

SPEAKER-INDEPENDENT CONNECTED LETTER RECOGNITION WITH A MULTI - STATE TIME DELAY NEURAL NETWORK

Hermann Hild and Alex Waibel

Universität Karlsruhe, Germany
Carnegie Mellon University, Pittsburgh, USA

ABSTRACT

We present a Multi-State Time Delay Neural Network (MS-TDNN) for speaker-independent, connected letter recognition. Our MS-TDNN achieves 98.5/92.0% word accuracy on speaker dependent/independent English letter tasks [7, 8]. In this paper we will summarize several techniques to improve (a) continuous recognition performance, such as sentence level training, and (b) phonetic modeling, such as network architectures with "internal speaker models", allowing for "tuning-in" to newspeakers. We also present results on our large and still growing German Letter data base, containing over 40.000 continuously spelled by 55 speakers.

Spelled Letter Recognition, Speaker-
independence, MS-TDNN

1. INTRODUCTION

Recognition of spelled strings of letters is essential for all applications involving special vocabularies, such as names or addresses. Despite its small vocabulary, the task is quite difficult because the English or German letters are easily confused. Even humans often need further inquiry to distinguish between the similar sounds of (for example) the letters **M** and **N**, or **D** and **T**. Throughout this text, we will use the terms "letter" and "word" interchangeably. The term "sentence" refers to a string of letters.

The Baseline MS-TDNN [5, 8] is a time shift invariant architecture of a neural network with a time delay neural network (TDNN) classifier. Figure 1 shows the architecture of the recognizer.

DIWLayer. Instead of phonemes, the output are now words, and error derivatives are backpropagated from the word units through the alignment paths and the front-end TDNN.

The choice of sensible objective functions is of great importance. For training on the phoneme level, there is an output vector $Y = (y_1, \dots, y_n)$ and a corresponding target vector $T = (t_1, \dots, t_n)$ for each frame in time. T represents the correct phoneme j in a "1-out-of- n " coding, i.e. $t_i = \delta_{ij}$. Standard *Mean Square Error* ($MSE = \sum_{i=1}^n (y_i - t_i)^2$) is problematic for "1-out-of- n " coding for large n ($n > 50$ in our case); consider that for a target $(1, 0, \dots, 0)$, the error for $(1.0, 0.2, \dots, 0.2)$ has only half the error than for $(1.0, 0.2, \dots, 0.2)$. This

McTellan

which (like cross-correlation) is not a

dependent penalty $Pen_w(d) = \log(k + prob_w(d))$, where
the pdf $prob_w(d)$ is approximated from the training data
is a small constant to avoid zero probabilities.
led to the accumulated score AS of the
* $Pen_w(d)$, whenever a word
re 2(b). The ra-
f the

4 . 2 . G E R M A N L E T T E R S

Ware in the process of creating a large data base of German spelled letters. At this time, more than 40,000 letters from 55 speakers (table 2) were collected and labeled. Volunteers are asked to spell a set of 50 to 150 sentences per, without artificial pauses between letters. A different set, consisting of 100 sentences, is randomly drawn from a