

Phone Duration Modeling for Speaker Age Estimation in Children

Prashanth Gurunath Shivakumar, *Student Member, IEEE*, Somer Bishop,
Catherine Lord, Shrikanth Narayanan, *Fellow, IEEE*

Abstract—Automatic inference of important paralinguistic information such as age from speech is an important area of research with numerous spoken language technology based applications. Speaker age estimation has applications in enabling personalization and age-appropriate curation of information and content. However, research in speaker age estimation in children is especially challenging due to paucity of relevant speech data representing the developmental spectrum, and the high signal variability especially intra age variability that complicates modeling. Most approaches in children speaker age estimation adopt methods directly from research on adult speech processing. In this paper, we propose features specific to children and focus on speaker’s phone duration as an important biomarker of children’s age. We propose phone duration modeling for predicting age from child’s speech. To enable that, children speech is first forced aligned with the corresponding transcription to derive phone duration distributions. Statistical functionals are computed from phone duration distributions for each phoneme which are in turn used to train regression models to predict speaker age. Two children speech datasets are employed to demonstrate the robustness of phone duration features. We perform age regression experiments on age categories ranging from children studying in kindergarten to grade 10. Experimental results suggest phone durations contain important development-related information of children. Phonemes contributing most to estimation of children speaker age are analyzed and presented.

Index Terms—Phone Duration, Children speech, Age Estimation, Speaker Age Regression

I. INTRODUCTION

SPEECH contains important paralinguistic information including speaker’s age, gender, emotions, and other behavior constructs [1]. Knowledge of such information can help improve spoken language technologies (SLT) such as speech recognition, speaker recognition, and enhance the experience of SLT based applications by providing robustness against variability along those dimensions. Inference of age, gender from children speech can help better tailor conversational interfaces such as education and learning platforms, entertainment, interactive gaming, tutoring, social networking for different age-gender demographics. Speech-based biomarkers are also increasingly used in supporting health applications[2], including related to developmental disorders [3].

Knowledge of an individual’s age is an important meta-data for several applications. Automatic recognition of speaker

age is valuable especially when speech is the only form of data available. Age information enables targeted information dispersal and better personalization, including ensuring privacy and security, thereby enhancing the experiences supported by the speech technology applications. Arguably, child centric applications have more to benefit by using paralinguistic information in enabling novel approaches in safeguarding children and enforcing age appropriate content.

Most of the past research in speaker age estimation is based on adult speech. Earlier research involved training classifiers on statistical functionals of speech descriptors such as loudness, pitch, jitter, shimmer, mel-frequency cepstral coefficients (MFCC) [4, 1]. Gaussian mixture models (GMM) based systems trained on MFCCs have been a popular choice for speaker age prediction [5]. Maximum a-posteriori (MAP) adaptation, discriminative training using maximum mutual information (MMI) have been shown to be successful additions to GMMs [5, 6]. [6] proposed joint factor analysis (JFA) with a GMM back-end for age classification. Later, i-vector with total variability modeling trained on MFCC features significantly advanced the performance of age regression achieving 7.6 years of mean absolute error (MAE) [7]. Supervised i-vectors further improved the performance by decreasing the MAE by a relative 2.4% [8]. Within class covariance normalization (WCCN) was found to be useful both in the case of i-vector and supervised i-vector [9, 8]. Cosine distance scoring is typically used for classification with i-vectors and support vector regression in case of age regression task. [10] reported improvements using i-vectors by adopting shallow artificial neural networks as backend for regression.

More recently, deep neural networks have been employed for speaker age estimation. In [11], the hybrid acoustic DNN-HMM from an automatic speech recognition (ASR) system is used to extract phonetically-aware senone posterior i-vector, instead of the typical GMM-Universal background model (UBM). In [12], bottleneck features are extracted from a hybrid DNN-HMM phone recognition system and subsequently used to train the i-vector, yielding better performance. End-to-end deep neural network architectures have also been explored for age estimation [13, 14, 15]. One such system, popularly termed as x-vectors, comprises several layers of time delay neural network followed by a pooling layer that computes mean and standard deviation over time. The statistics are concatenated and propagated through several feed forward layers to finally output softmax distribution over predefined, binned, age categories. With large amounts of training data or data augmentation, x-vectors have been shown to generally

P. Shivakumar and S. Narayanan are with the Department of Electrical and Computer Engineering, University of Southern California, Los Angeles, California 90089, USA (e-mail: pgurunat@usc.edu;shri@sipi.usc.edu). S. Bishop is with Department of Psychiatry, University of California, San Francisco, USA e-mail:(somer.bishop@ucsf.edu). C. Lord is with the Semel Institute of Neuroscience and Human Behavior, University of California, Los Angeles, USA e-mail:(CLord@mednet.ucla.edu).

outperform i-vectors for age estimation [13]. Recurrent and convolutional neural network architectures have also been explored [14, 16].

Although, there has been interest in automatic recognition of paralinguistic information from speaker data, there has been considerably less research focused on children speech where there is significant age-dependent developmental variability [17, 18]. Most of the work involving children treat them as a broad sub-population group and perform classification across broad age groups such as children, youth, adult and senior adults [5, 4, 6, 9, 15]. [19] proposed GMM supervectors (GMM-UBM) and support vector machines (SVM) for classification and regression of children’s age. [20] proposed fuzzy based strategy to aggregate the output of multiple classifiers each trained using MFCC features pertaining to vowels. Extreme learning machine and SVM was used for classification among children of 6 age classes (7 to 12 years). In [21] and [22], children ranging from age 4 to 14 years are categorized into 3 groups based on their age and classification is performed demonstrating the performance advantage of the i-vector system trained on MFCC features and linear discriminant analysis (LDA) against GMM-UBM, GMM-SVM systems. In [9], age regression is performed among the children sub-population, however the mean absolute error and the correlation in case of children is found to be poor. [23] performed age classification task among children categorizing 4 to 14 years using deep neural network with TDNN-LSTM architecture trained on raw speech waveform. The OGI Kids corpus was employed and data augmentation is performed using amplitude and speed perturbation to increase the training data for DNN.

