

# Recognition of Alzheimer's Dementia From the Transcriptions of Spontaneous Speech Using fastText and CNN Models

Amit Meghanani , Anoop C. S. \* and Angarai Ganesan Ramakrishnan

*MILE Laboratory, Department of Electrical Engineering, Indian Institute of Science, Bengaluru, India*

Correspondence\*:  
Anoop C. S.  
anoopcs@iisc.ac.in

## 2 ABSTRACT

3 Alzheimer's dementia (AD) is a type of neurodegenerative disease that is associated with a  
4 decline in memory. However, speech and language impairments are also common in Alzheimer's  
5 dementia patients. This work is an extension of our previous work, where we had used  
6 spontaneous speech for Alzheimer's dementia recognition employing log-Mel spectrogram and  
7 Mel frequency cepstral coefficients (MFCC) as inputs to deep neural networks (DNN). In this work,  
8 we explore the transcriptions of spontaneous speech for dementia recognition and compare the  
9 results with several baseline results. We explore two models for dementia recognition - i) fastText  
10 and ii) convolutional neural network (CNN) with a single convolutional layer, to capture the n-gram  
11 based linguistic information from the input sentence. The fastText model uses a bag of bigrams  
12 and trigrams along with the input text to capture the local word orderings. In the CNN based  
13 model, we try to capture different n-grams (we use  $n = 2,3,4,5$ ) present in the text by adapting the  
14 kernel sizes to  $n$ . In both fastText and CNN architectures, the word embeddings are initialized  
15 using pre-trained GloVe vectors. We use bagging of 21 models in each of these architectures  
16 to arrive at the final model using which the performance on the test data is assessed. The best  
17 accuracies achieved with CNN and fastText models on the text data are 79.16% and 83.33%,  
18 respectively. The best root mean square errors (RMSE) on the prediction of mini-mental state  
19 examination (MMSE) score are 4.38 and 4.28 for CNN and fastText, respectively. The results  
20 suggest that the n-gram based features are worth pursuing, for the task of AD detection. fastText  
21 models have competitive results when compared to several baseline methods. Also, fastText  
22 models are shallow in nature and have the advantage of being faster in training and evaluation,  
23 by several orders of magnitude, compared to deep models.

24 **Keywords:** fastText, CNN, Alzheimer's, dementia, MMSE

## 1 INTRODUCTION

25 Dementia is a syndrome characterised by the decline in cognition that is significant enough to interfere with  
26 one's independent, daily functioning. Alzheimer's disease contributes to around 60–70% of dementia cases.  
27 Towards the final stages of Alzheimer's Dementia (AD), the patients lose control of their physical functions  
28 and depend on others for care. As there are no curative treatments for dementia, the early detection is  
29 critical to delay or slow down the onset or progression of the disease. The mini-mental state examination

30 (MMSE) is a widely used test to screen for dementia and to estimate the severity and progression of  
31 cognitive impairment.

32 AD affects the temporal characteristics of spontaneous speech. Changes in the spoken language are  
33 evident even in mild AD patients. Subtle language impairments such as difficulties in word finding and  
34 comprehension, usage of incorrect words, ambiguous referents, loss of verbal fluency, speaking too much  
35 at inappropriate times, talking too loudly, repeating ideas, and digressing from the topic are common in  
36 the early stages of AD (Savundranayagam et al., 2005) and they turn extreme in the moderate and severe  
37 stages. Szatlóczy et al. (2015) show that AD can be detected with the help of a linguistic analysis more  
38 sensitively than with other cognitive examinations. Mueller et al. (2018b) analyzed the connected language  
39 samples obtained from simple picture description tasks and found that the speech fluency and the semantic  
40 content features declined faster in participants with early mild cognitive impairment. The language profile  
41 of AD patients is characterized by “empty speech”, devoid of content words (Nicholas et al., 1985). They  
42 tend to use pronouns without proper noun references and indefinite terms like “this”, “that”, “thing” etc.,  
43 more often (Mueller et al., 2018a). These results motivate us to believe that modeling the transcriptions of  
44 the narrative speech in the cookie-theft picture description task using n-gram language models can help in  
45 the detection of AD and prediction of MMSE score.

46 In this work we address the AD detection and MMSE score prediction problems using two natural  
47 language processing (NLP) based models - i) fastText and ii) convolutional neural network (CNN). These  
48 models have the advantage that they can be easily structured to capture the linguistic cues in the form of  
49 n-grams from the transcriptions of the picture description task, provided with the Alzheimer’s Dementia  
50 Recognition through Spontaneous Speech (ADReSS) dataset (Luz et al., 2020). CNNs, though originated  
51 in computer vision, have become popular for NLP tasks and have achieved great results in sentence  
52 classification (Kim, 2014), semantic parsing (tau Yih et al., 2014), search query retrieval (Shen et al., 2014),  
53 and other traditional NLP tasks (Collober et al., 2011). Our convolutional neural network model draws  
54 inspiration from the work on sentence classification using CNNs (Kim, 2014). The fastText (Joulin et al.,  
55 2017) is a simple and efficient model for text classification (eg. tag prediction and sentiment analysis). The  
56 fundamental idea in the fastText classifier is to calculate the n-grams of an input sentence and append them  
57 to the end of the sentence. Our choice of fastText model is also motivated by its ability to often outperform  
58 deep learning classifiers in terms of accuracy and training/evaluation times (Joulin et al., 2017).

