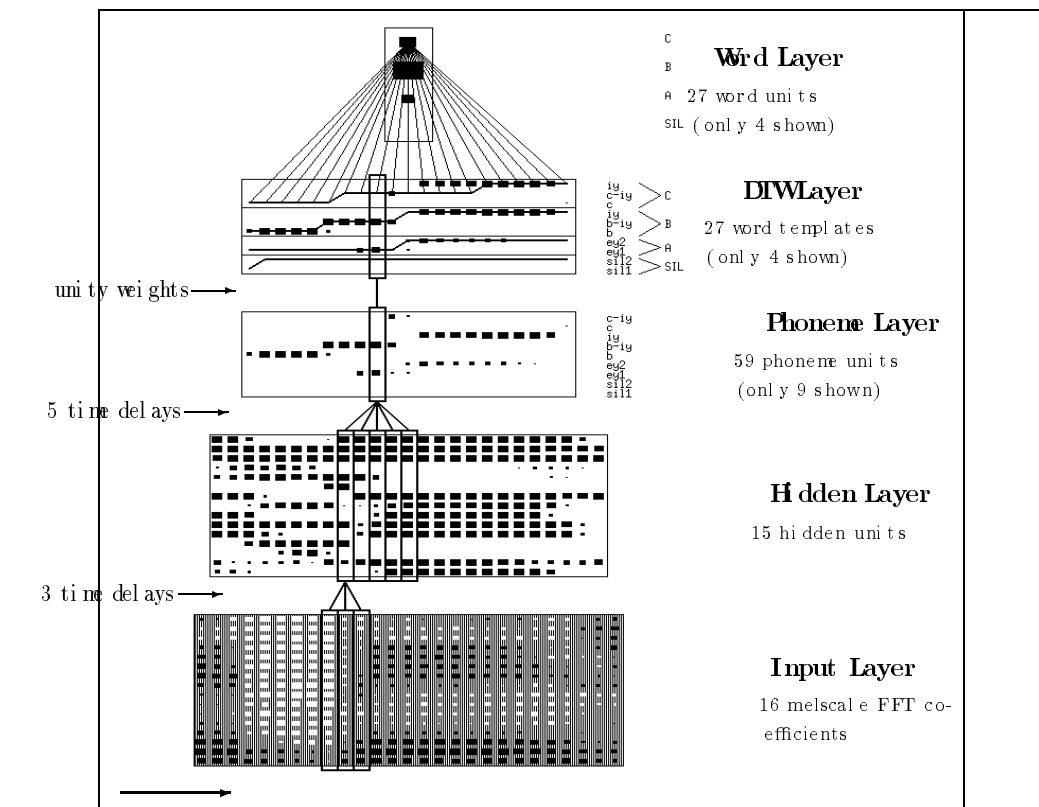


Figure 1: The MS-TDNN recognizing the excerpted word 'B'. Only the activations for the words 'SIL', 'A', 'B', and 'C' are shown.

classified by another network. In this paper, we present the MS-TDNN as a connectionist speech recognition system for connected letter recognition. After describing the baseline architecture, training techniques aimed at improving sentence level performance and architectures with gender-specific subnets are introduced.

Baseline Architecture. Time Delay Neural Networks (TDNNs) can combine the robustness and discriminative power of Neural Nets with a time-shift architecture to form high accuracy phoneme classifiers [WH⁺89].

MS-TDNN (MS-TDNN) [HFV⁺91, Haf92, HV⁺92], an extension of the TDNN, is used for classifying words (represented as sequences of phonemes) by a dynamic time warping (DTW) procedure (DIW) into the TDNN architecture.

In an MS-TDNN in the process of recognizing the excerpted word 'B',

16 mel-scale FFT coefficients at a 10-msec frame rate are processed by

to compute a score for each phoneme (static) in the "Phoneme Layer". In

the "Phoneme Layer", the word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

The word 'B' is modeled by a sequence of phonemes.

Connected Letter Recognition with a Multi-State Time Delay Neural Network

Hermann Hild and **Alex Wilbel**
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213-3891, USA

Abstract

The Multi-State Time Delay Neural Network (MS-TDNN) integrates a nonlinear time alignment procedure (DTW) and the high-accuracy phoneme spotting capabilities of a TDNN into a connectionist speech recognition system with word-level classification error backpropagation. We present an MS-TDNN for recognizing continuously spelled letters, a task characterized by a highly confusable vocabulary. Our MS-TDNN achieves word accuracy on speaker dependent/independent data, forming previously reported results on the task. We propose training techniques aimed at improving performance, including free alignment, ratio modeling and error backpropagation. We show that the MS-TDNN performs better than the word level. Results are presented and discussed on a subset of the data.

1 INTRODUCTION

The recognition