



2 **Automatic assessment of children’s oral reading using speech**  
3 **recognition and prosody modeling**

4 **Kamini Sabu<sup>1</sup> · Preeti Rao<sup>1</sup>**

5 Received: 16 May 2018 / Accepted: 7 June 2018  
6 © CSI Publications 2018

Author Proof

7 **Abstract** This work targets building an oral reading “tu-  
8 tor” that provides automatic and reliable feedback to  
9 children learning to read. The work uses state-of-the-art  
10 speech recognition technology coupled with prosody  
11 modeling. The system is tested on available datasets of  
12 children’s readings in English as L2. The expected chal-  
13 lenges relate to dealing with children’s speech with a  
14 variety of skill levels. Both word decoding accuracies and  
15 prosody attributes like phrasing and prominence are con-  
16 sidered for assessment. The relation between different  
17 acoustic features computed from the speech signal and the  
18 perceived quality will be investigated. The goal is to have a  
19 system that can provide feedback and evaluation that is  
20 highly correlated with that of human judges such as lan-  
21 guage teachers.

22 **Keywords** Prosody · Phrasing · Prominence · Reading  
23 assessment · Child speech

24 **1 Introduction**

25 As per the Annual Survey Education Reports [1] by Pra-  
26 tham over last 4–5 years, basic literacy skills of children in  
27 rural India are below par. Almost 70% of the students in  
28

5th standard can’t read even a simple English sentence in 29  
2nd standard curriculum book. In spite of ‘education for 30  
all’ movement by government, the reading skills of chil- 31  
dren have not improved—two major reasons being the 32  
shortage of qualified teachers and dismal student–teacher 33  
ratio. In this scenario, we aim to build an automatic reading 34  
evaluation system to assess children’s oral reading in terms 35  
of reading miscues, speech rate and fluency [2, 3]. We 36  
further aim to emulate the teacher or expert in providing 37  
motivating but sound feedback which will help students 38  
improve their reading skill. 39

Reading skills, as per the U.S. National Reading Panel 40  
[4], comprise of word decoding accuracy, reading rate and 41  
fluency. Word decoding accuracy is indicated through a 42  
count of detected miscues in terms of word substitutions, 43  
omissions or insertions. Speech rate is measured in terms of 44  
number of words read per minute or number of words 45  
correctly read per minute. Fluency has two components— 46  
lexical fluency and prosodic fluency. Lexical fluency 47  
indicates false starts and hesitations, while prosodic fluency 48  
refers to expressiveness which can further be considered in 49  
terms of phrasing and prominence. The ability of chunking 50  
the complete sentence into meaningful word sequences is 51  
called phrasing. The emphasis of words indicating new 52  
information in order to direct a listener’s attention to this is 53  
called prominence. Prosody is known to be assistive and 54  
indicative of comprehension of the read text [5, 6] 55

TBall [7], Listen [8] and FLORA [9] are some major 56  
research groups working in this area. TBall dealt with 57  
kindergarten students reading isolated word lists. They 58  
used specially designed language model in ASR to detect 59  
lexical disfluency. However, lexical disfluency is not 60  
enough to rate the oral reading fluency. The Listen group 61  
correlates the prosodic contours (pitch, energy, intensity 62  
and latency) of child speech with corresponding adult 63

A1 MeitY Visvesvaraya PhD scheme.

A2 ✉ Kamini Sabu  
A3 kaminisabu@iitb.ac.in

A4 Preeti Rao  
A5 prao@ee.iitb.ac.in

A6 <sup>1</sup> Department of Electrical Engineering, Indian Institute of  
A7 Technology Bombay, Mumbai, India

64 speech to rate fluency of 7–10 years old children reading  
65 short stories. However, the adult speech may not always be  
66 available or usable, especially when student makes sub-  
67 stitution or omission. The FLORA group grades overall  
68 literacy for 1 min paragraph reading by children of grades  
69 1–4. They use lexical as well as prosodic features for  
70 giving overall literacy score. Component-wise grading  
71 may, however, prove to be more beneficial in giving  
72 feedback to students.