59 The rest of the paper is organised as follows. Section 2 discusses the ADReSS dataset in detail. Section  
60 3 discusses the baseline results in AD detection. Section 4 discusses our proposed NLP based models  
61 followed by the listing of results in section 5. Our results and conclusions are discussed in section 6.

## 2 ADDRESS DATASET

62 The ADReSS dataset (Luz et al., 2020) is designed to provide Alzheimer’s research community with a  
63 standard platform for AD detection and MMSE score prediction. The dataset is acoustically pre-processed  
64 and balanced in terms of age and gender. It consists of audio recordings and transcriptions (in CHAT  
65 format (Macwhinney, 2009)) of the Cookie Theft picture description task, elicited from subjects in the age  
66 group of 50-80 years. The training set consists of data from 108 subjects, 54 each from AD and non-AD  
67 classes. The test set has data from 48 subjects, again balanced with respect to AD and non-AD classes.  
68 More information on the ADReSS dataset can be found in the ADReSS challenge baseline paper (Luz  
69 et al., 2020).

### 3 REVIEW OF BASELINE METHODS

70 This section provides a brief overview of the various approaches for AD detection and MMSE score  
 71 prediction on ADReSS dataset. These approaches can be broadly classified into 3 types based on the type  
 72 of the features used in the problem- i) acoustic feature based, ii) linguistic feature based and iii) a fusion of  
 73 acoustic and linguistic features. The performance of different approaches on the AD detection and MMSE  
 74 score prediction tasks are compared using the accuracy and root mean square error (RMSE) measures  
 75 computed on the ADReSS test set.

$$Accuracy = \frac{TN + TP}{N} \quad (1)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N}} \quad (2)$$

76 where  $N$  is the total number of subjects involved in the study,  $TP$  the number of true positives and  $TN$   
 77 the number of true negatives.  $\hat{y}_i$  and  $y_i$  are the estimated and target MMSE scores for  $i^{th}$  test sample. The  
 78 results of different approaches on the ADReSS dataset are summarized in Table 1.

#### 79 3.1 Acoustic feature-based methods

80 Luz et al. (2020), explore several acoustic features like extended Geneva minimalistic acoustic parameter  
 81 set (eGeMAPS) (Eyben et al., 2016), emobase, ComParE-2013 (Eyben et al., 2013), and multi-resolution  
 82 cochleagram (MRCG) (Chen et al., 2014) feeding the traditional machine learning algorithms like linear  
 83 discriminant analysis, decision trees, nearest neighbour, random forests and support vector machines.  
 84 In our previous work (Meghanani et al., 2021), we have used CNN/ResNet + long short-term memory  
 85 (LSTM) networks and pyramidal bidirectional LSTM + CNN networks trained on log-Mel spectrogram and  
 86 Mel-frequency cepstral coefficient (MFCC) features extracted from the spontaneous speech. Pompili et al.  
 87 (2020), exploit the pre-trained models to produce i-vector and x-vector based acoustic feature embeddings.  
 88 They evaluate x-vector, i-vector, and statistical speech-based functional features. Rhythmic features are  
 89 proposed in (Campbell et al., 2020), as lower speaking fluency is a common pattern in patients with AD.  
 90 Koo et al. (2020), use VGGish (Hershey et al., 2017) trained with Audio Set (Gemmeke et al., 2017) for  
 91 audio classification. They have proposed a modified version of convolutional recurrent neural network  
 92 (CRNN), where an attention layer is the forefront layer of the network, and fully connected layers follow  
 93 the recurrent layer.

#### 94 3.2 Linguistic feature-based methods

95 Recently, there have been multiple attempts on the AD detection problem based on text based features  
 96 and models. Searle et al. (2020), use traditional machine learning techniques like support vector machines  
 97 (SVMs), gradient boosting decision trees (GBDT), and conditional random fields (CRFs). They also try deep  
 98 learning transformer based models, specifically, bidirectional encoder representations from transformers  
 99 (BERT) (Devlin et al., 2019), RoBERTa (Liu et al., 2019) and DistilBERT/DistilRoBERTa (Sanh et al.,  
 100 2019). Pompili et al. (2020), encode each word of the clean transcriptions into 768-dimensional context  
 101 embedding vector using a frozen English BERT model pre-trained with 12-layers. Three different neural  
 102 models are trained on top of contextual word embeddings: (i) global maximum pooling, (ii) bidirectional  
 103 long short-term memory (BLSTM) based recurrent neural networks (RNN) provided with an attention  
 104 module, and (iii) the second model augmented with part-of-speech (POS) embeddings. In the work  
 105 of Campbell et al. (2020), authors have used the manual transcripts to extract linguistic information  
 106 (interventions, vocabulary richness, frequency of verbs, nouns, POS-tagging, etc.) for creating the input  
 107 features of the classifier. They use another sequential deep learning based classifier, which classifies

108 directly from the sequence of Global Vectors (GloVe) based word embeddings. Koo et al. (2020), use  
 109 transformer (Vaswani et al., 2017) based language models, generative pretraining (GPT) (Radford et al.,  
 110 2018), RoBERTa (Liu et al., 2019), and transformer-XL (Dai et al., 2020) to get textual features and  
 111 perform classification and regression tasks using a modified convolutional recurrent neural network based  
 112 structure.

