

Talker change detection: A comparison of human and machine performance

Neeraj Kumar Sharma,^{1, a)} Shobhana Ganesh,² Sriram Ganapathy,² and Lori L. Holt¹

¹*Department of Psychology, Carnegie Mellon University, 5000 Forbes Avenue,
Pittsburgh 15213, USA*

²*Department of Electrical Engineering, CV Raman Road, Indian Institute of Science,
Bangalore 560012, India*

(Dated: 21 December 2018)

1 The automatic analysis of conversational audio remains difficult, in part due to the
2 presence of multiple talkers speaking in turns, often with significant intonation vari-
3 ations and overlapping speech. The majority of prior work on psychoacoustic speech
4 analysis and system design has focused on single-talker speech, or multi-talker speech
5 with overlapping talkers (for example, the cocktail party effect). There has been much
6 less focus on how listeners detect a change in talker or in probing the acoustic fea-
7 tures significant in characterizing a talker’s voice in conversational speech. This study
8 examines human talker change detection (TCD) in multi-party speech utterances us-
9 ing a novel behavioral paradigm in which listeners indicate the moment of perceived
10 talker change. Human reaction times in this task can be well-estimated by a model
11 of the acoustic feature distance among speech segments before and after a change in
12 talker, with estimation improving for models incorporating longer durations of speech
13 prior to a talker change. Further, human performance is superior to several on-line
14 and off-line state-of-the-art machine TCD systems.

^{a)} nsharma2@andrew.cmu.edu

15 I. INTRODUCTION

16 Everyday speech communication involves more than extracting a linguistic message¹. Lis-
17 teners also track paralinguistic indexical information in speech signals, such as talker iden-
18 tity, dialect, and emotional state². Indeed, in natural speech communication, linguistic and
19 indexical information are likely to interact since conversations typically involve multiple
20 talkers who take turns of arbitrary duration, with gaps on the order of only 200 ms³. On the
21 listener’s side, the perception of conversational speech demands quick perception of talker
22 changes to support communication.

23 Perceptual learning of talker identity enhances speech intelligibility in both quiet⁴ and
24 acoustically-cluttered environments^{5,6}. This suggests that sensitivity to talker attributes
25 affects speech recognition both in clear speech and under adverse listening conditions. Fur-
26 ther, talker dependent adaptability in perception can be induced from exposure to just a few
27 sentences⁷. These benefits hint at listeners’ ability to track talkers in conversational speech,
28 even in the absence of visual or spatial cues.

29 Detecting a change in talker would seem to rely upon an ability to track regularities in the
30 perceived features specific to a voice, and to detect changes from these upon a talker change.
31 Lavner et al. (2009)⁸ suggest that the talkers are identified by a distinct group of acoustic
32 features. Yet, Sell et al. (2015)⁹ argue that a combination of vocal source, vocal tract, and
33 spectro-temporal receptive field¹⁰ features fail to explain perceived talker discrimination in a
34 listening test with simple single-word utterances. In a similar way, Kimberly et al. (2003)¹¹
35 have described inattention to talker changes in the context of listening for comprehension as

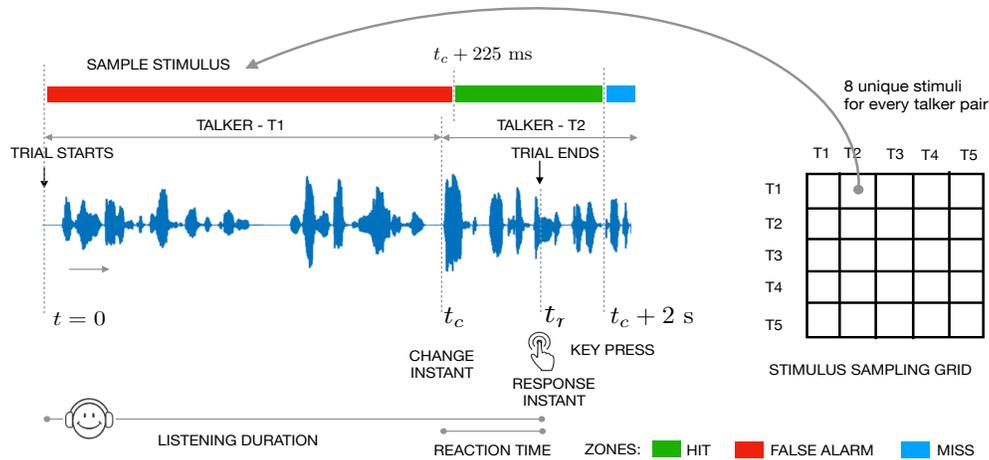


FIG. 1. (color online) Illustration of the proposed talker change detection (TCD) paradigm used in the present listening test study.

36 a form of talker change deafness¹². They suggest that voice information is not continuously
 37 monitored at a fine-grain level of acoustic representation, and conversational expectations
 38 may shape the way listeners direct attention to voice characteristics and perceive differences
 39 in voice. In fact, Neuhoff et al. (2014)¹³ found improved voice change detection when the
 40 language is unfamiliar to the listener, suggesting that there may be interactions between
 41 linguistic and indexical information.

42 Acknowledging that conversational demands¹⁴ in natural speech will often shift attention
 43 toward acoustic features that signal linguistic rather than indexical information, listeners'
 44 ability to detect talker changes does suggest that they track the variability in acoustics fea-
 45 tures associated with a talker's voice. Yet, despite the importance of indexical characteristics
 46 of speech to communication, quite little is known about the nature of the detailed acoustic
 47 features across which talkers differ, the distributions of information characterizing different
 48 talkers along these acoustic features^{15,16}, and the listeners' ability to detect a change in

49 talker. This is especially true for fluent, connected speech, as opposed to isolated words.
50 In this paper, we aim to advance understanding of the information human listeners use to
51 track the change in talker in continuous multi-party speech. We first develop and test a novel
52 experimental paradigm to examine human talker change detection (TCD) performance. We
53 next model the listeners' reaction time (RT) to respond to a talker change in relationship to
54 multiple acoustic features as a means of characterizing the acoustic feature space that listen-
55 ers may track in a voice. We then relate these human perceptual results with performance
56 of state-of-the-art on-line and off-line machine systems implementing TCD.

57 **A. Reaction time as a measure of TCD**

58 We developed a novel paradigm with the goal of obtaining a continuous behavioral measure
59 of listeners' ability to detect a change in talker across relatively fluent, continuous speech.
60 Each stimulus was composed of two concatenated 9 – 14 s utterances sourced from audio
61 books, and spoken by either a single male talker or two different male talkers. The utterances
62 were always drawn from different stories, or parts of a story, so that semantic continuity
63 did not provide a clue to talker continuity. Listeners responded with a button press upon
64 detecting a talker change, thus providing a continuous reaction time measure of how much of
65 an acoustic sample was needed to detect a change in talker. Figure 1 provides an illustration
66 of the paradigm. To the best of our knowledge, this is the first application of a RT change-
67 detection approach to examine human TCD performance.

68 Past studies have used RT to analyze perception of simpler acoustic attributes. For
69 example, studies of tone onset detection¹⁷ and broadband sound onset^{18,19} have reported

70 an inverse relationship between RT and stimulus loudness / spectral bandwidth. Studies
71 guiding the design of warning sounds^{20,21} haven shown faster detection of natural, compared
72 to synthetic, stimuli. Particularly relevant to the present study, recent research characterizes
73 human listeners’ ability to track distributional noise statistics in an audio stream by asking
74 listeners to report a change in statistics with a button press²². Detection of a change was
75 facilitated when listeners heard longer noise samples. This prior work illustrates the promise
76 of the RT to detect a change as a means by which to examine how listeners build a model
77 of incoming sound regularities, although speech signals are inherently more non-stationary
78 and the distributional acoustic information that contributes to talker identity is not well
79 understood.