There are additional challenges involved in handling and modeling children speech which complicate the process of automatic age estimation in children. Collection of children speech data is relatively more expensive. The data scarcity of child speech resources poses additional challenges in statistical modeling. Typical data augmentation techniques such as speech rate, pitch perturbations and inclusion of adult speech data which are effective tools in children speech recognition, may prove less helpful in the case of age estimation task. One of the reasons is because children speech is characterized by age-dependent shifts in overall spectral content and formant frequencies [24]. The inclusion of adult speech data is less likely to aid in the performance improvement of children age estimation because of the wide mismatch of spectral parameters between children and adult speech. Moreover, human perceptual evaluations indicate that speech rate influences speaker age estimation [25]. Faster speech rate is associated with lower age estimates and vice-versa [25]. Error in age estimation is also linked to misclassification of gender [26]. For example, perception of gender as male portrayed tendencies toward lower age estimates [26]. These observations make data augmentation techniques such as speech rate and pitch perturbations unsuitable for the task of automatic age estimation.

From a speech modeling perspective, child speech is relatively more complex with high signal variability due to the developmental changes along various aspects including structural (e.g., vocal tract anatomy), motoric (e.g., speech

related movements), cognitive (e.g., linguistic knowledge) and social (e.g., affect expressions) [17, 18]. High within-speaker variability is observed across all children ages through adulthood [17]. Substantial variation in growth rates of children reflects in substantial variation of vocal tract structure for children of same age and for a specific child at different ages [27] which further complicates modeling of children age from speech. High inter-speaker variability observed across age groups [28, 17] poses further difficulties in estimating efficient within-age class boundaries. It is well documented in the literature that children speech recognition is significantly less accurate [24, 28, 29, 30, 31] underscoring the modeling difficulties associated with children speech.

From a psycho-acoustics perspective, the perception of children’s age is particularly distinct. Humans tend to incorporate assumptions about a child speaker’s gender in estimating child’s age [32, 26, 33]. In general speaker gender inference is relatively poor in case of children compared to adult speakers since speech of both female and male children is characterized with higher F0 values [17] which potentially manifests as errors in age estimation. Humans were found to persistently underestimate age for older girls [33]. These trends pose further challenges in children speaker age estimation.

Our work is motivated from the investigations of variations of temporal parameters in children across age categories [17]. [17] found phoneme durations to be associated with speech development in children. In this study, we propose features derived from phone durations for the task of age estimation from children speech. Manual transcriptions are forced-aligned with speech data to obtain phone duration distributions and are subsequently used to train regression models. Our work remains one of the very few works in children speech domain to employ regression for speaker age estimation. Although, there have been a few past works that incorporate durations in terms of pauses and overall length of utterances in speaker age prediction, our work is distinct in its explicit modeling of phone duration in children speech as a biomarker for speaker age estimation. To the best of our knowledge, this work is a unique attempt at modeling speaker age information purely based on temporal variations in speech by means of phoneme duration modeling.

The rest of the paper is organized as follows: Section II describes the proposed phone-duration features and section II-B presents the regression model. Section III lists the speech databases employed in our study. The experimental setup is described in section IV. The results are presented and discussed in section V. Finally, the study conclusions are provided in section VI and possible future directions discussed.

II. PHONE DURATION FEATURES

Duration is a critical descriptor of speech signals. Varying resolution of duration can convey a variety of information ranging from low level descriptors such as speaking rate and pauses in speech to more abstract information about cognitive process, emotion and conversational dynamics. [17] studied the crucial role of variability of duration in children’s speech. Analysis of durations of 10 vowels and fricative

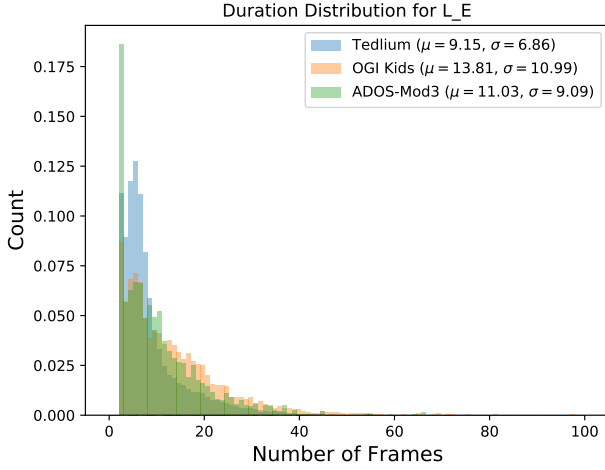


Fig. 1: Phone duration distribution for phoneme /L/ (end position) - Adult (TEDLIUM) vs. Children (OGI Kids & ADOS-Mod3)

portion of /S/ established significant effect of age on duration descriptors. Younger children especially of age 5 and 6 years exhibited significantly longer mean vowel durations compared to older children, with age-dependent duration values reaching a minimum around the age 15 years. Increased intra and inter speaker variation in duration is observed across age groups but a trend which is found to reduce with increasing age. Both inter and intra speaker variation patterns approach adult levels for children of 13 years and above. It was found that younger children tend to exaggerate long vowels including /Y/, /AE/, /AA/ and /ER/. The authors also established that the effect of gender on the duration is not significant.

Few studies have established correlation between mean durations of predefined set of syllables and children speaker’s true (chronological) age as well as human perceived age [32]. Phrase and word durations as well as the inter-word pause durations decrease with increase in age [34]. Phonological studies in children link phoneme durations to developing speech articulatory and neuro-motor timing control in growing children [35, 36, 37]. [17] also found significant correlation between children age and sentence duration.

Psycholinguistic studies have found that speech duration is related to cognition in children, i.e., speech patterns revealed children take longer time to express utterances with higher cognitive demand [38]. Cognitive processes such as selection, retrieval and planning also reflect in temporal speech pause patterns found in children [39].

Moreover, children speech is associated with increased mispronunciations, disfluencies, frequent pauses, non-vocal verbalization [24, 40]. Children’s speech is characterized with repetitions and revisions which are reflective of language development [41]. Phone duration distributions can implicitly encode such speech characteristics found in children and has the ability to capture several para-linguistic patterns including as reflected in speaking rate and stress markers.