73 This work is targeted towards children’s oral reading  
74 assessment. Children are known to be highly inconsistent  
75 in articulation. This causes difficulties in modeling their  
76 pronunciations and prosody. Further, these are native  
77 speakers of Marathi with English being the second lan-  
78 guage and the language of instruction in school. The native  
79 language is seen to have notable influence on their English  
80 accent. Since they are new learners of the language, hesi-  
81 tations and filled-pauses are often observed in recordings.  
82 This makes the problem scenario unique. Large training  
83 speech data in the same context may improve the perfor-  
84 mance. This is, however, not possible beyond a certain  
85 extent. The target use case is within the school environ-  
86 ment. This further adds highly variable background noise  
87 that must be handled to ensure system reliability.

88 As observed in the U.S. National Reading Panel report  
89 [4], prosody is neglected during reading assessment in  
90 schools. Only the word reading accuracy and sometimes  
91 reading speed are given importance. The unawareness of  
92 prosody results in variety of intonations and accents. Pro-  
93 sody is also highly affected by the regional accent and  
94 prosody style.

95 We therefore propose a system to assess lexical as well  
96 as prosodic miscues. The lexical miscues are evaluated in  
97 terms of the number of detected insertions, deletions and  
98 substitutions. The prosodic miscues are indicated in terms  
99 of whether the expected words are made prominent and  
100 whether phrase breaks are realized where desired. The  
101 dataset collections and annotation procedure is described in  
102 Sect. 2. The proposed system architecture is explained in  
103 Sect. 3. Finally, discussion on results is included in Sect. 4.

## 104 2 Data acquisition and annotation

105 The data is specifically collected from an urban school  
106 where students are used to English as language of  
107 instruction at school. Students of age group 10–14 years  
108 are considered as target class. These students were asked to  
109 read 10 sentence stories printed on paper. The story texts  
110 were normal English text writings, not annotated in any  
111 way for indicating phrasal breaks or important words. The  
112 readings were recorded through an Android voice recorder  
113 application using a headset microphone. All the recordings

114 were made at 16 kHz sampling frequency. The recordings  
115 are then stored and used for further analysis and assess-  
116 ment. As the surroundings were relatively quiet, the SNR  
117 (Signal to Noise Ratio) is around 20–30 dB. Twenty stu-  
118 dents with good English reading ability are selected. Ten  
119 stories read by each student make a total of 200 recordings.  
120 Table 1 gives detailed statistics of the data used for this  
121 study.

122 The authors manually removed the noisy segments in  
123 the data, if any, and aligned data at sentence level. Gib-  
124 berish and unintelligible words were removed from the  
125 audio and the correct transcriptions were written corre-  
126 sponding to the intelligible words. Prosody labeling is done  
127 by three experts in the form of phrasing and prominence.  
128 For this, they were provided audios with corresponding  
129 correct text transcriptions. These transcriptions did not  
130 include any punctuation marks or capital letters. The raters  
131 marked phrasal breaks (inter-sentence and intra-sentence)  
132 and highlighted those words which were perceived as being  
133 emphasized.

### 134 2.1 Inter-rater agreement

135 Prosody annotation is a subjective task for which inter-rater  
136 agreement is used to estimate the reliability of the ratings.  
137 Inter-rater agreement/reliability is usually indicated in  
138 terms of Fleiss’ Kappa.

139 For prominence prediction task, Fleiss’ Kappa in  
140 vicinity of 0.4 has been reported in [10, 11], while for  
141 phrase boundary estimation, Fleiss’ Kappa in the range  
142 0.54 to 0.62 has been reported in [10]. A similar trend is  
143 observed in our task, where agreement ( $\kappa = 0.3$ ) for  
144 prominence is lower than that for phrasing ( $\kappa = 0.5$ ).