80 **B. Comparing machine and human performance in TCD**

81 With an increasing repository of conversational audio data^{23,24}, automatic detection of talker
82 change is considered to be an essential preprocessing in machine speech recognition²⁵. For
83 example, the availability of timestamps corresponding to talker change instants in a record-
84 ing benefits both speaker identification²⁶ and speech recognition²⁷ tasks. Recent results
85 from the machine systems literature on TCD²⁸⁻³² suggest that talker change detection is
86 difficult. The state-of-the-art in TCD can be categorized into two approaches: metric-based
87 and classification-based. The metric-based approach relies on computing a distance, in some
88 feature space, between successive short-time segments (such as 25 ms segments, with suc-
89 cessive shifts 10 ms) of a speech recording. A talker change is flagged upon detection of a
90 distance that exceeds a preset threshold. The accuracy of TCD in metric-based approaches

91 is dependent on the choice of features (such as vocal tract features^{33,34} and vocal source
92 features^{30,35–37}), segment duration (usually $> 3 - 5$ s), and the distance metric (such as
93 likelihood ratio³⁸, Bayesian information criterion (BIC)^{33,39}, or improved BIC⁴⁰). An alter-
94 native is the classification-based approach whereby a binary classifier is trained to classify
95 short-time speech segments as talker change or no talker change. The accuracy of the ap-
96 proach is dependent on the feature space, the complexity of the classifier (such as support
97 vector machines⁴¹, neural networks³⁰, and deep neural networks^{31,42}), and the amount of
98 training data. In the present work, we evaluate the performance of a few state-of-the-art
99 machine systems across the same stimuli and with the same performance metrics as used in
100 the human perceptual paradigm. These comparisons may help in identifying whether there
101 are performance gaps in human and machine TCD across fairly continuous, fluent speech.

102 II. METHODS

103 A. Participants

104 A total of 17 participants in the age group of 18 – 36 (university students and one staff
105 member, median age 26) took part in the listening test. All listeners reported normal
106 hearing with good fluency in speaking and understanding English. All participants provided
107 informed consent to participate in the study, and the study protocol was approved by the
108 Carnegie Mellon University Institutional Review Board and the Indian Institute of Science
109 Review Board.

110 **B. Stimuli**

111 The stimuli were composed of concatenated utterances from two talkers, talker T_x and talker
 112 T_y . The utterances were taken from audio books drawn from the LibriSpeech corpus⁴³, a
 113 public-domain corpus of about 1000 hours of audio data by approximately 1000 distinct
 114 talkers. Each trial involved a stimulus that was a concatenation of two utterances drawn
 115 from the same, or two different, talkers. The utterances corresponded to sentences read out in
 116 the audio book, and featured natural speech intonation and a rise and fall in speech envelope
 117 at the start and end. The sentences were chosen randomly from the respective talker’s audio
 118 book. To avoid an obvious talker change detection due to gender attributes, we chose all

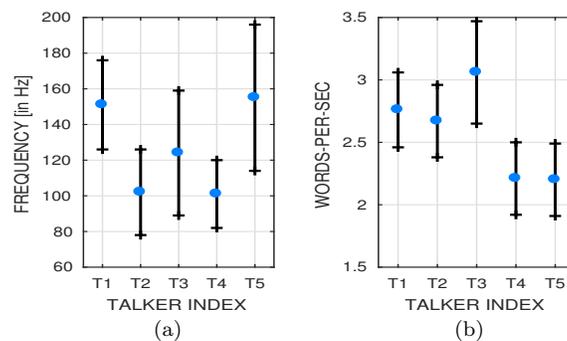


FIG. 2. A comparison of talker attributes with respect to (a) fundamental frequency variation, and (b) word speaking rate. The vertical bars indicate one standard deviation spread around the mean value.

119

120

121 male talkers. Based on an informal pilot experiment aimed at finding a set of perceptually
 122 separable voices, we chose five talkers from the corpus (IDs 374, 2843, 5456, 7447, 7505) for
 123 the listening test stimulus design (here referred to as T1, T2, etc.). The average fundamental
 124 frequency (estimated using STRAIGHT⁴⁴) and speaking rate expressed as words spoken per

125 second (estimated from the audio book transcripts) of the five talkers are depicted in Fig. 2,
126 with significant overlap across talkers to make for challenging TCD.

127 To make a stimulus, talker T_x was chosen from the list of N talkers, and a sentence
128 utterance was retrieved from the corresponding talker’s audio book. A short utterance from
129 another talker T_y was chosen, and this was concatenated to the utterance from T_x . As
130 the utterances were natural speech, there were natural pauses. Owing to this, the silent
131 interval between T_x ’s end and T_y ’s start after concatenation was random and ranged from
132 200 – 1000 ms. In any stimulus, speech corresponding to T_x was between 5 – 10 s and that
133 corresponding to T_y was 4 s. A sample stimulus is shown in Fig. 1.

134 For each pair of T_x - T_y talkers there were $M = 8$ unique stimuli. The stimulus set can
135 be represented as sampled from a grid, as shown in Figure 1. This resulted in a total of
136 $M \times N^2 = 200$ distinct speech stimuli, each 9 – 14 s in duration.

137 C. Protocol and apparatus

138 A graphical user interface (GUI) for stimulus presentation was made using Gorilla¹, a soft-
139 ware platform for designing behavioral science tests. A test session comprised three tasks
140 carried out in sequence, namely, Task-1, 2 and 3. Task-1 measured RT for noise-to-tone
141 detection and Task-2 measured RT for tone frequency change detection (see Supplementary
142 material²). The acoustic attributes associated with a change in the stimuli in these first two
143 tasks were easily recognizable. As a result, Task-1 and Task-2 served as benchmarks against
144 which to compare human RT on Task-3, associated with TCD.

145 In each task, listeners were instructed to press a button (`space bar`) immediately upon
 146 detection of a change in the stimulus. The audio stopped after the button press and visual
 147 feedback indicating \checkmark (or \times) for correct (or incorrect) detection appeared immediately.
 148 Participants were seated in sound-attenuated booths wearing Sennheiser headphones³(with

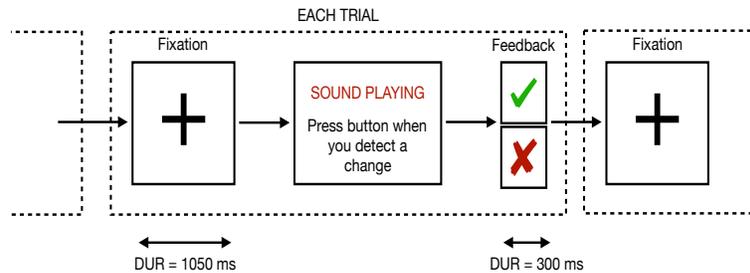


FIG. 3. An illustration of a listening test trial.

149
 150

151 flat spectral response from 60 – 8000 Hz), with diotic stimulus presentation. In order to
 152 prevent fatigue, participants were allowed to take breaks after blocks of 25 trials. Each
 153 participant listened to a few (8 – 10) different Task-3 TCD stimuli (not used in the test) to
 154 become conversant with the paradigm. On average, the complete test session was 45 mins
 155 (with Task-1, Task-2, and Task-3 taking 5, 10 and 30 mins, respectively).

156 D. Performance measures

157 For each trial in which there was a button press, the RT for change detection was obtained
 158 as the difference between the response instant (denoted by t_r) and the ground-truth acoustic
 159 change instant (denoted by t_c), that is, $RT = t_r - t_c$. An illustration is provided in Figure 1.
 160 The lower limit for RT for change perception in sound attributes is on the order of $RT <$
 161 250 ms¹⁹. Hence, RTs in the range 0 – 250 ms are likely to be associated with speech heard

162 prior to the change instant t_c . The upper bound on RT (2000 ms) was chosen based on prior
 163 research²².

164 We analyzed the hits, misses, and false alarms (FA) in the responses. The 200 trials in
 165 Task-3 per subject were categorized into two pools for analyses:

166 *Pool A*, involving trials with T_x a different talker from T_y (two-talker trials) and either
 167 $RT > 225$ ms or no button press.

168 *Pool B*, involving trials with $T_x = T_y$ and trials with $T_x \neq T_y$ but $RT < 225$ ms. These all
 169 are single-talker trials (i.e the trials in which the subject's response was based on attention
 170 to only one talker).

171 From these pools of data, we defined the following detection measures:

172 • **Hit rate:** A hit corresponds to a trial in *Pool A* with $225 \text{ ms} < RT < 2000$ ms. Hit
 173 rate is the ratio of number of hits to the number of trials in *Pool A*.