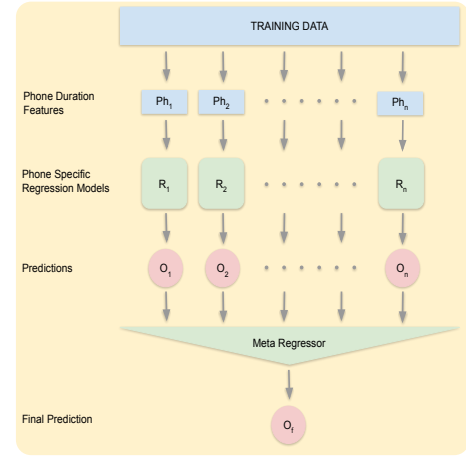


Fig. 2: Proposed Age Regression Model Architecture

A. Proposed Phone-Duration Features

Motivated and supported from the findings in prior literature, in this work, we propose features explicitly engineered to model phone duration distribution in children speech to determine a speaker’s age. First, the speech data is forced-aligned with the manual transcriptions. Later, the temporal occupancy distribution for the following set of phones are computed:

- Position dependent phones: to capture temporal patterns of phones depending on their position in the word (beginning, intermediate or end) or in isolation.
- Position independent phones: obtained from aggregated statistics of position independent phones.
- Lexical stress marked phones: vowels carrying either no stress, primary stress or secondary stress.
- Silence phones: to model the pauses and speaking traits.
- Special phones such as spoken noise: to model and capture hesitancy, disfluencies and filled pauses.
- Global distributions: set of all non-silence phones, consonants and vowels.

Finally, statistical functionals are computed from the duration distributions for each phone, i.e., eight distribution descriptors namely *mean*, *variance*, *minimum*, *maximum*, *skewness*, *kurtosis*, *entropy* and *mean absolute deviation*.

Figure 1 shows an example of duration distribution for the position dependent phoneme L_E, where “_E” indicates the occurrence of the phone /L/ at the end of the word. The figure comprises histograms comparing phone duration distribution for 3 different corpora: one adult (Tedlium) and two children (OGI Kids and ADOS-Mod3; see section III for descriptions). It is evident from the figure, that adult speech is associated with significantly smaller durations compared to children speech. Children phone durations have typically higher means and standard deviations. Such developmental trends with phoneme durations are prevalent among children of different age categories. Figure 3 presents the phone duration distributions of phoneme /T_I/ (I represents intermediate position of phoneme in words) for each age category ranging from kindergarten to children studying in 10th grade. The

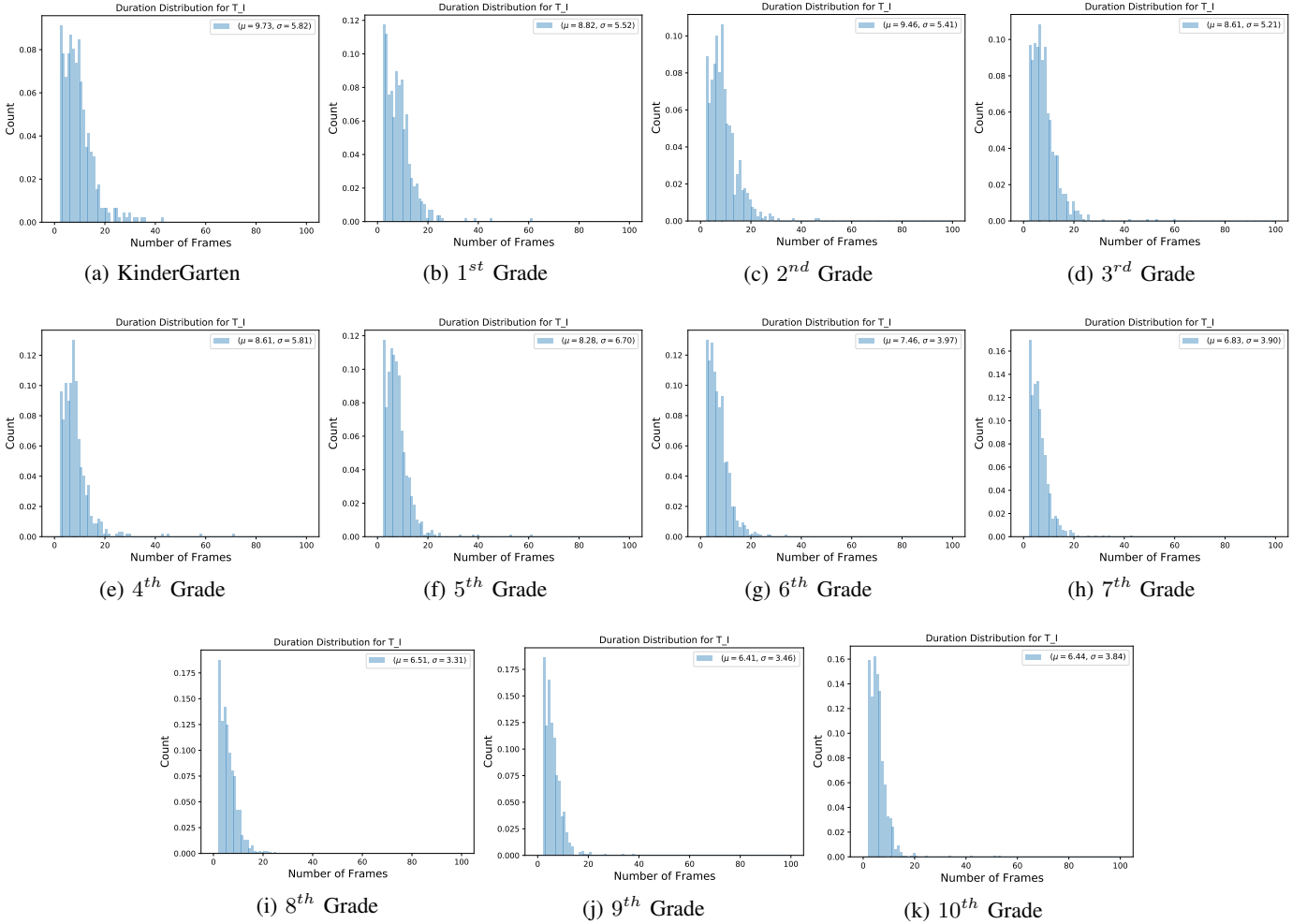


Fig. 3: Phone Duration Distribution for phone T_I over different ages

distribution of all the phoneme durations are found to be Gaussian in nature and unimodal.