145 Ground truth markings from different annotators differ  
146 at many places. Taking the opinions of different annotators  
147 into account wisely is also important. Different approaches  
148 for determining the ground truth reference like majority  
149 voting or some specified number or all the raters agree  
150 condition [12]. We mark the ground truth for prominence  
151 for any given word as in [13, 14]. If any of the raters marks  
152 a word as prominent, the word will be labeled prominent.

**Table 1** Dataset statistics

Parameter	Rating
Number of students	20
Number of stories	10
Number of recordings	200
Total duration of recordings	2.2 h
Average duration of recordings	40 s

153 The ground truth for phrasal break detection is also  
154 determined in the same way (Table 2).

155 There is large variability in reading abilities among  
156 students in the same school grade. While the word  
157 decoding accuracy is good for most of the students con-  
158 sidered at the given text difficulty level, the prosody is  
159 highly variable. As per the feedback from raters, only a few  
160 students appear to comprehend the text while reading as  
161 revealed by their expressiveness via the prosody. Many  
162 students neglect punctuation marks completely while  
163 reading. Most of the students don't know the importance of  
164 stressing particular words as can be judged from the  
165 cadence or rhythm in their reading. Some students have  
166 tendency to read each and every word accurately, con-  
167 sciously and separately, while some other students read  
168 flatly in monotonous voice giving no stress at all.  
169 Accordingly, we made two groups of students and calcu-  
170 lated Fleiss' kappa values for prominence separately for 11  
171 students with good prosody and 9 students with poor pro-  
172 sody. The agreement values were 0.28 and 0.25 respec-  
173 tively. This indicates that raters agree slightly better for  
174 students with good prosody. This may be due to the  
175 reduced inconsistency between the acoustic cues and  
176 expected top-down syntax cues in good prosody.

### 177 3 System implementation

178 The overall proposed system is as shown in Fig. 1. The  
179 audio recording is passed through a voice activity detector  
180 to obtain sentence-level utterances. The test utterance is  
181 then enhanced for noise suppression using GAN (Genera-  
182 tive Adversarial Network) [15]. The enhanced utterance is  
183 passed through the ASR decoder which converts the speech  
184 into a text hypothesis using previously trained acoustic and  
185 language models. The hypothesis is compared with  
186 canonical text to estimate the speech rate and yield  
187 detected miscues in the form of omission and substitution.  
188 The ASR decoder also provides word-level alignments  
189 which are used to get word-level prosodic features. Pro-  
190 sodic features related to pitch, intensity and spectral bal-  
191 ance contours are used. These features are fed to the  
192 prosodic attribute prediction classifier to get prosody rating  
193 estimates in the form of phrasing and prosody.

**Table 2** Statistics of expert raters' agreement for phrasing and prominence marking

No. of raters	0	1	2	3
No. of prominent words	7982	4197	2982	1328
No. of phrase breaks	7102	3802	3012	2573

### 3.1 Automatic speech recognition 194

For this task, we use state-of the art ASR system discussed 195  
below. 196

#### 3.1.1 Acoustic model 197

The acoustic model is trained on 5 h of children reading <sup>AQ3</sup> 98  
data recorded by fluent English and Hindi speakers. Models 199  
for 47 Hindi and English phones are trained. Hindi phones 200  
are included to take care of native language phones which 201  
appear in our dataset, but are not part of English language. 202  
First speaker normalized features are obtained from MFCC 203  
features using SAT GMM-HMM model. These are passed 204  
through DNN Tandem model to get bottleneck features. 205  
The bottleneck features and speaker normalized features 206  
are appended to get the final feature vectors. A SAT GMM- 207  
HMM model is trained with these feature vectors. Finally 208  
MAP adaptation is performed with some other story 209  
reading data to tune the acoustic model parameters for 210  
current task [16]. The ground truth transcriptions based 211  
word segmentation are obtained through forced-alignments 212  
using these acoustic models. 213