113 Graph based representation of word features (Tomás and Radev, 2012), (Cong and Liu, 2014), which  
 114 have shown promise in classifying texts (De Arruda et al., 2016) are also employed for detection of mild  
 115 cognitive impairments. Santos et al. (2017) model transcripts as complex networks and enrich them with  
 116 word embedding to better represent short texts produced in neuro-psychological assessments. They use  
 117 metrics of topological properties of complex networks in a machine learning classification approach to  
 118 distinguish between healthy subjects and patients with mild cognitive impairments. Such graph based  
 119 techniques have also been used in the word sense disambiguation (WSD) tasks to identify the meaning of  
 120 words in a given context for specific words conveying multiple meanings.(Corra et al., 2018). They suggest  
 121 that a bipartite network model with local features employed to characterise the context can be useful in  
 122 improving the semantic characterization of written texts without the use of deep linguistic information.

### 123 3.3 Bimodal Methods

124 Methods with bimodal input features (both acoustic and linguistic) are also used for AD recognition in  
 125 various studies like (Pompili et al., 2020), (Campbell et al., 2020), (Sarawgi et al., 2020b), (Koo et al.,  
 126 2020), (Sarawgi et al., 2020a), and (Rohanian et al., 2020). However, in this work, we restrict ourselves to  
 127 the NLP-based approaches.

## 4 PROPOSED NLP-BASED METHODS

### 128 4.1 Data Preparation

129 In this work, we explore the linguistic features for AD detection and hence only the textual transcripts in  
 130 the ADReSS dataset are used. The transcripts contain the conversational content between the participant  
 131 and the investigator. This include pauses in speech, laughter and discourse markers such as ‘um’ and ‘uh’.  
 132 Each transcript is considered as a single data point with their corresponding AD label and MMSE score.  
 133 We create two transcription level datasets after pre-processing the transcripts as in Searle et al. (2020) -  
 134 1) PAR: containing the utterances of participant alone, 2) PAR+INV: containing utterances from both the  
 135 participant and the investigator. In addition to the preprocessing performed in Searle et al. (2020), we keep  
 136 PAR and INV tags as well in the data (which defines whether the utterance is spoken by the participant or  
 137 the investigator).

### 138 4.2 CNN Model

139 Language impairments like difficulties in lexical retrieval, loss of verbal fluency, and breakdown in  
 140 comprehension of higher order written and spoken languages are common in AD patients. Hence the  
 141 linguistic information like the n-grams present in the input sentence, may provide good cues for AD  
 142 detection. Any  $n \times d$  CNN filter, where  $n$  is the number of sequential words looked over by the filter and  $d$   
 143 is the dimension of word embedding, can be viewed as a feature detector looking for a specific n-gram in  
 144 the input that can capture the language impairments associated with AD.

145 We describe the details of the CNN model from the work (Kim, 2014) as follows. Let  $z_i \in R^d$  be a  
 146  $d$ -dimensional word vector corresponding to the  $i$ -th word in the sentence. A sentence of length  $L$  is  
 147 represented as  $\{z_1, z_2, \dots, z_L\}$ . Let  $z_{i:i+j}$  represents the concatenation of the words  $z_i, z_{i+1}, \dots, z_{i+j}$ . A  
 148 convolution operation involves a filter  $w \in R^{nd}$ , which is applied to a window of  $n$  words to produce a

149 new feature as shown in equation 3, where  $s_i$  is generated from a window of words  $z_{i:i+n-1}$  by

$$s_i = f(w \cdot z_{i:i+n-1} + b) \quad (3)$$

150 In the equation 3,  $f$  is a non-linear function and  $b$  is the bias term. A feature map  $\mathcal{S}$  is obtained by applying  
151 the filter to all possible windows of words in the sentence  $[z_{1:n}, z_{2:n+1}, \dots, z_{L-n+1:L}]$ .

$$\mathcal{S} = [s_1, s_2, \dots, s_{L-n+1}] \quad (4)$$

152 A max-pool over time (Collober et al., 2011) is performed over the feature map to get  $s_{max} = \max \mathcal{S}$   
153 as the feature corresponding to that filter. This corresponds to the n-gram that is “most relevant” in the  
154 AD recognition task. The weights of the filters, which in turn determine the “most relevant” feature,  
155 are learnt using backpropagation. CNNs are trained with just one layer of convolution. Variable length  
156 sentences are automatically handled by the pooling scheme. We use pre-trained 100-dimensional GloVe  
157 word vectors (Pennington et al., 2014) for word embedding. Multiple kernels of sizes  $2 \times 100$ ,  $3 \times 100$ ,  
158  $4 \times 100$  and  $5 \times 100$  are employed to have a look at the bigrams, trigrams, 4-grams and 5-grams within  
159 the text. We use 100 filters each with height 2, 3, 4 and 5. Multiple configurations with filter sizes [2,3,4],  
160 [3,4,5] and [2,3,4,5] are applied which are referred to as CNN-bi+tri+4 gram, CNN-tri+4+5 gram, and  
161 CNN-bi+tri+4+5 gram in our tables. The outputs of the filter are concatenated together to form a single  
162 vector. Dropout with probability  $p = 0.5$  is applied on the concatenated filter output and the results are  
163 passed through a linear layer for the final prediction task. The linear layer weights up the evidences from  
164 each of these n-grams and make a final decision. Fig. 1 shows the basic CNN operation over an example  
165 sentence.