B. Proposed Age Regression System

The proposed regression model architecture is illustrated in Figure 2. The architecture is based on stacking two layers of regressors. An ensemble of individual estimators, handling one phoneme each, are trained independently on the eight aforementioned distribution feature descriptors to predict the speaker’s age. A final, meta regressor operates on the outputs of the individual estimators to give the final age prediction. The final estimator is trained on the predictions of individual estimators using cross-validation. In this work, we employ two regression models, support vector regressor and random-forest based AdaBoost regressor. The final, meta regressor model is of the same class as the base individual estimator. The choice of regression models are based on the following factors: (i) the amount of training data available; less data makes DNN based models unsuitable, (ii) support vector based models are popular choice among prior literature, and (iii) decision tree based model for inferring feature importance.

The proposed stacking ensemble learning has certain advantages over a single regression model. Ensemble learning helps

in achieving low bias and low variance in final predictions. The stacked estimators also help handle high feature dimension (2912 features, i.e., 364 phonemes, 8 features each) efficiently in contrast to typical dimension reduction alternatives. The stacked architecture for duration modeling helps in alleviating over-fitting issues since the final estimator is trained on the cross-validated predictions of the base estimators. It enables implicit feature selection among different phonemes, since the meta classifier operates on top of outputs of base estimators pertaining to each phoneme. It also has the added advantage to enable easier assessment of contribution of information provided by duration distribution of underlying phone towards age estimation.

III. CHILDREN SPEECH DATABASES

In this study, we perform experiments on two children speech corpora to assess the transferability and robustness to different domains and acoustic conditions.

A. OGI Kids Speech Corpus

We employ OGI Kids speech corpus [42] as the primary database for our experiments due to its wide age distribution

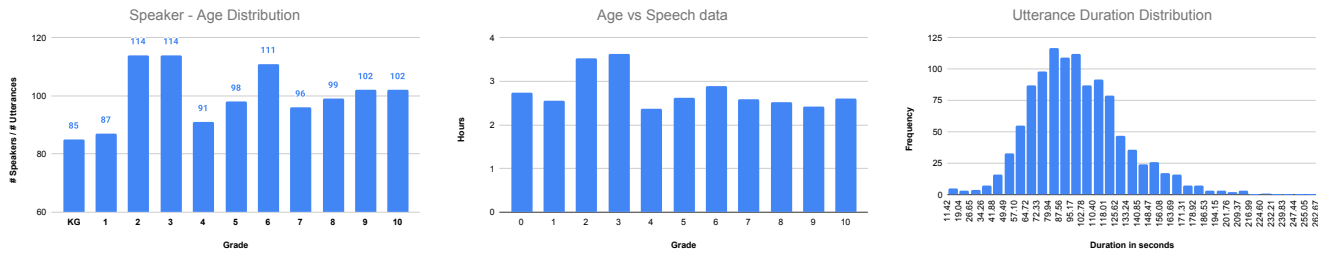


Fig. 4: OGI Kids Corpus Statistics

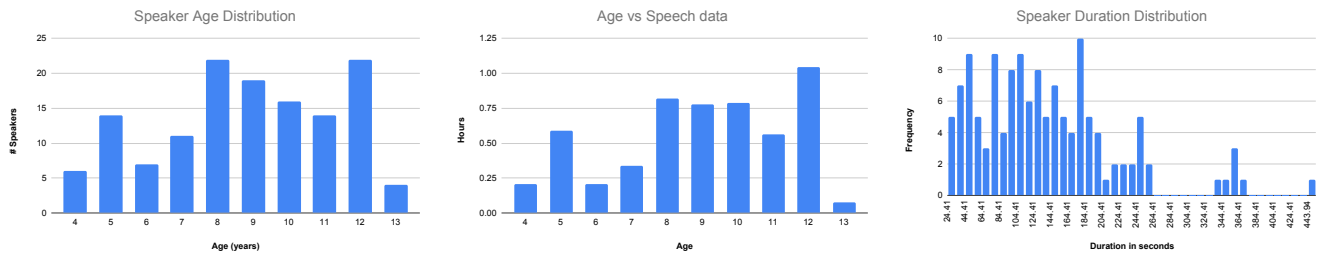


Fig. 5: ADOS-Mod3 Data Statistics

demographics among children. We make use of the spontaneous speech subset of the corpus comprising adult interviewer asking a series of questions and eliciting a spontaneous response from children. The corpus consists of 1100 distinct children speakers with ages ranging from children studying in kindergarten to 10th grade. It includes a total of approximately 30.5 hours of speech. Each speaker utters a single sentence and the mean duration per utterance is approximately 100 seconds. The speaker age distribution, amount of speech data and utterance duration distribution statistics are presented in figure 4.