174 • **Miss rate:** A miss corresponds to a trial in *Pool A* with $RT > 2000$ ms. Miss rate
 175 is the ratio of number of misses to the number of trials in *Pool A*. Note that the miss
 176 rate is $100 - \text{hit rate}$.

177 • **False alarm rate:** A false alarm (FA) corresponds to a trial in *Pool B* featuring a
 178 button press. False alarm rate is the ratio of number of FAs to the sum of trials in
 179 *Pool B* and *Pool A* (this equals 200).

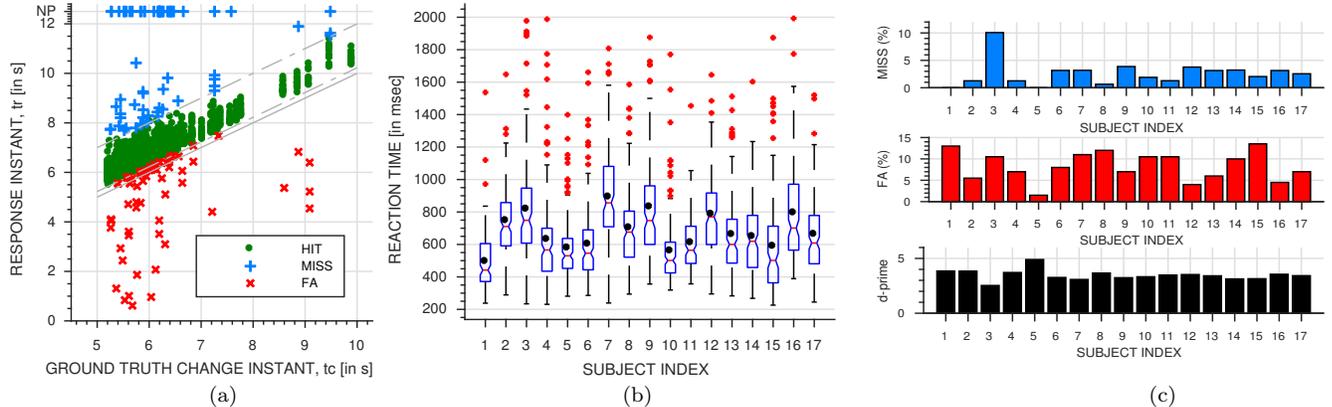


FIG. 4. (color online) (a) Illustration of human reaction time (RT) versus the ground-truth talker change instant (t_r vs t_c) across a total of 2720 trials (with $T_x \neq T_y$) over 17 subjects. The three inclined gray lines from bottom to top correspond to $t_r = t_c$, $t_c + 225$, $t_c + 2000$, respectively. NP stands for no button press. (b) Subject-wise summary using a boxplot of RTs in trials with hits. The black dots correspond to means. (c) Subject-wise miss and false alarm rates, and d-prime obtained from 200 trials for each subject.

180 III. RESULTS: HUMAN TCD EXPERIMENT

181 A. RT distribution

182 Figure 4(a) depicts the distribution of TCD reaction time t_r as a function of ground-truth
 183 talker change instant t_c for all trials which have a talker change. (taken from *Pool A* and
 184 *Pool B*). As seen, the majority (approx. 95%) of responses fall in the hit zone, that is,
 185 $t_c + 225 < t_r < t_c + 2000$ ms. Analyzing the hit trials from *Pool A*, the subject-wise
 186 RT summary is shown in Figure 4(b). Across subjects, the response time to detect a talker
 187 change tended to require mostly under a second of speech from the true change instant, with
 188 subject-dependent distributions of average RT and variability across quantiles. Analyzing

189 the detection parameters, the subject-wise hit, miss and FA rates are shown in Figure 4(c).
 190 The hit, miss, and false alarm rates averaged across all subjects were 97.38%, 2.62%, and
 191 8.32%, respectively. Listeners performed the TCD task very accurately; the average d-prime
 192 across subjects was 3.48 (d-prime is defined as $\mathcal{Z}(\text{hit rate}) - \mathcal{Z}(\text{FA rate})$, where function
 193 $\mathcal{Z}(p)$, $p \in [0, 1]$, is the inverse of the cumulative distribution function of the Gaussian
 194 distribution)

195 The distribution of RT corresponding to hits from all subjects is shown in Figure 5. The
 196 non-Gaussian nature of the data is evident from the histogram and the normal probability
 197 plot. To improve the Gaussianity, we applied a log transformation ($\log_{10} RT$) on the RT
 198 data; the resulting distribution is shown in the same figure. We used this transformed data
 199 in the regression analysis, which is presented next.

200 B. Dependence of RT on speech duration

201 We examined the extent to which the duration of speech experienced prior to a talker change
 202 impacted TCD. We probed this using linear regression analysis on log-RT (averaged across
 203 subjects) versus speech duration before change instant t_c . As the stimulus set is composed of
 204 naturally spoken sentences from story books, the grid along speech duration is non-uniformly
 205 sampled in the dataset. Hence, we performed the regression analysis only for stimuli with
 206 $t_c < 7000$ ms, as beyond this t_c value, the grid was sparsely sampled. The result is shown
 207 in Figure 6. The high variability in log-RT may be attributed to the complicated feature
 208 space associated with speech signals. This is unlike what is observed in tone frequency (or
 209 noise-to-tone) change detection for which the variability in RT correlates with the simple

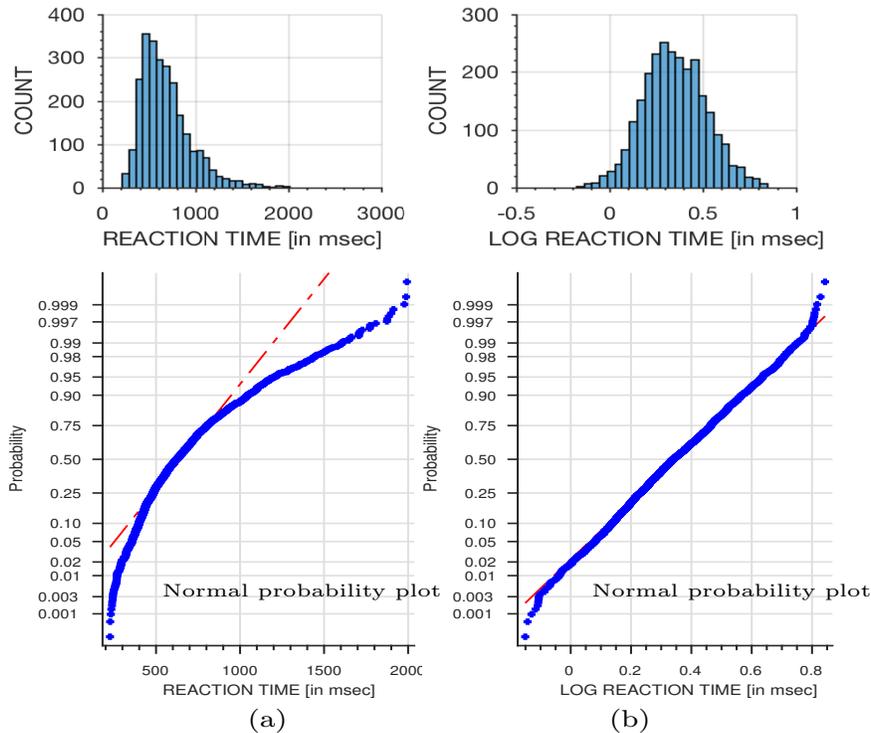


FIG. 5. (color online) Illustration of the distribution (obtained as a histogram) of the RT data for trials on which there was a hit for (a) raw RT data and (b) log-transformed RT data to improve the fit to a normal distribution.

210 spectral difference (see Supplementary material and also²²). Despite considerable variability,
 211 we observe a trend such that TCD is slower with shorter samples of speech from the initial
 212 talker (a decrease in the t_c value). Next, we attempt to model RT as a function of different
 213 acoustic features extracted from the stimulus before and after change instant to understand
 215 the acoustic dimensions listeners may track in TCD.