#### 166 4.2.1 Training Details

167 For the classification task, training is performed for 100 epochs with a batch size of 16. Adam optimizer  
168 is used with a learning rate of 0.001. Model with the lowest validation loss is saved and used for prediction.  
169 Since AD classification is a two class problem, binary cross entropy with logits loss is used as the loss  
170 function. For the MMSE score prediction task, the output layer is a fully connected layer with linear  
171 activation function. In the regression task the network is trained for 1500 epochs with the objective to  
172 minimize the mean squared error.

173 We use bootstrap aggregation of models known as bagging Breiman (1996) to predict the final  
174 labels/MMSE scores for test samples. Bootstrap aggregation is an ensemble technique to improve the  
175 stability and accuracy of machine learning models. It combines the prediction from multiple models. It  
176 also reduces variance and helps to avoid overfitting. We fit 21 models and the outputs are combined by a  
177 majority voting scheme for final classification. In the regression task, the outputs of these bootstrap models  
178 are averaged to arrive at the final MMSE score.

#### 179 4.3 fastText

180 fastText based classifiers calculate the n-grams of an input sentence explicitly and append them to the end  
181 of the sentence. In this work, we use bigrams and trigrams. We conducted the experiments with 4-grams  
182 as well, but the results did not show any improvement over the use of trigrams. This bag of bigrams and  
183 trigrams acts as additional features to capture some information about the local word order.

184 Figure 2 shows the architecture of fastText model. The fastText model has 2 layers, an embedding layer  
185 and a linear layer. The embedding layer calculates the word embedding (100-dimensional) for each word.  
186 The average of all these word embeddings is calculated and fed through the linear layer for final prediction  
187 as described in Fig. 2. fastText models are faster for training and evaluation by many orders of magnitude,



188 compared to the “deep” models. As mentioned in the work (Joulin et al., 2017), fastText can be trained on  
189 more than one billion words in less than ten minutes using a standard multicore CPU, and classify half a  
190 million sentences among 312K classes in less than a minute.

#### 191 4.3.1 Training Details

192 All training details are the same as mentioned in section 4.2.1. The only difference is that dropout is not  
193 used in this model. Here also we use 21 bootstrapping models and the outputs are combined as described in  
194 section 4.2.1.

## 5 RESULTS

195 We have performed 5-fold cross-validation, to estimate the generalization error. One of the folds has 20  
196 validation samples and the remaining four have 22 validation samples. The results of cross-validation on  
197 CNN and fastText models trained on PAR and PAR+INV sets are listed in Table 2. The best performing  
198 model for classification during the cross validation was fastText with bigrams on the PAR+INV set, which  
199 yields an average cross validation accuracy of 86.09%. Among the CNN models, tri+4+5 grams give the  
200 best accuracy in both PAR (77.54%) and INV+PAR (81.27%) sets. As far as accuracy is concerned, both  
201 the CNN and fastText models seem to benefit with the inclusion of utterances from the investigator. For  
202 the prediction of MMSE score, CNN with bi+tri+4+5 grams (RMSE of 4.38) was the best. The fastText  
203 models seem to get a clear advantage in RMSE with the addition of the utterances from the investigator.  
204 However such a large difference in RMSE is not observable between the CNN models using PAR and  
205 INV+PAR sets. The cross-validation results confirmed our belief that the n-grams from the transcriptions  
206 of the picture description task could be useful in the detection of AD.

207 Table 3 lists the classification accuracy and RMSE in the prediction of MMSE score on the test set of the  
208 ADRess corpus. The table also lists the precision, recall and  $F_1$  score for each class. They are computed  
209 as precision  $\pi = \frac{TP}{TP+FP}$ , recall  $\rho = \frac{TP}{TP+FN}$ , and  $F_1$  score  $= \frac{2\pi\rho}{\pi+\rho}$ , where  $TP$ ,  $FP$ ,  $TN$  and  $FN$  are  
210 the number of true positives, false positives, true negatives and false negatives, respectively. The listed  
211 results are obtained after bootstrapping with 21 samples. The best classification accuracy is 83.33% which  
212 is achieved using fastText model with appended bigrams and trigrams. The accuracies are similar in both  
213 PAR and PAR+INV sets using the fastText model. The maximum accuracy obtained with CNN models is  
214 79.16%, which is achieved on the INV+PAR set using bi+tri+4 grams or tri+4+5 grams. In the detection  
215 task, the CNN models seem to get some advantage by the addition of utterances from the investigator. Also  
216 the accuracies seem to degrade when bigrams, trigrams, 4-grams and 5-grams are considered together. This  
217 behaviour is consistent across the PAR and PAR+INV sets. The best RMSE in the prediction of MMSE  
218 score, is 4.28 which is obtained on the PAR+INV set using fastText model employing only bigrams. In  
219 the regression task using fastText, the use of bigrams achieve slightly better RMSE compared to the use  
220 of both bigrams and trigrams. Also the fastText models seem to benefit from the use of utterances from  
221 the investigator. In contrast, CNN models do not seem to get any specific advantage with the inclusion of  
222 investigator’s utterances. The performance of the CNN models remain almost the same across the use of  
223 bi+tri+4, tri+4+5, and bi+tri+4+5 grams.