B. Autism Diagnostic Observation Schedule - Module 3 (ADOS-Mod3)

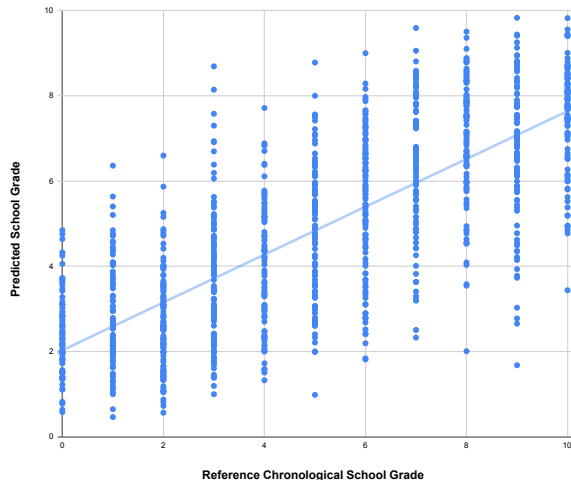
The ADOS-Mod3 corpus [43] comprises child-adult dyadic conversations involving semi-structured, standardized assessment of communication and social interactions. The children in the corpus are diagnosed with autism spectrum disorder (ASD), attention deficit hyperactivity disorder (ADHD) and various other developmental disorders including language disorder. The speech sessions were collected at two different locations including University of Michigan Autism and Communication Disorder Center and Cincinnati Children’s Medical Center. We make use of speech data from children only for speaker age estimation and omit adult speech data. The corpus consists of 179 children, out of which we consider a subset of 135 children for whom we had generated good quality automatic phonetic alignments (see section IV for details regarding alignments) for the age regression analysis of this study. The age of children range from 43 to 158 months (4 to 13 years). The corpus contains a total of 5.4 hours of manually-transcribed speech. Each speaker has a mean duration of approximately 144 seconds. The speaker age distribution, amount of speech data per age and speaker durations distribution are presented in Figure 5.

Several factors associated with this dataset add additional complexity for the task of speaker age estimation. First, the differences in the neuro-developmental condition due to ASD, ADHD and other developmental disorders possibly reflected in the speech complicates the age estimation from speech [44]. Second, the speech data are recorded in far-field conditions with a single distant microphone leading contributing to acoustic variability. Third, the data analyzed are from two different locations with different room and channel characteristics adding to the complexity of speech modeling. Finally, the corpus consists significantly less data (17%) compared to OGI speech corpus. The above challenges help us evaluate the robustness of the proposed phone duration model.

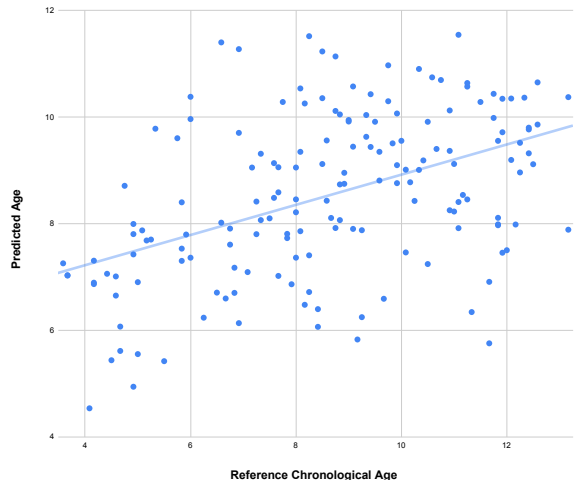
IV. EXPERIMENTAL SETUP

In this section, we provide details of our experimental setup for reproducibility purposes. For forced-alignment we employ the KALDI speech recognition toolkit [45]. The feature pipeline consists of extracting 13-dimensional mel-filter cepstral coefficients (MFCC) using a window size of 25ms and a shift of 10ms. Linear discriminant analysis (LDA) is performed on top of the MFCC features by considering left and right context of 3 frames. Furthermore, Maximum likelihood linear transform is applied on top of LDA features. Finally, feature-space maximum likelihood linear regression (fMLLR) based speaker adaptive training is used to train a Gaussian mixture model - hidden markov model (GMM-HMM) based acoustic model. The resulting acoustic model is used for forced-alignment to obtain phoneme level alignments. Later, the statistics are accumulated for each phoneme under consideration and their 8 functional descriptors namely *mean*, *variance*, *minimum*, *maximum*, *skewness*, *kurtosis*, *entropy* and *mean absolute deviation* are computed.

Two separate acoustic models are trained for forced alignments, one for the OGI Kids Corpus and the other for the



(a) OGI Kids



(b) ADOS-Mod3

Fig. 6: Age Regression Scatter Plot

ADOS-Mod3 corpus. Since both OGI Kids and ADOS-Mod3 corpus are fairly small, we include additional speech data in favor of constructing better acoustic models and thereby obtaining better quality alignments. The acoustic model employed for OGI Kids is trained by including 198 hours of My Science Tutor (MyST) children speech corpus [46] involving children studying in grades 3, 4 and 5. The MyST corpus comprises conversational style speech of children recorded under low noise and close talk conditions similar to OGI Kids. The acoustic model employed for ADOS-Mod3 is trained by including all the 173 children with the addition of adult speech of clinicians conducting the Autism diagnostic assessment. The addition of adult speech is known to yield better quality of acoustic models under low data scenarios [29]. We do not include the MyST data since we believe the addition of adult speech data under similar recording conditions i.e., far-field, high reverberation environment in case of ADOS-Mod3 is more beneficial. The total number of phones in the OGI Kids corpus is 364. Whereas, the number of phones in case of ADOS-Mod3 corpus is restricted to 185 phones (excludes lexical stress markers) to better handle the smaller size of the training corpus.

For the age estimation task, we directly perform regression to predict the (reference) school grade of children (as proxy for age) in case of OGI Kids Corpus. In case of ADOS-Mod3, the age of children were converted from months to years. The performances are hence directly comparable between the two models. In this work, we experiment with two regression models, i.e., support vector machine regressor (SVR) and the decision tree based random forest AdaBoost regressor. The choice of SVR is due to its popularity and proven effectiveness for age estimation in prior works. Whereas, the decision tree based AdaBoost model has the benefit of providing feature importance which helps us analyze the contributions of phonemes and their discriminative power towards age estimation. Given the small size of the speech corpora, we perform leave-one-

speaker-out (LOSO) cross-validation. The hyper-parameter tuning of the regression models are handled implicitly through nested cross-validation. For performance evaluation, we report mean absolute error (MAE), R^2 score and Pearson correlation. In this work, we employ a simple baseline model that predicts the mean of the age of the speakers for comparison purposes.

V. RESULTS

Table I presents the results of children speaker age estimation on OGI Kids and ADOS-Mod3 through regression. Our proposed phone duration model achieves a mean absolute error of 1.62 and a correlation of 0.76 with SVR on OGI Kids corpus. The results are significantly better than the baseline system based on mean age prediction. Moreover, the correlation results are comparable to correlation between human perceived speaker age and the chronological age which is believed to be approximately 0.7 [33]. We observe that the SVR model outperforms the AdaBoost model.