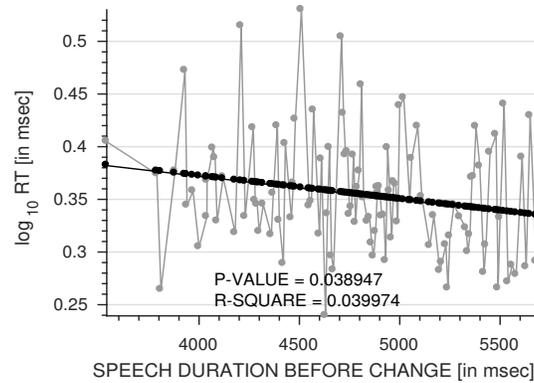


FIG. 6. Dependence of average RT on speech duration before the change instant. The black line is the linear regression fit.

216 C. Modeling RT with acoustic features

217 Past studies^{48,49} have used RT modeling to reveal visual input features associated with object
 218 recognition. Motivated by this, here we model RT to detect a change in talker across hit
 219 trials as a function of the difference in the acoustic feature space between speech segments
 220 sampled before and after the change instant. We consider a collection of acoustic features
 221 that may be important in tracking voice identity, and examine their ability to estimate
 222 human participants' TCD RT.

223 The approach is illustrated in Figure 7(a,b). Let D_b and D_a denote segments of the
 224 stimulus before and after change instant t_c , respectively. We hypothesize that a listener es-
 225 timates acoustic features from these segments, summarizes them and compares the summary
 226 using a distance measure. The resulting distance serves as strength of evidence for change
 227 detection, and by Piéron's law⁵⁰ this should impact the RT. The greater the distance, the
 228 faster listeners may detect a talker change and thus the smaller the RT value. We assume

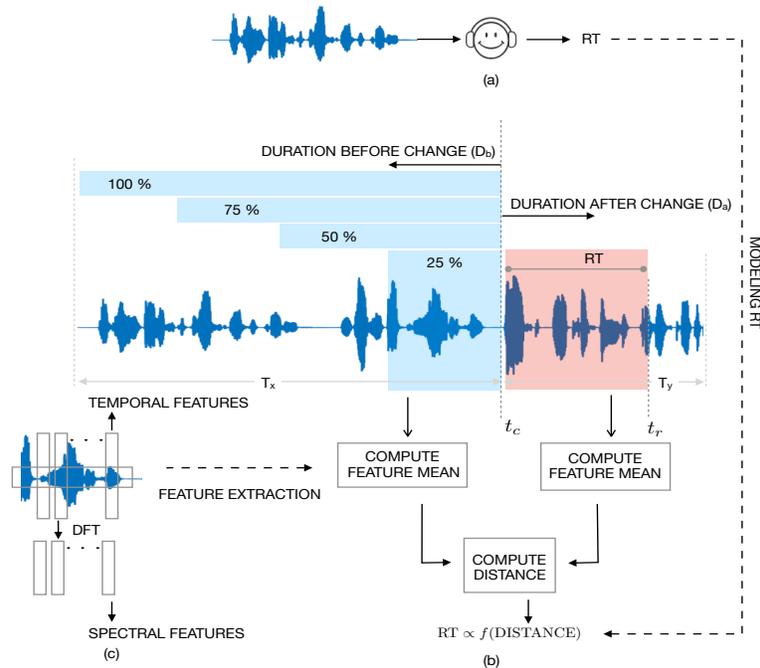


FIG. 7. (color online) Proposed approach to model RT using acoustic features before and after change instant.

229 D_a corresponds to the segment from t_c to t_r . For D_b , we model different segment durations
 230 before the change instant (that is, 25%, 50%, 75%, and 100% of t_c duration segment, before
 231 change instant t_c).

232 1. Feature computation

233 A wide range of acoustic features can be computed from D_b and D_a segments. Here, we
 234 consider the 9 sets of features described in Table I. These are chosen to sample from a
 235 range of possible temporal and spectral acoustic features (illustrated in Figure 7(c)). From
 236 the perspective of speaker attributes, the feature set can be further grouped into those
 237 capturing pitch (F0), vocal tract formants (LSFs, MEL, MFCCs), rate of temporal variation
 238 of vocal tract features (temporal derivatives⁵¹ of MFCCs, namely, MFCC-D and MFCC-

TABLE I. Acoustic features used in the regression model analysis.

FEATURE SET	FEATURES	TYPE	DIMENSION	TIME SCALE
F0	Fundamental Frequency	Spectral	1×1	25 ms
LSF	Line spectral frequencies	Spectral	10×1	25 ms
MEL	Mel-spectrogram	Spectral	40×1	25 ms
MFCC	Mel-frequency cepstral coefficients	Spectral	12×1	25 ms
MFCC-D	First-order temporal derivative of MFCCs	Spectral	12×1	25 ms
MFCC-DD	Second-order temporal derivative of MFCCs	Spectral	12×1	25 ms
ENGY-D	Derivative of short-time energy	Temporal	1×1	25 ms
PLOUD	Loudness strength, sharpness, and spread	Spectral	3×1	25 ms
SPECT	Spectral flatness, Spectral flux, Spectral roll-off, Spectral shape, Spectral slope	Spectral	8×1	25 ms

239 DD), perceived loudness (PLOUD⁵²), spectral timbre (SPECT), and the rate of temporal
 240 variation in short-time energy (ENGY-D). These features are computed every 10 ms with
 241 *Hanning* windowed short-time segments of 25 ms. All features were extracted using the
 242 *Yaafe*⁵³ Python package, an efficient open-source code library for speech and audio analysis.

243 2. Regression on feature distance

244 For each feature set, we summarized the segments D_b and D_a using the mean of the features
 245 in each segment. The PLOUD, and SPECT feature set were characterized by a combination
 246 of different features. Hence, we mean- and variance-normalized these feature sets over the
 247 whole duration prior to segment-wise mean computation. Following this, we computed the
 248 Euclidean distance between the obtained means. Owing to significant variability in RT
 249 across subjects (see Figure 4(b)), we modeled each subject’s RT separately.

250 We defined the number of trials with a hit for the p^{th} subject to be denoted by N_p .
 251 Corresponding to these trials, we have N_p RTs, and for each RT we compute a distance

252 between the mean features extracted from segments D_b and D_a . There are 9 such distances
 253 based on the choice of features, and we denote these by d_k , $k = 1, \dots, 9$, that is, one distance
 254 for each feature set. To evaluate the impact of the feature distances on RT to detect a talker
 255 change, we perform a linear regression on feature distances to estimate the RT using a
 256 regression model for the k^{th} feature set

$$\underbrace{\begin{bmatrix} \log RT_1 \\ \log RT_2 \\ \vdots \\ \log RT_{N_p} \end{bmatrix}}_{\mathbf{r}} = \underbrace{\begin{bmatrix} 1 & d_{k,1} \\ 1 & d_{k,2} \\ \vdots & \vdots \\ 1 & d_{k,N_p} \end{bmatrix}}_{\mathbf{D}} \underbrace{\begin{bmatrix} w_0 \\ w_k \end{bmatrix}}_{\mathbf{w}} \quad (1)$$

257 where w_0 and w_k are the model parameters representing the mean RT and slope of the
 258 regression line, respectively, and RT_i and $d_{k,i}$ denote the RT and feature distance in the i^{th}
 259 trial, respectively. We solve for the model in (5) using minimum mean square error (MMSE).
 260 That is,

$$\hat{\mathbf{w}} = \min_{\mathbf{w}} \|\mathbf{r} - \mathbf{D}\mathbf{w}\|_2 = (\mathbf{D}^T\mathbf{D})^{-1}\mathbf{D}^T\mathbf{r} \quad (2)$$

$$\hat{\mathbf{r}} = \mathbf{D}\hat{\mathbf{w}}. \quad (3)$$

261 The MMSE optimized values of w_k for all subjects and for the 9 features sets are shown in
 262 Figure 8(a). These are obtained for D_b set to 100% of t_c . For a majority of subjects, w_k is
 263 negative for all the feature sets. This signifies that, on average, RT decreases with increased
 264 feature distance. Some feature sets have a greater negative slope than others. For example,
 265 the slope is maximally negative for MFCC-D and MFCC-DD feature sets. To quantify the

266 modeling performance, we used the r-square measure⁵⁴, computed as

$$\text{r-square} = 1 - \frac{\|\mathbf{r} - \hat{\mathbf{r}}\|_2^2}{\|\mathbf{r} - \bar{\mathbf{r}}\|_2^2} \quad (4)$$

267 where $\bar{\mathbf{r}}$ is the mean of elements in \mathbf{r} . The r-square is also referred to as the “explained
268 variance” by the model; a value close to 100% indicates good modeling performance; that
269 is, the proposed model is able to better explain the observed variance in the data.