## 6 DISCUSSION AND CONCLUSIONS

224 In this work, we explore two models - CNN with a single convolution layer and fastText, to address  
225 the problem of AD classification and prediction of MMSE score from the transcriptions of the picture  
226 description task. The choice of these models were based on our initial belief that modeling the transcriptions  
227 of the narrative speech in the picture description task using n-grams could give some indication on the  
228 status of AD. The chosen models are also shallow. The number of parameters are much less than the usual

229 deep learning architectures and hence they can be trained and evaluated quite fast. Yet, the performance of  
230 these models is competitive with the baseline results reported with complex models (refer Table 1). The  
231 results suggest that the n-gram based features are worth pursuing, for the task of AD detection.

232 Among the considered models, fastText model with bigrams and trigrams appended to the input, achieves  
233 the best classification accuracy (83.33%). In the regression task, the best results (RMSE of 4.28) are  
234 achieved using fastText model with only the bigrams appended to the input. The fastText models have a  
235 clear edge over CNN in the classification task. Empirical evidences suggest that fastText models benefit  
236 from the inclusion of utterances from the investigator in the regression task, though they do not make much  
237 difference in the classification task. The CNN models on the other hand perform better on the PAR+INV  
238 sets in the classification task. In the regression task, their performance is similar across the PAR and  
239 PAR+INV sets. Bigrams have an edge over bi+tri grams in fastText, when used for prediction of MMSE  
240 score. However, the performance of the CNN models remain almost the same across the use of bi+tri+4,  
241 tri+4+5, and bi+tri+4+5 grams, in the regression task.

## REFERENCES

- 242 Breiman, L. (1996). Bagging predictors. *Mach. Learn.* 24, 123–140. doi:10.1007/BF00058655
- 243 Campbell, E. L., Docío-Fernández, L., Raboso, J. J., and García-Mateo, C. (2020). Alzheimer's dementia  
244 detection from audio and text modalities. *arXiv* 2008.04617
- 245 Chen, J., Wang, Y., and Wang, D. (2014). A feature study for classification-based speech separation at  
246 low signal-to-noise ratios. *IEEE/ACM Transactions on Audio Speech and Language Processing* 22,  
247 1993–2002. doi:10.1109/TASLP.2014.2359159
- 248 Collober, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural  
249 language processing (almost) from scratch. *Journal of Machine Learning Research* 12, 2493–2537.  
250 doi:10.5555/1953048.2078186
- 251 Cong, J. and Liu, H. (2014). Approaching human language with complex networks. *Physics of Life*  
252 *Reviews* 11, 598 – 618. doi:https://doi.org/10.1016/j.plrev.2014.04.004
- 253 Corra, E. A., Lopes, A. A., and Amancio, D. R. (2018). Word sense disambiguation. *Inf. Sci.* 442, 103–113.  
254 doi:10.1016/j.ins.2018.02.047
- 255 Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q., and Salakhutdinov, R. (2020). Transformer-XL:  
256 Attentive language models beyond a fixed-length context. *Proceedings of the 57th Annual Meeting of*  
257 *the Association for Computational Linguistics*, 2978–2988doi:10.18653/v1/P19-1285
- 258 De Arruda, H., Costa, L., and Amancio, D. (2016). Using complex networks for text classification:  
259 Discriminating informative and imaginative documents. *EPL* 113. doi:10.1209/0295-5075/113/28007
- 260 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional  
261 transformers for language understanding. *Proceedings of the 2019 Conference of the North American*  
262 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1*  
263 *(Long and Short Papers)*, 4171–4186doi:10.18653/v1/N19-1423
- 264 Edwards, E., Dognin, C., Bollepalli, B., and Singh, M. (2020). Multiscale System for Alzheimer's  
265 Dementia Recognition Through Spontaneous Speech. *Proc. Interspeech 2020*, 2197–2201doi:10.21437/  
266 Interspeech.2020-2781
- 267 Eyben, F., Scherer, K. R., Schuller, B. W., Sundberg, J., André, E., Busso, C., et al. (2016). The Geneva  
268 minimalistic acoustic parameter set (GeMAPS) for voice research and affective computing. *IEEE*  
269 *Transactions on Affective Computing* 7, 190–202
- 270 Eyben, F., Weninger, F., Gross, F., and Schuller, B. (2013). Recent developments in openSMILE, the  
271 Munich open-source multimedia feature extractor. *Proceedings of the 2013 ACM Multimedia Conference*