The results with ADOS-Mod3 are slightly worse compared to OGI Kids. This is expected due to two factors: (i) the neuro-developmental disorders can have an impact on cognitive development, and reflected in speech production differences; this in turn can lead to differences in speaker age perception and predictions from the acoustic speech signal, and (ii) the ADOS-Mod3 corpus (5.4 hours) has significantly less data compared to OGI Kids (30.5 hours). However, the results

Database	Model	MAE	R^2 score	Correlation
OGI Kids	Baseline	2.69	0.0	0.0
	SVR (RBF)	1.62	0.58	0.76
	AdaBoost	1.82	0.48	0.71
ADOS-Mod3	Baseline	2.07	0.0	0.0
	SVR (RBF)	1.79	0.24	0.49
	AdaBoost	1.74	0.29	0.54

TABLE I: Results: Children Speaker Age Estimation

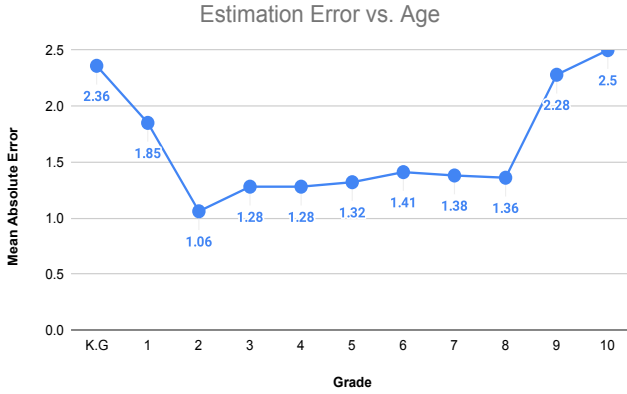


Fig. 7: MAE across different age categories - OGI Kids Corpus

are significantly better than the mean baseline. In the case of ADOS-Mod3, we observe that the AdaBoost regression model outperforms the SVR. The results on ADOS-Mod3 corpus further attests to the robustness of the proposed phone duration modeling. Figures 6a and 6b illustrate scatter plots of chronological age versus the predicted speaker age on the OGI Kids and ADOS-Mod3 corpora, respectively. Note, the scatter plots for OGI Kids are quantized to the school grade of children, whereas the ADOS-Mod3 has age in terms of months. Overall, the results suggest that the proposed phone duration features contain developmental information in children speech. This emphasizes that speaker age in children can be robustly estimated by modeling the temporal variation via phone duration distributions.

Next, we perform additional analysis to assess the factor of age on the performance of the system. Figure 7 plots the mean absolute error in each age category of OGI kids corpus using the SVR model. From the results, we observe that the error is low for children ages ranging from 2nd grade to 8th grade categories, reaching minimum for the 2nd grade group. However, the error increases sharply for younger (kindergarten and 1st grade) and older children in the corpus (9th - 10th grade). One possible explanation for the observed trend is as follows: (i) in case of younger children studying in kindergarten and 1st grade, although the inter and intra age variations are expected to be very high, the intra-age variation dominates, thereby resulting in high error, (ii) in case of children studying among 2nd and 8th grade, the inter-age variation dominates resulting in lower error rates, and (iii) in case of elder children (9th and 10th grade) both the intra and inter age variations are significantly less which results in relatively low performance (comparable to that of adults).

We also perform feature importance analysis to gain insights on the contributions of each phoneme toward children speaker age estimation. We derive impurity based feature importance that are accessible from tree-based algorithms, in our case the random forest based AdaBoost model trained on OGI Kids corpus. The importance measures are computed based on total reduction of the optimization criterion, often referred to as the Gini importance. We compute the feature importance only on the final, meta estimator that operates on the output of the

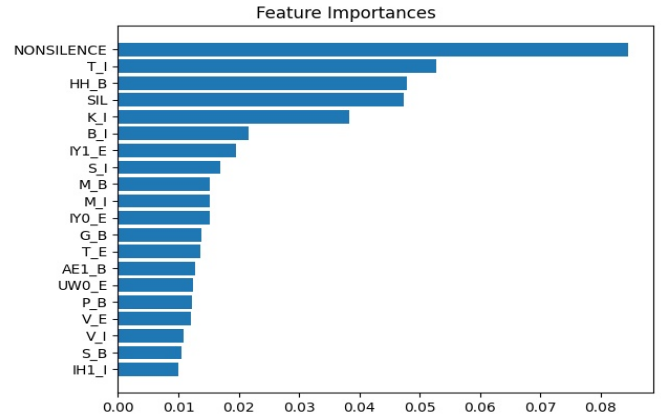


Fig. 8: Phoneme-wise feature importance - OGI Kids Corpus

phoneme specific base estimators. This allows us to assess the contribution of input features on the phoneme level rather than statistical functionals. Figure 8 shows the bar plot of the top-20 most contributing phonemes computed on OGI Kids corpus. Higher values translate to more importance. The following observations are made with our experimental setup:

- 1) /NONSILENCE/ (aggregated duration of *all* position dependent non-silent phones) comes out as the most critical contributor for speaker age estimation.
- 2) Position dependent phones /T_I/, /HH_B/ and /K_I/ contribute critically to determining speaker age in children.
- 3) /SIL/ (silence) capturing inter-word pauses, speaking rate, and disfluencies also helps in determining children age.
- 4) The above 5 acoustic categories are more than twice as important than the durations of other categories in the speech acoustic inventory.
- 5) The appearance of position dependent phones among the top-20 indicates that duration of phones appearing in different parts of a word carry discriminating information on children growth.