271 Figure 8(b) shows the obtained r-square (shown in % as explained variance) for different
272 feature sets. For clarity of presentation, the depicted percentage is the average percentage
273 across all subjects. At the level of individual features, the MFCC-D outperforms all other
274 feature sets. Moreover, a majority of feature sets fall below 10%, thereby failing to explain
275 a significant portion of the variance in RT.

276 To examine combinations of features sets, we performed a multiple linear regression by
277 combining the feature distances from all feature sets as follows,

$$\underbrace{\begin{bmatrix} \log RT_1 \\ \log RT_2 \\ \vdots \\ \log RT_{N_p} \end{bmatrix}}_{\mathbf{r}} = \underbrace{\begin{bmatrix} 1 & d_{1,1} & \dots & d_{9,1} \\ 1 & d_{1,2} & \dots & d_{9,2} \\ \vdots & \vdots & \dots & \vdots \\ 1 & d_{1,N_p} & \dots & d_{9,N_p} \end{bmatrix}}_{\mathbf{D}} \underbrace{\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_9 \end{bmatrix}}_{\mathbf{w}} \quad (5)$$

278 This gave the best r-square (see Figure 8(b)). A p-value (from hypothesis testing) illustration
279 depicting the relevance of each feature set in the multiple linear regression is shown in
280 Figure 8(c). The MFCC-D and MFCC-DD are most relevant, $p < 0.05$, across all subjects.
281 A scatter plot of RT versus estimated RT , pooling all subjects’ trials, is shown in Figure 8(d).

Talker change detection

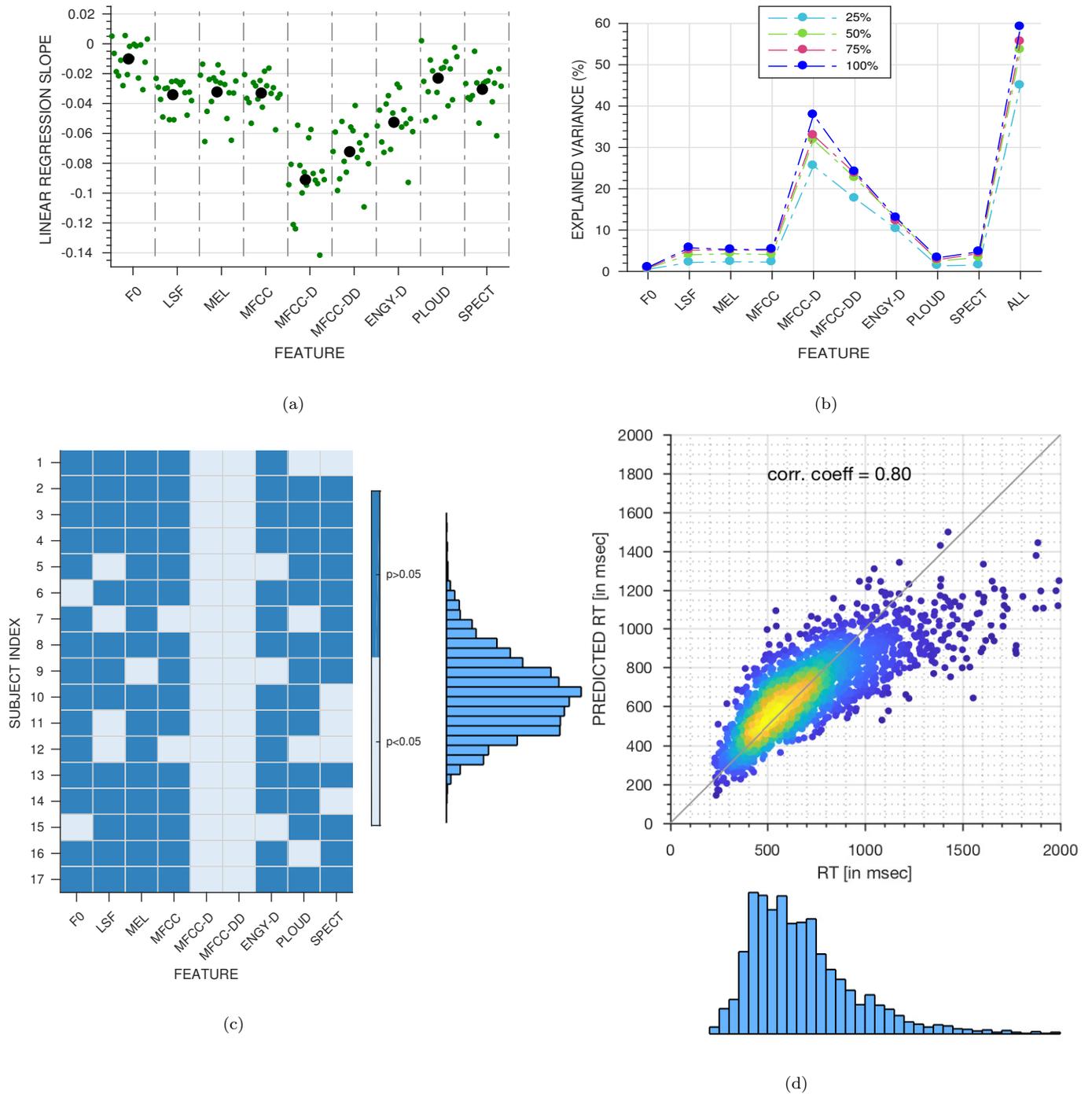


FIG. 8. (color online) (a) Linear regression slope indicating the performance of individual feature sets in accounting for the variance in human TCD reaction times. (b) The percentage of explained variance (r-square) for the linear regression model, averaged across subjects. (c) Illustration of significance (p-value) of different feature sets in multiple linear regression. (d) Scatter plot (polling all subjects' trials) of true RT versus estimated RT obtained from a multiple linear regression model computed across all feature sets.

282 The Pearson correlation coefficient amounted to 0.8, indicating a good overall estimation
 283 accuracy. Also, it can be seen that a major portion of true RT fall in the range 400–1000 ms,
 284 with few trials with $RT > 1000$ ms. In this range, the estimation is also concentrated along
 285 the $y = x$ line.

286 We also tested the modeling performance with decreasing duration of segment D_b , that
 287 is duration set to 75, 50, or 25% of t_c duration before the change instant (a shown in
 288 Figure 7(b)). The result shown in Figure 8(b) depicts that using 100% of the segment
 289 duration better models human performance, with systematic decreases as the sample of the
 290 talker’s speech decreases in duration.

291 D. Talker pair-wise analysis

292 To analyze the variability in TCD performance across the talker pairs, we examined talker-
 293 wise performance in TCD RT (see Figure 9). Most of the talker pairs have the average RT
 294 (computed across subjects) in the same range, except $T_1 - T_5$ and $T_5 - T_1$. Also, the miss
 295 rate was found to be higher for these pairs of talkers. This suggests that these pairs may
 296 be overlapping a lot in the perceived talker space. Comparing the FA rate, averaged across
 297 subjects, talker T_3 had the highest FA rate.

