

# 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 delays", allowing for "tuning-in" to news speaker. We present results on our large and still growing letter data base, containing over 40.000 letter-spelled by 55 speakers.

Letter Recognition, Speaker-

## Hidden Layer

15 hidden units  
I ON

of letters is essential  
Input Layer  
vocabularies, such as  
English and German let-

ter need further  
e similar sounds of (for  
D and T. Throughout  
er" and "word" in-  
s to a string of

integrates the time-

[12] and a nonlin-

a word-level

process

by 16

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

choice of sensible objective functions is of great  
importance. In training on the phoneme level, there is

$(y_1, \dots, y_n)$  and a corresponding

$= (t_1, \dots, t_n)$  for each frame in time.  $T$

is a  $j$  in a "1-out-of- $n$ " code-

Mean Square Error ( $MSE =$

metric for "1-out-of- $n$ " codings

in case); consider for example

$(0, \dots, 0, 0)$ , the output  $(0.0, \dots, 0.0)$

is a more desirable output

avoided by

$$-(y_k - t_i)^2$$

"outliers" with an er-

error approaching 1.0.

achieved best re-

the "Classification

to maximize

$y_k$  and

absolute

bits:



A B)

A A B)

A proper treat-

ment is important for a

task at word

level

### Realizing of Insertion

on Errors. In continuous recognition

instead of looking at word units the well-known

algorithm [11] is used to find an opti-

mally specified sequence of words. The

misspellings cause many word

errors such as "TE" vs. "T"

and the proper duration model-

is [6], minimum phoneme

recognition". In ad-

dition, word ( $w$ )

dependent penalty  $Pen_w(d) = \log(k + prob_w(d))$ , where the pdf  $prob_w(d)$  is approximated from the training data and  $k$  is a small constant to avoid zero probabilities. This is added to the accumulated score  $AS$  of the word  $w$  with  $AS = AS + \lambda_w * Pen_w(d)$ , whenever a word  $w$  is predicted in figure 2(b). The ratio  $\lambda_w$  indicates the degree of influence of the word  $w$  on the degree of freedom of the model. Our approach is mathematically exact and does not require any change of the "weight" of the word  $w$ .

**the Sentence Level.**

On the sentence level, we trained to classify excerpted sentences from previously spoken sentences. This method to extend training shows the alignment between the original and the forced alignment (enforced), post-path, while the training is different.

both speaker letters

Models (ISMs) is measured, each of which is specific for a group of speakers, male/female speakers. They vary from simple

gender-specific sub-words

6  
8  
36.9  
91.0

(Res. Man. Spell - mode)	
1000 train, 11/900 test speaker/words	
	our M-TINN
	gender specific
90.8	92.0

(in % on the test set) on speaker connected letter tasks.

## German 2.5 GERMAN LETTERS

43/34823 train, 12/8206 test speaker/words

Size | exceeded the processing of a large data base of Ger-  
 ayer) | train spelled letters in this test, more than 40.000 let-  
 its | 976 from 52 speakers (table 2) were collected and labeled.  
 ts | Volunteers asked to spell 90 of 350 to 150 sentences

3: Word accuracy in % on German Spelling  
 in a natural manner, without artificial pauses between let-  
 ters. Each individual spells a different set, consisting of

three categories: proper names, drawn randomly from a

genius of 100.000 names, some randomness and

pseudo-random letter sequences. The latter subset

carefully acknowledge support by the NSF.

ed to increase the percentage of the less frequent

Creating the large German data.

as Q or X, to make sure there is a reasonable

without the tireless work

ing data for all letters. For example, after

her, Henrike Iness, the ratio of Q to

3 : 1000 to 75 : 1000 .

and phoneme boundaries, the data

ng the IANS IQ recognizer

. After a initial training

used to relabel the

the error rate on

### Systems

all	
Spr	letter
43	34823
1865	128206
3	11936
	55
	43029

ell Data Base.

addition to "Silence" and the

alphabet, the German alpha-

ls ("Unhaute") Ä, Ö,

ely, there are several

the official version "es-

", "scharfes-S", or in

also used. Since con-

en (- ) was also

erent possible

"Cedanken-

elling