#### 3.1.2 Language model 214

The lexical dictionary used for language model has all the 215  
words observed in the transcribed recordings. Besides, 216  
some other possible variations expected in reading are also 217  
added manually. A language model is built separately for 218  
each story. First a universal garbage model is trained as a 219  
general unigram model. The training is performed on 220  
canonical text of all the available stories. This LM will 221  
have only those words which occurred at least once in the 222  
training text. Other words which did not occur in this text 223  
are also added in the LM as parallel paths. The newly 224  
added unseen words are given the probability same as the 225  
words which occurred only once in the training text. The 226  
probability values for other words are then adjusted 227  
accordingly using Good Turing method for smoothing and 228  
discounting [17]. Now a trigram model is trained on the 229  
canonical text of the target story. This trigram model, with 230  
the garbage model in parallel, forms the final language 231  
model. 232

### 3.2 Prosody evaluation 233

For prosody evaluation, we need to compute certain 234  
acoustic-prosodic features. First the pitch, intensity and 235  
spectral tilt are calculated every 10 ms across the utterance 236  
to get prosodic contours. Statistical measures are calcu- 237  
lated—mean, minimum, maximum, median, span and 238  
standard deviation—for these contours across every word 239

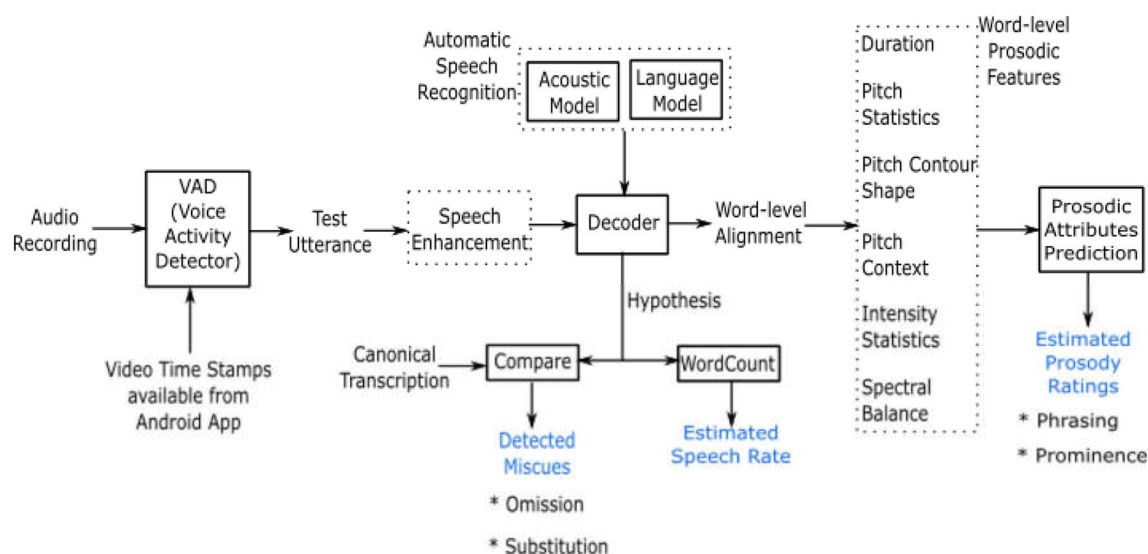


Fig. 1 System block diagram

240 interval. Short time magnitude spectrum is also computed  
 241 for every word with 10 ms window. Then fixed band  
 242 energy contours are computed for four frequency bands—  
 243 0–500 Hz, 500–1000 Hz, 1000–2000 Hz and  
 244 2000–4000 Hz. Total and relative energy are calculated for  
 245 these bands and above referred six statistical measures are  
 246 also computed. Context based features are calculated cor-  
 247 responding to the above features as difference with values  
 248 for previous and successive words. Features indicative of  
 249 pitch contour shape involve correlation with some ideal  
 250 contours—rising, falling, peak, valley and Gaussian with  
 251 five different variance values—0.2, 0.5 1.0, 2.0 and 5.0.  
 252 These ideal contours are of same length as the word  
 253 duration. First four contours are taken such that their  
 254 maximum and minimum also coincide with the word pitch  
 255 contour. Further, duration related features are also com-  
 256 puted which include average, maximum and minimum  
 257 syllable duration for the word and pause/silence durations  
 258 before and after the word.