- 272 , 835–838doi:10.1145/2502081.2502224
- 273 Gemmeke, J., Ellis, D., Freedman, D., Jansen, A., Lawrence, W., Moore, R., et al. (2017). Audio set: An  
274 ontology and human-labeled dataset for audio events. *2017 IEEE International Conference on Acoustics,  
275 Speech and Signal Processing (ICASSP)* , 776–780doi:10.1109/ICASSP.2017.7952261
- 276 Hershey, S., Chaudhuri, S., Ellis, D., Gemmeke, J., Jansen, A., Moore, R. C., et al. (2017). CNN  
277 architectures for large-scale audio classification. *2017 IEEE International Conference on Acoustics,  
278 Speech and Signal Processing (ICASSP)* , 131–135
- 279 Joulin, A., Grave, E., Bojanowski, P., and Mikolov, T. (2017). Bag of tricks for efficient text classification.  
280 *Proceedings of the 15th Conference of the European Chapter of the Association for Computational  
281 Linguistics: Volume 2, Short Papers* , 427–431
- 282 Kim, Y. (2014). Convolutional neural networks for sentence classification. *Proceedings of the 2014  
283 Conference on Empirical Methods in Natural Language Processing (EMNLP)* , 1746–1751doi:10.3115/  
284 v1/D14-1181
- 285 Koo, J., Lee, J. H., Pyo, J., Jo, Y., and Lee, K. (2020). Exploiting Multi-Modal Features from Pre-Trained  
286 Networks for Alzheimer’s Dementia Recognition. *Proc. Interspeech 2020* , 2217–2221doi:10.21437/  
287 Interspeech.2020-3153
- 288 Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., et al. (2019). RoBERTa: A robustly optimized  
289 BERT pretraining approach. *ArXiv abs/1907.11692*
- 290 Luz, S., Haider, F., de la Fuente, S., Fromm, D., and MacWhinney, B. (2020). Alzheimer’s Dementia  
291 Recognition Through Spontaneous Speech: The ADReSS Challenge. *Proc. Interspeech 2020* , 2172–  
292 2176doi:10.21437/Interspeech.2020-2571
- 293 Macwhinney, B. (2009). The CHILDES Project Part 1: The CHAT Transcription Format doi:10.1184/R1/  
294 6618440.v1
- 295 Meghanani, A., Anoop, C. S., and Ramakrishnan, A. G. (2021). An exploration of log-mel spectrogram  
296 and MFCC features for Alzheimer’s dementia recognition from spontaneous speech. *Accepted in The  
297 8th IEEE Spoken Language Technology Workshop (SLT), January 19-22, 2021*
- 298 Mueller, K. D., Hermann, B., Mecollari, J., and Turkstra, L. S. (2018a). Connected speech and language  
299 in mild cognitive impairment and Alzheimer’s disease: A review of picture description tasks. *J Clin Exp  
300 Neuropsychol.* 40, 917–939. doi:10.1080/13803395.2018.1446513
- 301 Mueller, K. D., Kosciak, R. L., Hermann, B., Johnson, S. C., and Turkstra, L. S. (2018b). Declines in  
302 connected language are associated with very early mild cognitive impairment: Results from the Wisconsin  
303 registry for Alzheimer’s prevention. *Frontiers in Aging Neuroscience* 9. doi:10.3389/fnagi.2017.00437
- 304 Nicholas, M., Obler, L. K., Albert, M., and Helm-Estabrooks, N. (1985). Empty speech in Alzheimer’s  
305 disease and fluent aphasia. *Journal of speech and hearing research* 28, 405–410. doi:https://doi.org/10.  
306 1044/jshr.2803.405
- 307 Pappagari, R., Cho, J., Moro-Velázquez, L., and Dehak, N. (2020). Using State of the Art Speaker  
308 Recognition and Natural Language Processing Technologies to Detect Alzheimer’s Disease and Assess  
309 its Severity. *Proc. Interspeech 2020* , 2177–2181doi:10.21437/Interspeech.2020-2587
- 310 Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation.  
311 *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* ,  
312 1532–1543doi:10.3115/v1/d14-1162
- 313 Pompili, A., Rolland, T., and Abad, A. (2020). The INESC-ID Multi-Modal System for the ADReSS 2020  
314 Challenge. *Proc. Interspeech 2020* , 2202–2206doi:10.21437/Interspeech.2020-2833
- 315 Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). Improving language understanding by  
316 generative pre-training. <https://www.cs.ubc.ca/amuham01/LING530/papers/radford2018improving.pdf>