VI. CONCLUSION & FUTURE WORK

In this work, we investigate features solely based on phone durations in speech (i.e., acoustic realizations of phonemes) for the task of speaker age estimation in children speech. Phoneme occupancy distribution is derived by forced-aligning manual transcripts with speech signal. Statistical functionals describing the distributions are extracted for each phone which constitutes the features. A double layer stacking regressor architecture is employed with a meta estimator operating on top of multiple base estimators, each trained on statistical functional features corresponding to each phoneme. The results suggests that phone durations contain critical developmental information helpful in predicting speaker age among children. The results indicate that a speaker's age among children can be effectively predicted by looking only at the temporal variations in speech signal. The best performing phone duration model yields mean absolute error of 1.62 and a correlation of 0.76. The estimation of speaker age among

children is associated with high error among young and older children, while yielding minimum estimation error among children studying among 2nd grade to 8th grade. We find that aggregated phone durations of non-silence phones is the most important feature. Among the other phonemes, particularly /T_I/, /HH_B/ and /K_I/ play important role. We also find that inter-speech silence duration also play an important role in predicting child speaker age. Subsequent experiments on additional speech corpora, ADOS-Mod3, comprising speech data from children with ASD/ADHD diagnosis, further underscores the robustness of the phone duration features.

In the future, we plan to combine the phone duration features along with other speech based features including spectral features such as MFCC and voice quality features such as jitter, shimmer to explore complementary information for improved age estimation. Additionally, combination of phone duration features and unsupervised total variability modeling based i-vectors or the supervised features derived through DNNs such as x-vectors can potentially complement and improve performance especially under low data scenarios. We would also like to explore scenarios when manual speech transcripts are unavailable and thus alignments derived from an ASR is the only option, i.e., exploring the effect of automated transcriptions on the performance of speaker age estimation.

REFERENCES

- [1] B. Schuller, S. Steidl, A. Batliner, F. Burkhardt, L. Devillers, C. Müller, and S. S. Narayanan, "Paralinguistics in speech and language—state-of-the-art and the challenge," *Computer, Speech, and Language*, vol. 27, no. 1, pp. 4–39, Jan 2013.
- [2] D. Bone, C.-C. Lee, T. Chaspari, J. Gibson, and S. Narayanan, "Signal processing and machine learning for mental health research and clinical applications," *IEEE Signal Processing Magazine*, vol. 34, no. 5, pp. 189–196, September 2017.
- [3] D. Bone, T. Chaspari, and S. Narayanan, "Behavioral signal processing and autism: Learning from multimodal behavioral signals," in *Autism Imaging and Devices*. CRC Press, 2017, pp. 335–360. [Online]. Available: <https://www.taylorfrancis.com/books/e/9781315371375/chapters/10.1201/9781315371375-21>
- [4] B. Schuller, S. Steidl, A. Batliner, F. Burkhardt, L. Devillers, C. Müller, and S. S. Narayanan, "The interspeech 2010 paralinguistic challenge," in *Eleventh Annual Conference of the International Speech Communication Association*, 2010.
- [5] M. Li, K. J. Han, and S. Narayanan, "Automatic speaker age and gender recognition using acoustic and prosodic level information fusion," *Computer Speech & Language*, vol. 27, no. 1, pp. 151–167, 2013.
- [6] M. Kockmann, L. Burget, and J. Černocký, "Brno university of technology system for interspeech 2010 paralinguistic challenge," in *Eleventh Annual Conference of the International Speech Communication Association*, 2010.
- [7] M. H. Bahari, M. McLaren, D. A. van Leeuwen *et al.*, "Speaker age estimation using i-vectors," *Engineering Applications of Artificial Intelligence*, vol. 34, pp. 99–108, 2014.
- [8] P. G. Shivakumar, M. Li, V. Dhandhanian, and S. S. Narayanan, "Simplified and supervised i-vector modeling for speaker age regression," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 4833–4837.
- [9] J. Grzybowska and S. Kacprzak, "Speaker age classification and regression using i-vectors." in *INTERSPEECH*, 2016, pp. 1402–1406.
- [10] A. Fedorova, O. Glembek, T. Kinnunen, and P. Matějka, "Exploring ann back-ends for i-vector based speaker age estimation," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [11] S. O. Sadjadi, S. Ganapathy, and J. W. Pelecanos, "Speaker age estimation on conversational telephone speech using senone posterior based i-vectors," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5040–5044.
- [12] A. A. Mallouh, Z. Qawaqneh, and B. D. Barkana, "New transformed features generated by deep bottleneck extractor and a gmm-ubm classifier for speaker age and gender classification," *Neural Computing and Applications*, vol. 30, no. 8, pp. 2581–2593, 2018.
- [13] P. Ghahremani, P. S. Nidadavolu, N. Chen, J. Villalba, D. Povey, S. Khudanpur, and N. Dehak, "End-to-end deep neural network age estimation." in *INTERSPEECH*, 2018, pp. 277–281.
- [14] R. Zazo, P. S. Nidadavolu, N. Chen, J. Gonzalez-Rodriguez, and N. Dehak, "Age estimation in short speech utterances based on lstm recurrent neural networks," *IEEE Access*, vol. 6, pp. 22 524–22 530, 2018.
- [15] Z. Qawaqneh, A. A. Mallouh, and B. D. Barkana, "Dnn-based models for speaker age and gender classification," in *International Conference on Bio-inspired Systems and Signal Processing*, vol. 5, 2017, pp. 106–111.
- [16] H. A. Sánchez-Hevia, R. Gil-Pita, M. Utrilla-Manso, and M. Rosa-Zurera, "Convolutional-recurrent neural network for age and gender prediction from speech," in *2019 Signal Processing Symposium (SPSymposium)*. IEEE, 2019, pp. 242–245.
- [17] S. Lee, A. Potamianos, and S. Narayanan, "Acoustics of children's speech: Developmental changes of temporal and spectral parameters," *The Journal of the Acoustical Society of America*, vol. 105, no. 3, pp. 1455–1468, 1999.
- [18] —, "Developmental acoustic study of american english diphthongs," *The Journal of the Acoustical Society of America*, vol. 136, no. 4, pp. 1880–1894, 2014.
- [19] T. Bocklet, A. Maier, and E. Nöth, "Age determination of children in preschool and primary school age with gmm-based supervectors and support vector machines/regression," in *International Conference on Text, Speech and Dialogue*. Springer, 2008, pp. 253–260.
- [20] S. M. Mirhassani, A. Zourmand, and H.-N. Ting, "Age estimation based on children's voice: a fuzzy-based decision fusion strategy," *The Scientific World Journal*, 2014.
- [21] S. Safavi, M. Russell, and P. Jančovič, "Identification