298 IV. MACHINE SYSTEM FOR TCD

299 We evaluated the performance of three machine TCD systems on the same stimulus materials
 300 used in the human TCD experiment. The first system was an adaptation of a state-of-the-
 301 art diarization system, designed to segment audio into distinct talker segments based on

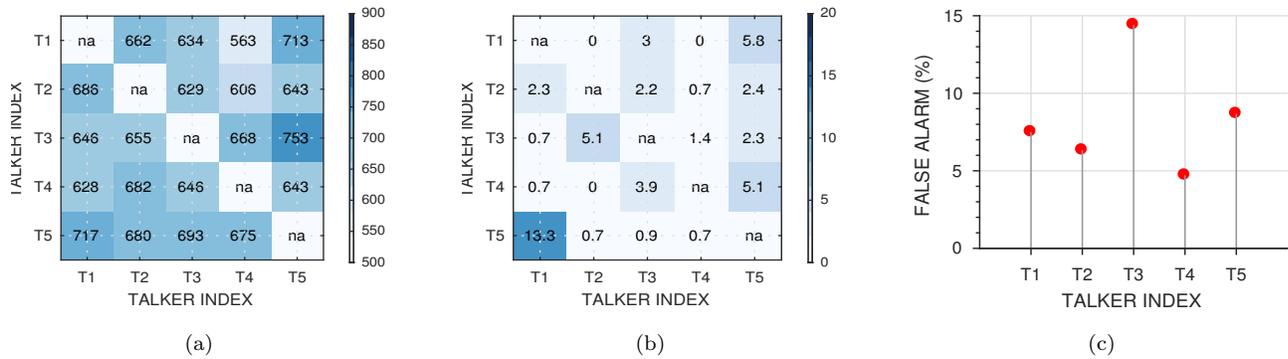


FIG. 9. (color online) Dependence of: (a) average RT on talker pairs (T_x-T_y), (b) average miss rate on talker pairs (T_x-T_y), and (c) average false alarm rate across talkers.

302 i-vector and probabilistic linear discriminant analysis (PLDA)⁵⁵. Subsequently, the talker
 303 change instants can be obtained as segment boundaries. Traditionally, these systems operate
 304 across a whole audio file. We refer to this mode of operation as off-line TCD.

305 In contrast to off-line machine systems, listeners in the human TCD task did not listen
 306 to the whole audio file. Instead, they performed an on-line change detection, pressing the
 307 button as soon as a talker change was perceived. To better model this aspect of human
 308 listening, we implemented an on-line variant of the off-line machine system. Here, the
 309 system sequentially operated on segments of increasing duration starting with an initial 1 s
 310 segment for a given audio file. The sequential operation was stopped as soon as the second
 311 talker was detected, with this instant corresponding to the response instant for a talker
 312 change in machine recognition. An illustration of the approach is shown in Figure 10. As
 313 we hypothesized for human listeners, the on-line system uses the acoustic features extracted
 314 from onset until the current segment duration in making a decision for talker segmentation.

315 The second machine system is the commercially available state-of-the-art IBM Watson
 316 Speech-to-Text (STT) system that incorporates a talker change detector module⁴. The
 317 third system is purely based on textual features, with no access to the acoustic detail of the
 318 speech. This text-based TCD system takes input from the transcript of the audio file and
 319 analyzes the semantic similarity between contiguous words for talker change detection. It
 320 thereby provides a control system with which to evaluate whether our approach to controlling
 321 semantic similarity in the novel TCD task (introducing a change in sentence context on both
 322 trials with and without a talker change) was successful.

323 For comparison of machine performance with human TCD behavior, the hit-rate, miss-
 324 rate and the false alarm rate were computed using the same definitions as those for human
 325 TCD experiments. The following subsections describe the systems and results.

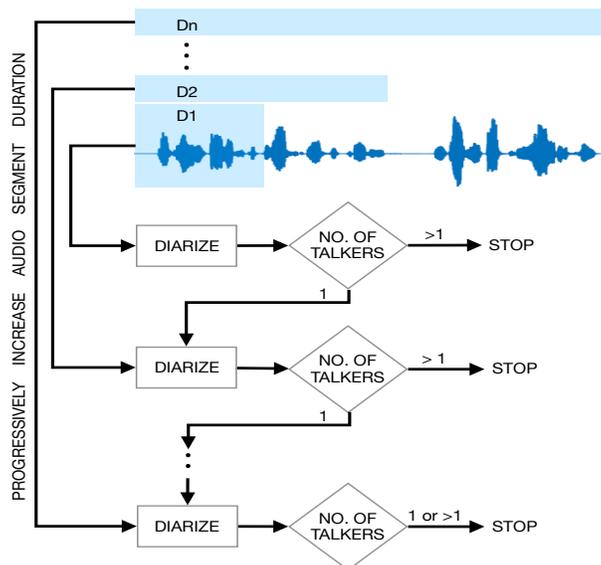


FIG. 10. Proposed sequential diarization system as an on-line system for TCD.

A. On-line and Off-line diarization-based TCD systems

We first segmented an audio file into small (750 ms) and temporally-overlapping segments (temporal shift of 500 ms). Each short segment was transformed into an i-vector representation⁵⁷. All pair-wise distances between the short-segment i-vectors were computed using a PLDA based scoring approach. In the literature, PLDA based scoring has been shown to impart robustness against small duration recordings (5 – 10 s) and variations in talking style⁵⁸, and this suit the present stimulus set. Using the PLDA scores, agglomerative hierarchical clustering was performed to identify short-segments belonging to the same talkers and subsequently to merge the short-segments of the same talker. The obtained output was a segmentation of the input audio file into distinct single talker segments.

We developed the diarization system in the following two modes:

Off-line Diarization Here, the diarization was performed on the complete audio file in one pass.

On-line Diarization Here, instead of doing diarization across the complete audio file in one pass, we began with an input of 1 s and then sequentially increased it by 1 s, until two talker segments were detected or the end of file was reached (illustrated in Figure 10).

The complete system setup was developed using the Kaldi toolkit⁵⁹. This involved training the i-vector extractor based on a universal background model composed of a 512-component Gaussian mixture model with a diagonal covariance matrix and trained on the LibriSpeech corpus⁴³. The system used 12-dimensional MFCC features, obtained from successive 25 ms (with temporal shifts of 10 ms) short-time segments derived from the audio files. The

347 MFCC features were mean- and variance-normalized using a 3 s running window. The i-
 348 vector representations were 128-dimensional. The audio files corresponding to talkers used
 349 in the listening test were removed from the training dataset.

350 Since the extraction of i-vectors involves an overlapping shift, the system can flag a change
 351 instant before the actual ground-truth change instant. In order to account for this look-
 352 ahead, for these systems a tolerance window $\delta = 500$ ms was given and the RT corresponded
 353 to a hit if $t_c - \delta < RT < t_c + 2000$ ms. The operating point for the system can be tuned
 354 by varying the segment clustering threshold. Thus, each operating point had its own miss
 355 and false-alarm-rates. Hence, we generated a detection error trade-off curve for each system.
 356 These are shown in Figure 12.

357 B. IBM Watson speech-to-text (STT) system

358 Along with an automatic speech recognizer (ASR), the commercially-available IBM Watson
 359 STT system also includes a follow-up TCD system. Built as a one-of-its-kind approach⁶⁰,
 360 the TCD system makes use of word boundaries from the ASR output in addition to acoustic
 361 feature information. For each acoustic segment corresponding to a pair of contiguous words,
 362 two separate Gaussian models are fit using MFCC features from left and right of the detected
 363 word boundary, respectively. A talker change is flagged based on the BIC algorithm or
 364 T^2 -criterion. Use of word annotations from ASR reduces false alarms in change detection
 365 within a word and in non-speech regions. We used the off-the-shelf implementation of this
 366 system available as a Web API. This system is pre-tuned and provides only one operating
 367 point. Like the diarization system, this system also was found to often give a change instant

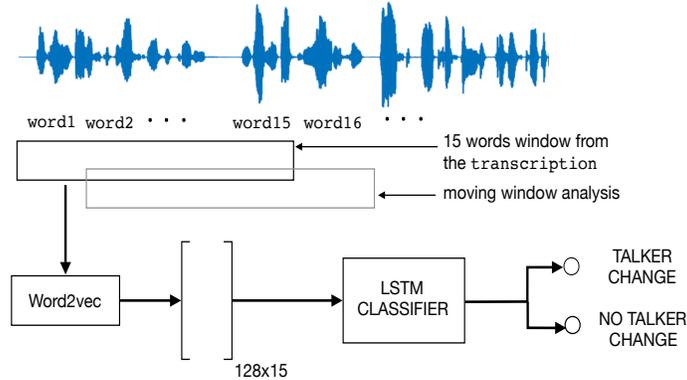


FIG. 11. Illustration of text-based context change classification.

368 before the ground-truth change instant. This may be due to the larger temporal context
 369 used in the ASR module. On post analysis for hit rate computation, the detected change
 370 instant was found to lie in $t_c - \delta < RT < t_c + 2000$ ms, with $\delta = 200$ ms. We considered all
 371 these responses as correct detections in hit rate computation. The resulting miss and false
 372 alarm rates are shown in Figure 12.