259 These features are then used for training random forest  
 260 classifiers [14, 18, 19] for phrasal break detection and for  
 261 prominent word detection. The detection decisions for the  
 262 given utterance can be used further to predict ratings for  
 263 high level prosodic attributes and hence to give feedback  
 264 whether an utterance was read “like a sentence” and  
 265 whether the intended important words were stressed and  
 266 intended phrasal breaks were realized.

## 4 Results and discussion

267 Performance of ASR is judged in terms of Word Error Rate  
 268 (WER). With leave one story out based cross validation,  
 269 we obtained WER of 3.44% for the system. In terms of  
 270 children’s reading assessment, number of word miscues is  
 271 considered as the measure of evaluating ASR performance.  
 272 Word miscues include substitution, insertion and deletion.  
 273 The precision-recall values for word miscue detection are  
 274 noted in Table 3.

275 For prosody assessment, 5-fold cross-validation is per-  
 276 formed. We noted individual feature importance in classi-  
 277 fication using Gini criteria [18]. The top important features  
 278 with best Gini values were used for final classification and  
 279 the corresponding results are shown in Table 3 in terms of  
 280 precision-recall for phrase boundary detection and promi-  
 281 nence detection.

282 From manual expert ratings and acoustic analysis of  
 283 audio chunks, we observed and further verified from  
 284 experiments that pitch trends are important in determining  
 285 intonations and signify proper realization of sentence  
 286 endings. Phrasal breaks can be estimated using features like  
 287 pitch contour shape across the word, maximum syllable  
 288 duration in the word and pause/silence duration after the  
 289

Table 3 Results of reading assessment in terms of precision and recall

	Precision (%)	Recall (%)
Word miscue detection	70.4	68.9
Prominent word detection	73.2	73
Phrasal break detection	59.2	80



290 word. The context based features, especially the difference  
291 values with the next word, are found to be important for  
292 phrasal breaks. Pitch span, maximum intensity and maxi-  
293 mum syllable duration are important features in the  
294 prominence detection task.

295 As can be seen from Table 3, precision for phrasal break  
296 estimation is still very low. We also need to work for  
297 further improvement in prominence detection task too.  
298 While analysing the results, we observe that the inter-de-  
299 pendence of prosody attributes is an important aspect that  
300 affects the phrasal break and prominent word detection  
301 performance. e.g. syllable lengthening (indicated by max-  
302 imum syllable duration) and pitch decline on phrase-final  
303 word (indicated by pitch span) are characteristics of proper  
304 phrasing. The same characteristics (duration elongation and  
305 high pitch span) are cues to prominent words as well. This  
306 causes false detection of phrase-final words as prominent  
307 and vice versa. Further, lexical disfluencies like hesitation  
308 and over-extended phones also have longer durations and  
309 lead to prominence marking.

310 The important cases where the proposed system fails,  
311 include the confusion of fricatives with noise by ASR  
312 leading to faulty word-alignment and pitch contour irreg-  
313 ularities. This clearly indicates the importance of accurate  
314 extraction of prosodic features (pitch, energy and duration)  
315 for predicting prosodic events (phrase boundary and  
316 prominent word). Different techniques to extract these  
317 features and different functionals obtained from them can  
318 improve/degrade the performance of prosodic events pre-  
319 diction system.

320 This is an ongoing research work towards automating  
321 the highly reliable and consistent assessment of student's  
322 reading skills on known text. To reach a reliable level of  
323 evaluation, we need to work further on improving both  
324 phrasing as well as prominence predictions. Attempts  
325 towards accurate prosodic feature extraction, use of other  
326 lexical features from ASR and training on larger dataset  
327 can prove to be helpful. Besides, different symmetric and  
328 asymmetric temporal contexts and different normalization  
329 techniques can also be tested. Regional accent related  
330 features may help in further improvements.