- 317 Rohanian, M., Hough, J., and Purver, M. (2020). Multi-Modal Fusion with Gating Using Audio, Lexical  
318 and Disfluency Features for Alzheimer's Dementia Recognition from Spontaneous Speech. *Proc.*  
319 *Interspeech 2020* , 2187–2191doi:10.21437/Interspeech.2020-2721
- 320 Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller,  
321 faster, cheaper and lighter. *ArXiv abs/1910.01108*
- 322 Santos, L., Corrêa Júnior, E. A., Oliveira Jr, O., Amancio, D., Mansur, L., and Aluísio, S. (2017).  
323 Enriching complex networks with word embeddings for detecting mild cognitive impairment from  
324 speech transcripts. *Proceedings of the 55th Annual Meeting of the Association for Computational*  
325 *Linguistics (Volume 1: Long Papers)* , 1284–1296doi:10.18653/v1/P17-1118
- 326 Sarawgi, U., Zulfikar, W., Khincha, R., and Maes, P. (2020a). Uncertainty-aware multi-modal ensembling  
327 for severity prediction of Alzheimer's dementia. *ArXiv abs/2010.01440*
- 328 Sarawgi, U., Zulfikar, W., Soliman, N., and Maes, P. (2020b). Multimodal inductive transfer learning for  
329 detection of Alzheimer's dementia and its severity. *arXiv preprint arXiv:2009.00700*
- 330 Savundranayagam, M., Hummert, M. L., and Montgomery, R. (2005). Investigating the effects of  
331 communication problems on caregiver burden. *The journals of gerontology. Series B, Psychological*  
332 *sciences and social sciences* 60 1, S48–55
- 333 Searle, T., Ibrahim, Z., and Dobson, R. (2020). Comparing natural language processing techniques for  
334 Alzheimer's dementia prediction in spontaneous speech. *Proc. Interspeech 2020* , 2192–2196doi:10.  
335 21437/Interspeech.2020-2729
- 336 Shen, Y., He, X., Gao, J., li Deng, and Mesnil, G. (2014). Learning semantic representations using  
337 convolutional neural networks for web search. *WWW 2014* , 373–374doi:10.1145/2567948.2577348
- 338 Syed, M. S. S., Syed, Z. S., Lech, M., and Pirogova, E. (2020). Automated Screening for Alzheimer's  
339 Dementia Through Spontaneous Speech. *Proc. Interspeech 2020* , 2222–2226doi:10.21437/Interspeech.  
340 2020-3158
- 341 Szatlóczki, G., Hoffmann, I., Vincze, V., Kálmán, J., and Pákási, M. (2015). Speaking in Alzheimer's  
342 disease, is that an early sign? Importance of changes in language abilities in Alzheimer's disease.  
343 *Frontiers in aging neuroscience* 7. doi:10.3389/fnagi.2015.00195
- 344 tau Yih, W., He, X., and Meek, C. (2014). Semantic parsing for single-relation question answering.  
345 *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2:*  
346 *Short Papers)* , 643–648doi:10.3115/v1/P14-2105
- 347 Tomás, D. R. M. and Radev, D. (2012). Graph-based natural language processing and information retrieval.  
348 *Machine Translation* 26, 277–280. doi:10.1007/s10590-011-9122-9
- 349 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., et al. (2017). Attention is all  
350 you need. *Proceedings of the 31st International Conference on Neural Information Processing Systems*  
351 2017-December, 5999–6009
- 352 Yuan, J., Bian, Y., Cai, X., Huang, J., Ye, Z., and Church, K. (2020). Disfluencies and Fine-Tuning  
353 Pre-Trained Language Models for Detection of Alzheimer's Disease. *Proc. Interspeech 2020* , 2162–  
354 2166doi:10.21437/Interspeech.2020-2516

**Table 1.** Baseline methods on ADReSS test set

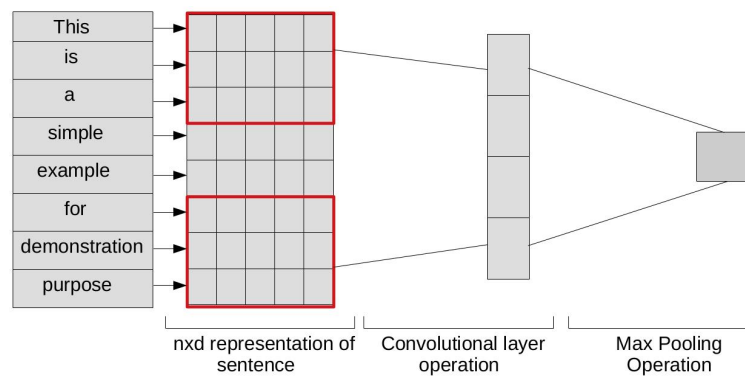
Model	Accuracy	RMSE
(Searle et al., 2020), DistilBERT	81.25%	4.58
(Searle et al., 2020), SVM+CRF	81.25%	5.22
(Pompili et al., 2020), x-vectors SRE	54.17%	–
(Pompili et al., 2020), Sentence embedding	72.92%	–
(Pompili et al., 2020), Fusion of system	81.25%	–
(Luz et al., 2020), linguistic	75.00%	5.20
(Sarawgi et al., 2020b), Ensemble	83.33%	4.60
(Koo et al., 2020), VGGish	72.92%	5.07
(Koo et al., 2020), Transformer-XL	81.25%	4.01
(Koo et al., 2020), VGGish+GloVe	77.08%	4.33
(Koo et al., 2020), VGGish+Transformer-XL	75.00%	3.74
(Koo et al., 2020), Ensembled Output	81.25%	3.77
(Campbell et al., 2020), Fusion II	75.00%	–
(Campbell et al., 2020), Fusion I	72.92%	–
(Campbell et al., 2020), RNN Model	75.00%	–
(Campbell et al., 2020), fluency	60.42%	–
(Campbell et al., 2020), x-vector	54.17%	–
(Sarawgi et al., 2020a), UA Ensemble	–	4.35
(Sarawgi et al., 2020a), UA Ensemble (weighted)	–	3.93
(Pappagari et al., 2020), Acoustic and Transcript	75.00%	5.37
(Rohanian et al., 2020), LSTM (Lexical+Dis)	72.92%	4.88
(Rohanian et al., 2020), LSTM with Gating (Acoustic+Lexical)	77.08%	4.57
(Rohanian et al., 2020), LSTM with Gating (Acoustic+Lexical+Dis)	79.17%	4.54
(Yuan et al., 2020), ERNIE3p	89.58%	–
(Syed et al., 2020)	85.42%	4.30
(Edwards et al., 2020), Phonemes and Audio	79.17%	–
(Meghanani et al., 2021), CNN-LSTM with MFCC	64.58%	6.24
(Meghanani et al., 2021), pBLSTM-CNN with log-Mel	52.08%	5.90
(Meghanani et al., 2021), ResNet-LSTM with log-Mel	62.50%	5.98