374 C. TCD based on textual features

375 Although we hypothesize that acoustic feature distributions are likely to play a major role
 376 in the listeners' model of talker identity, it is also likely that the listeners attended to the
 377 semantics of the words making the stimuli.

378 Designing an experimental paradigm that mitigates semantic contributions in order to
 379 evaluate the acoustic features that contribute to TCD is challenging. The present paradigm
 380 takes the approach of introducing a *semantic* change on every trial. Hence, listeners cannot
 381 rely on a semantic change as a reliable cue for TCD as some trials feature no talker change.
 382 The challenge to disentangle the contribution of text/content versus acoustic features led us

383 to examine a machine equivalent of TCD reliant upon only the information conveyed by the
384 text of our stimulus material^{61,62}. The proposed text-based TCD system is shown in Fig. 11.
385 Using the transcripts from the LibriSpeech corpus, we trained a `Word2Vec` model with the
386 Gensim Python package⁶³. The model represents every word with 128-dimensional vector.
387 This allows representing a sentence as a sequence of vectors obtained from the constituent
388 words. We built a talker change detector by analyzing the semantic similarity among sets
389 of words via analysis of the vector representations.

390 Specifically, the system is a classifier that takes N consecutive words as input and outputs
391 a class label C_0 if all words are by a single talker, and a class label C_1 otherwise. The value of
392 N in this experiment was chosen after analyzing the average number of words in a sentence
393 spoken by the first talker as well as the average number of words spoken by the second talker
394 within 2 s. Using features from 15 consecutive words from the corpus of talkers reading audio
395 books, we generated training examples for the two classes. The training set comprised of
396 two different talkers (label C_1) was created by taking the last 10 words of a sentence from
397 the transcript of an audio book read by one talker and the first 5 words from the transcript
398 of an audio book read by another talker. For the set of examples with no talker change
399 (class C_0), the dataset was created by taking 10 words from the end of one sentence and
400 the first 5 words from the beginning of the following sentence, both drawn from transcripts
401 of the audio book read by one talker. This dataset was created excluding the talkers in the
402 listening set. Using the training examples, we train a Long Short-Term Memory Network
403 (LSTM) to predict talker change based on context. The model was designed using the Keras
404 toolkit⁶⁴ comprising of one LSTM layer with 512 cells followed by three dense layers with

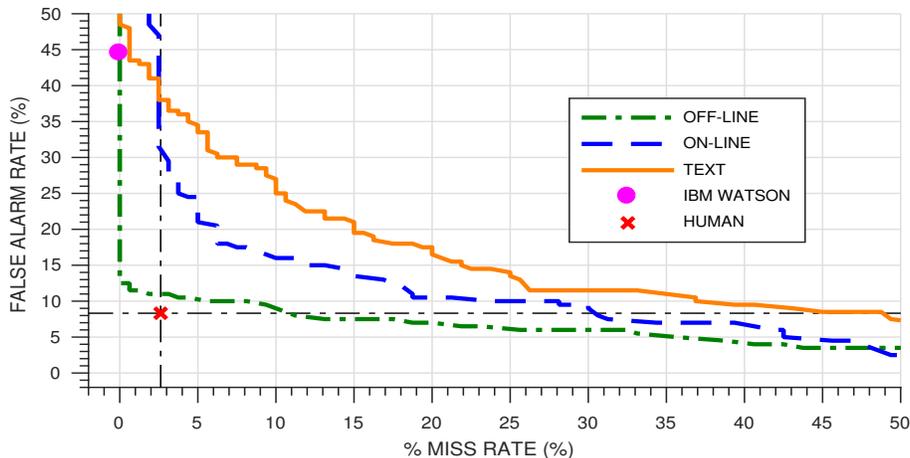


FIG. 12. (color online) Detection Error Trade-off (DET) curve for all the TCD systems evaluated on the stimuli set used in the listening experiments. The IBM Watson system is a single operating point. The human score is also marked in this figure for reference.

405 1024, 512 and 256 neurons, respectively. The test corpus was made from the transcripts of
 406 the sentences used in the listening test. Each sentence was input as a series of $N = 15$ words,
 407 to reflect the training set input, shifting a word by one to capture all words. A stimulus is
 408 marked as a change when it consists of at least one word from C_1 . For computing detection-
 409 error-trade-off curves, a threshold was used on the posteriors from the LSTM model. The
 410 resulting performance is shown in Figure 12.

411 V. DISCUSSION

412 Summarizing the findings, the results from the experiments make three primary contribu-
 413 tions towards understanding human and machine TCD performance.

414 A. Human performance for TCD

415 Human listeners performed the TCD task with very high accuracy, averaging only a 2.62%
416 miss rate and a 8.32% false alarm rate. The presence of false alarms signifies that human
417 listeners sometimes reported a talker change when there was none. This can be partly at-
418 tributed to varying stress, intonation, and speaking rate exhibited by each talker during the
419 course of reading the story book. At the same time, the high accuracy contrasts with some
420 prior studies reporting high rates of “change deafness” to a change in talker¹¹. However, it
421 is important to note that the present task required listeners to direct the attention to detect
422 a change in talker rather than to the comprehension alone. As suggested in prior research¹¹,
423 conversation expectations and task are likely to influence human listeners’ direction of at-
424 tention to fine-grain details of voice. At least in the context of an overt task requiring TCD,
425 the present results demonstrate human listeners’ ability to track acoustic features across
426 talkers and to differentiate the dimensions most relevant to detecting a change among male
427 voices. Here, in the interest of investigating the issue in a controlled task, we examined
428 read speech from audio books. The same approach might be applied to other stimulus sets
429 capturing even more natural conversational speech, although in these cases covariation with
430 context, semantic continuity, and other factors would be likely to complicate attempts to
431 understand listeners’ ability to track distributional information across acoustic dimensions
432 that are related to talker identities.

433 Speculating from average *RT* to detect a change, the perceptual load for TCD appears to
434 be greater than for simpler acoustic change detection scenarios such as noise-to-tone change

435 or tone frequency change in the same listeners (experiments reported in supplementary ma-
436 terial). Pooling trials across listeners, the average RT for change detection was 680 ms
437 (std. dev.= 274 ms), which was close to twice the average RT for a noise-to-tone change
438 detection, and falls within the average RT associated with a tone frequency change detec-
439 tion task (done with different magnitudes of change). We speculate that these differences
440 may be attributed to recruitment of different auditory sub-processes for change detection
441 in the context of speech and non-speech stimuli, specifically the need to accumulate distri-
442 butional acoustic information across the more complex, multi-dimensional acoustic features
443 that convey talker identity.

444 **B. Estimating RT using a simple regression model on feature distances**

445 We used a linear model to relate human listeners' $\log RT$ to detect a talker change and acous-
446 tic feature distances across a set of acoustic features. A simple Euclidean distance measure
447 between the mean of the acoustic feature measurement corresponding to speech segment
448 before and after change instant was used. Interestingly, we found the Piéron's law⁵⁰, stating
449 decrease in RT with increase in "strength of evidence" (such as loudness or frequency differ-
450 ence for tone stimuli), to hold for distance computed in the feature space (negative slopes in
451 Figure 8(a)) for talker change detection, as well. Quantifying the model performance in terms
452 of the percent of the human performance variance explained, the best fit was obtained in the
453 MFCC-D feature space, followed by MFCC-DD and ENGY-D features. The 12-dimensional
454 MFCC representation derived from MEL representation is a spectrally smoothed represen-
455 tation of the short-time spectrum, preserving the spectral peaks corresponding to formants,

456 with minimal pitch information. The MFCC-D representation, derived from the temporal
457 derivative of MFCC, captures the rate of variation in the spectrally smoothed representa-
458 tion. The improved model performance with the MFCC-D feature suggests that listeners
459 are likely using the spectro-temporal variability in formant frequency space while attending
460 to detect talker changes. We found a poor model fit with fundamental frequency (F0), often
461 implicated in talker differences. This may be because there was considerable overlap in the
462 fundamental frequency spread across our all-male set of talkers (shown in Figure 2). Like-
463 wise, it is notable that application of a multiple regression model using the 9 feature sets
464 as independent variables improved the model performance significantly ($r\text{-square} \approx 0.6$). A
465 significant contribution in this model came from MFCC-D and MFCC-DD ($p < 0.05$, see
466 Figure 8(c)). This suggests that TCD involved tracking acoustic information across a quite
467 complicated feature space. Visualizing the quality of estimated and true *RTs* (pooling from
468 all subjects) we found a correlation of 0.8 (Figure 8(d)). A majority of true RTs fall in the
469 range 400 – 900 msec, and a significant portion of this was mapped to the same range by the
470 proposed model. Focusing on the duration of speech segment before change instant used in
471 the model estimation showed a systematic improvement in performance with duration for all
472 the features (Figure 8(c)). This suggests that TCD response is less predictable using local
473 information around the change instant and likely the listener continuously builds a talker
474 model while attending to speech.