331 The long term goals include building a text independent  
332 system to give assessment based on one of the standard  
333 reading fluency scales and evaluating comprehension  
334 ability from the prosody. The system can be developed to  
335 work on other Indian accents as well by adapting to the  
336 desired regional accent. The real world scenario with dif-  
337 ferent noises should also be considered and efforts need to  
338 be made towards building a noise robust system.  
339

## References

1. ASER Centre (2016) ASER: the annual status of education report (Rural India). <http://img.asercentre.org/docs/Publications/ASER>. ASER Centre 341-343
2. Rao P, Swarup P, Pasad A, Tulsiani H, Das G (2016) Automatic assessment of reading with speech recognition technology. In: Proceedings of international conference on computers in education, Mumbai, India 344-347
3. Sabu K, Swarup P, Tulsiani H, Rao P (2017) Automatic assessment of children's L2 reading for accuracy and fluency. In: Proceedings of SLaTE, Stockholm, Sweden 348-349
4. Panel NR (2000) Teaching children to read: an evidence-based assessment of the scientific research literature on reading and its implications for reading instruction. Tech. rep., The Eunice Kennedy Shriver National Institute of Child Health and Human Development 350-355
5. Miller J, Schwanenflugel P (2008) A longitudinal study of the development of reading prosody as a dimension of oral reading fluency in early elementary school children. *Read Res Q* 43(4):336 356-358
6. Breen M, Fedorenko E, Wagner M, Gibson E (2010) Acoustic correlates of information structure. *Lang Cognit Process* 25(7-9):1044 359-361
7. Wang D, Narayanan S (2007) An acoustic measure for word prominence in spontaneous speech. *IEEE Trans Audio Speech Lang Process* 15(2):690 362-363
8. Duong M, Mostow J, Sitaram S (2011) Two methods for assessing oral reading prosody. *ACM Trans Speech Lang Process* 7(4):14.1 364-366
9. Bolanos D, Cole R, Ward W, Tindal G, Schwanenflugel P, Kuhn M (2013) Automatic assessment of expressive oral reading. *Speech Commun* 55(2):221 367-371
10. Mo Y (2008) Acoustic correlates of prosodic prominence for naive listeners of American English. In: Proceedings of Berkeley linguistics society 372-373
11. Kakouros S, Rasanen O (2014) Perception of sentence stress in English infant directed speech. In: Proceedings of INTER-SPEECH, Singapore 374-375
12. Tamburini F, Bertinetto CBP (2014) Prosodic prominence detection in Italian continuous speech using probabilistic graphical models. In: Proceedings of speech prosody, Dublin, Ireland 376-377
13. Kakouros S, Rasanen O (2016) 3PRO: an unsupervised method for the automatic detection of sentence prominence in speech. *Speech Commun* 82:67 378-381
14. Sabu K, Rao P (2018) Detection of prominent words in oral reading by children. In: Proceedings of speech prosody, Poznan, Poland 382-383
15. Sabu K, Nathani K, Rao P (2018) Improving the noise robustness of prominence detection for children's oral reading assessment. In: Proceedings of national conference on communications, Hyderabad, India 384-385
16. Swarup P (2017) Acoustic model training and adaptation for children's read speech recognition. Tech. rep, MTP Report, Department of Electrical Engineering, IIT Bombay 386-388
17. Church K, Gale W (1991) A comparison of the enhanced good-turing and deleted estimation methods for estimating probabilities of English bigrams. *Comput Speech Lang* 5:19 389-391
18. Breiman L (2001) Random forests. *Mach Learn* 45(1):5 392-393
19. Christodoulides G, Avanzi M (2014) An evaluation of machine learning methods for prominence detection in French. In: Proceedings of INTERSPEECH, Singapore 394-395