**Table 2.** Average 5-fold cross-validation results for AD classification and RMSE values

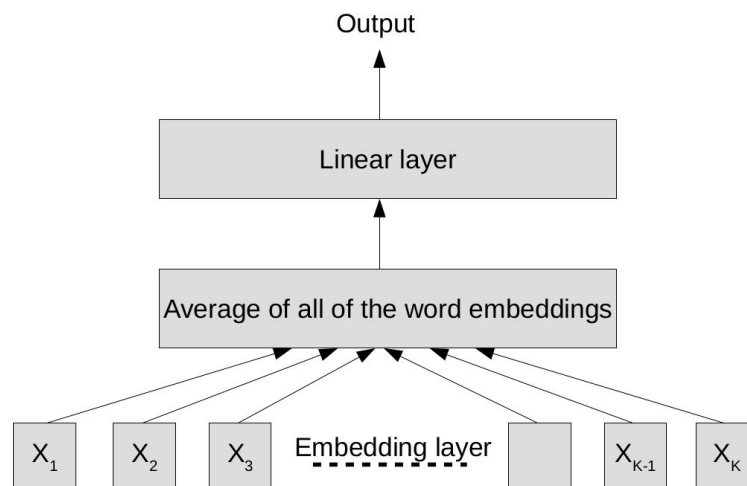
Dataset	Model	Accuracy	RMSE
PAR	CNN, bi+tri+4 gram	73.91%	4.55
PAR	CNN, tri+4+5 gram	77.54%	4.41
PAR	CNN, bi+tri+4+5 gram	76.54%	4.65
PAR	fastText, bigram	80.54%	5.43
PAR	fastText, bi+trigram	82.36%	5.40
PAR+INV	CNN, bi+tri+4 gram	80.18%	4.63
PAR+INV	CNN, tri+4+5 gram	81.27%	4.53
PAR+INV	CNN, bi+tri+4+5 gram	80.36%	4.38
PAR+INV	fastText, bigram	86.09%	4.66
PAR+INV	fastText, bi+trigram	85.90%	4.81

**Table 3.** Results on ADReSS test set

Dataset	Model	Class	Precision	Recall	F1 Score	Accuracy	RMSE
PAR	CNN, bi+tri+4 gram	Non-AD	0.74	0.71	0.72	72.91%	4.38
		AD	0.72	0.75	0.73		
PAR	CNN, tri+4+5 gram	Non-AD	0.76	0.67	0.71	72.91%	4.46
		AD	0.70	0.79	0.75		
PAR	CNN, bi+tri+4+5 gram	Non-AD	0.71	0.71	0.71	70.83%	4.42
		AD	0.71	0.71	0.71		
PAR	fastText, bigram	Non-AD	0.78	0.88	0.82	81.25%	4.51
		AD	0.86	0.75	0.80		
PAR	fastText, bi+trigram	Non-AD	0.81	0.88	0.84	<b>83.33%</b>	4.87
		AD	0.86	0.79	0.83		
PAR+INV	CNN, bi+tri+4 gram	Non-AD	0.77	0.83	0.80	79.16%	4.48
		AD	0.82	0.75	0.78		
PAR+INV	CNN, tri+4+5 gram	Non-AD	0.77	0.83	0.80	79.16%	4.47
		AD	0.82	0.75	0.78		
PAR+INV	CNN, bi+tri+4+5 gram	Non-AD	0.74	0.71	0.72	72.91%	4.44
		AD	0.72	0.75	0.73		
PAR+INV	fastText, bigram	Non-AD	0.78	0.88	0.82	81.25%	<b>4.28</b>
		AD	0.86	0.75	0.80		
PAR+INV	fastText, bi+trigram	Non-AD	0.79	0.92	0.85	<b>83.33%</b>	4.47
		AD	0.90	0.75	0.82		



**Figure 1.** Demonstration of CNN over text for an example sentence.



**Figure 2.** fastText model (Joulin et al., 2017) with appended n-gram features ( $X_1, X_2, X_3, \dots, X_{K-1}, X_K$ ) as input.