- of age-group from children’s speech by computers and humans,” in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [22] —, “Automatic speaker, age-group and gender identification from children’s speech,” *Computer Speech & Language*, vol. 50, pp. 141–156, 2018.
- [23] M. Sarma, K. K. Sarma, and N. K. Goel, “Children’s age and gender recognition from raw speech waveform using dnn,” in *Advances in Intelligent Computing and Communication*. Springer, 2020, pp. 1–9.
- [24] A. Potamianos and S. Narayanan, “Robust recognition of children’s speech,” *IEEE Transactions on speech and audio processing*, vol. 11, no. 6, pp. 603–616, 2003.
- [25] S. Skoog Waller, M. Eriksson, and P. Sörqvist, “Can you hear my age? influences of speech rate and speech spontaneity on estimation of speaker age,” *Frontiers in psychology*, vol. 6, p. 978, 2015.
- [26] P. F. Assmann, M. R. Kapolowicz, D. A. Massey, S. Barreda, and T. M. Nearey, “Links between the perception of speaker age and sex in children’s voices,” *The Journal of the Acoustical Society of America*, vol. 138, no. 3, pp. 1811–1811, 2015.
- [27] H. K. Vorperian, S. Wang, M. K. Chung, E. M. Schimek, R. B. Durtschi, R. D. Kent, A. J. Ziegert, and L. R. Gentry, “Anatomic development of the oral and pharyngeal portions of the vocal tract: An imaging study,” *The Journal of the Acoustical Society of America*, vol. 125, no. 3, pp. 1666–1678, 2009.
- [28] M. Gerosa, D. Giuliani, S. Narayanan, and A. Potamianos, “A review of asr technologies for children’s speech,” in *Proceedings of the 2nd Workshop on Child, Computer and Interaction*, 2009, pp. 1–8.
- [29] P. G. Shivakumar and P. Georgiou, “Transfer learning from adult to children for speech recognition: Evaluation, analysis and recommendations,” *Computer speech & language*, vol. 63, p. 101077, 2020.
- [30] P. G. Shivakumar, A. Potamianos, S. Lee, and S. S. Narayanan, “Improving speech recognition for children using acoustic adaptation and pronunciation modeling,” in *WOCCI*, 2014, pp. 15–19.
- [31] P. G. Shivakumar and S. Narayanan, “End-to-end neural systems for automatic children speech recognition: An empirical study,” *arXiv preprint arXiv:2102.09918*, 2021.
- [32] S. Barreda and P. F. Assmann, “Modeling the perception of children’s age from speech acoustics,” *The Journal of the Acoustical Society of America*, vol. 143, no. 5, pp. EL361–EL366, 2018.
- [33] P. Assmann, S. Barreda, and T. Nearey, “Perception of speaker age in children’s voices,” in *Proceedings of Meetings on Acoustics ICA2013*, vol. 19, no. 1. Acoustical Society of America, 2013, p. 060059.
- [34] L. Singh, P. Shantisudha, and N. C. Singh, “Developmental patterns of speech production in children,” *Applied acoustics*, vol. 68, no. 3, pp. 260–269, 2007.
- [35] B. L. Smith, “Temporal aspects of english speech production: A developmental perspective,” *Journal of Phonetics*, vol. 6, no. 1, pp. 37–67, 1978.
- [36] R. D. Kent and L. L. Forner, “Speech segment durations in sentence recitations by children and adults,” *Journal of phonetics*, vol. 8, no. 2, pp. 157–168, 1980.
- [37] R. D. Kent, “Anatomical and neuromuscular maturation of the speech mechanism: Evidence from acoustic studies,” *Journal of speech and hearing Research*, vol. 19, no. 3, pp. 421–447, 1976.
- [38] J. Dillon, “Cognitive complexity and duration of classroom speech,” *Instructional Science*, vol. 12, no. 1, pp. 59–66, 1983.
- [39] A. Esposito, M. Marinaro, and G. Palombo, “Children speech pauses as markers of different discourse structures and utterance information content,” in *Proceedings of the International Conference: From sound to sense*, vol. 50, 2004, pp. 10–13.
- [40] A. Potamianos, S. Narayanan, and S. Lee, “Automatic speech recognition for children,” in *Fifth European Conference on Speech Communication and Technology*, 1997.
- [41] T. M. Gallagher, “Revision behaviors in the speech of normal children developing language,” *Journal of Speech and Hearing Research*, vol. 20, no. 2, pp. 303–318, 1977.
- [42] K. Shobaki, J.-P. Hosom, and R. A. Cole, “The ogi kids’ speech corpus and recognizers,” in *Sixth International Conference on Spoken Language Processing*, 2000.
- [43] C. Lord, S. Risi, L. Lambrecht, E. H. Cook, B. L. Leventhal, P. C. DiLavore, A. Pickles, and M. Rutter, “The autism diagnostic observation schedule—generic: A standard measure of social and communication deficits associated with the spectrum of autism,” *Journal of autism and developmental disorders*, vol. 30, no. 3, pp. 205–223, 2000.
- [44] D. K. Oller, P. Niyogi, S. Gray, J. A. Richards, J. Gilkerson, D. Xu, U. Yapanel, and S. F. Warren, “Automated vocal analysis of naturalistic recordings from children with autism, language delay, and typical development,” *Proceedings of the National Academy of Sciences*, vol. 107, no. 30, pp. 13 354–13 359, 2010.
- [45] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, “The kaldi speech recognition toolkit,” in *IEEE 2011 workshop on automatic speech recognition and understanding*, no. CONF. IEEE Signal Processing Society, 2011.
- [46] W. Ward, R. Cole, D. Bolanos, C. Buchenroth-Martin, E. Svirsky, S. V. Vuuren, T. Weston, J. Zheng, and L. Becker, “My science tutor: A conversational multimedia virtual tutor for elementary school science,” *ACM Transactions on Speech and Language Processing (TSLP)*, vol. 7, no. 4, pp. 1–29, 2011.