475 C. Machine performance on the TCD Task

476 The off-line diarization system provided better performance than the on-line system machine
477 TCD system. This can be attributed to improved clustering of the i-vectors into two clusters
478 because of availability of more data after change instant compared to the on-line case. The
479 text-based TCD system (IBM Watson) exhibited relatively poor performance, indicating
480 that, for the stimulus set, the text features did not suffice for TCD. This provides assurance
481 that our novel paradigm for assessing human TCD (in which a change in sentence context
482 was introduced on each trial independent of a talker change) was sufficient to reduce any
483 bias arising from the semantic content of the utterance before and after the change instant.
484 The IBM Watson has a false alarm rate too high for successful use of the system in natural
485 speech conversations, and underscores the importance of modeling the acoustic distributional
486 characteristics of talkers voices in supporting successful diarization. In fact, documentation
487 of the system does caution that the presence of speaking style variations in the utterances
488 may lead to high false alarm rates. Comparing human and machine systems we found
489 a considerable gap in performance, with humans significantly outperforming state-of-the-
490 art machine systems (see Figure 12). In all, the present results highlight the complexity
491 of the acoustic feature space that listeners must navigate in detecting talker change and
492 underscores that machine systems have yet to incorporate the optimal feature set to model
493 human behavior.

494 VI. CONCLUSIONS

495 The key contributions from the paper can be summarized as follows,

496 (i) Developing a novel paradigm for probing the human talker change detection (TCD)
497 across short-duration natural speech utterances.

498 (ii) Characterizing human TCD performance using various parameters - reaction time
499 (RT), hit-rate, miss-rate and false alarm-rate.

500 (iii) Building a simple linear regression based model that estimates the human RT in TCD
501 using the distance between mean acoustic features from speech segments. This model
502 revealed the significance of MFCC-D and MFCC-DD acoustic features in TCD.

503 (iv) Comparing and benchmarking the machine TCD system performance implemented
504 using principles of speaker diarization and textual features with human TCD perfor-
505 mance.

506 VII. ACKNOWLEDGMENT

507 The authors would like to acknowledge the generous support of Kris Gopalakrishan, Brain-
508 Hub, and Carnegie Mellon Neuroscience Institute that fostered the collaboration between
509 Indian Institute of Science and Carnegie Mellon University to pursue this work, Prachi
510 Singh for the help with implementation of the machine systems, and all the volunteers who
511 participated in this study.

512 **VIII. REFERENCES**513 **REFERENCES**

514 ¹S. D. Goldinger, “Echoes of echoes? An episodic theory of lexical access,” *Psychological*
515 *Review* **105**(2), 251–279 (1998).

516 ²J. D. M. Laver, “Voice quality and indexical information,” *British Journal of Disorders of*
517 *Communication* **3**(1), 43–54 (1968).

518 ³S. C. Levinson, “Turn-taking in human communication - Origins and implications for lan-
519 *guage processing,” Trends in Cognitive Sciences* **20**(1), 6 – 14 (2016).

520 ⁴L. C. Nygaard and D. B. Pisoni, “Talker-specific learning in speech perception,” *Perception*
521 *& Psychophysics* **60**(3), 355–376 (1998).

522 ⁵P. T. Kitterick, P. J. Bailey, and A. Q. Summerfield, “Benefits of knowing who, where, and
523 *when in multi-talker listening,” The Journal of the Acoustical Society of America* **127**(4),
524 2498–2508 (2010).

525 ⁶I. S. Johnsrude, A. Mackey, H. Hakyemez, E. Alexander, H. P. Trang, and R. P. Carlyon,
526 “Swinging at a cocktail party: Voice familiarity aids speech perception in the presence of a
527 *competing voice,” Psychological Science* **24**(10), 1995–2004 (2013).

528 ⁷M. J. Sjerps, H. Mitterer, and J. M. McQueen, “Listening to different speakers: On the
529 *time-course of perceptual compensation for vocal-tract characteristics,” Neuropsychologia*
530 **49**(14), 3831 – 3846 (2011).

531 ⁸Y. Lavner, I. Gath, and J. Rosenhouse, “The effects of acoustic modifications on the iden-
532 tification of familiar voices speaking isolated vowels,” *Speech Communication* **30**(1), 9 – 26
533 (2000).

534 ⁹G. Sell, C. Suied, M. Elhilali, and S. Shamma, “Perceptual susceptibility to acoustic ma-
535 nipulations in speaker discrimination,” *The Journal of the Acoustical Society of America*
536 **137**(2), 911–922 (2015).

537 ¹⁰T. Chi, P. Ru, and S. A. Shamma, “Multiresolution spectrotemporal analysis of complex
538 sounds,” *J. Acoust. Soc. America* **118**(2), 887–906 (2005).

539 ¹¹K. M. Fenn, H. Shintel, A. S. Atkins, J. I. Skipper, V. C. Bond, and H. C. Nusbaum,
540 “When less is heard than meets the ear: Change deafness in a telephone conversation,”
541 *Quarterly Journal of Experimental Psychology* **64**(7), 1442–1456 (2011).

542 ¹²M. S. Vitevitch, “Change deafness: The inability to detect changes between two voices,”
543 *Journal of Experimental Psychology: Human Perception and Performance* **29**(2), 333
544 (2003).

545 ¹³J. G. Neuhoff, S. A. Schott, A. J. Kropf, and E. M. Neuhoff, “Familiarity, expertise, and
546 change detection: Change deafness is worse in your native language,” *Perception* **43**(2-3),
547 219–222 (2014).

548 ¹⁴D. A. Coker and J. Burgoon, “The nature of conversational involvement and nonverbal
549 encoding patterns,” *Human Communication Research* **13**(4), 463–494.

550 ¹⁵J. Kreiman and D. Sidtis, *Foundations of voice studies: An interdisciplinary approach to*
551 *voice production and perception* (John Wiley & Sons, 2011).

- 552 ¹⁶M. Latinus, P. McAleer, P. E. Bestelmeyer, and P. Belin, “Norm-based coding of voice
553 identity in human auditory cortex,” *Current Biology* **23**(12), 1075–1080 (2013).
- 554 ¹⁷L. E. Humes and J. B. Ahlstrom, “Relation between reaction time and loudness,” *Journal*
555 *of Speech, Language, and Hearing Research* **27**(2), 306–310 (1984).
- 556 ¹⁸J. Schlittenlacher, W. Ellermeier, and G. Avci, “Simple reaction time for broadband sounds
557 compared to pure tones,” *Attention, Perception, & Psychophysics* **79**(2), 628–636 (2017).
- 558 ¹⁹D. S. Emmerich, D. A. Fantini, and W. Ellermeier, “An investigation of the facilitation
559 of simple auditory reaction time by predictable background stimuli,” *Perception & Psy-*
560 *chophysics* **45**(1), 66–70 (1989).
- 561 ²⁰C. Suied, P. Susini, and S. McAdams, “Evaluating warning sound urgency with reaction
562 times,” *Journal of experimental psychology: applied* **14**(3), 201 (2008).
- 563 ²¹C. Suied, P. Susini, S. McAdams, and R. D. Patterson, “Why are natural sounds detected
564 faster than pips?,” *The Journal of the Acoustical Society of America* **127**(3), EL105–EL110
565 (2010).
- 566 ²²Y. Boubenec, J. Lawlor, U. Górska, S. Shamma, and B. Englitz, “Detecting changes in
567 dynamic and complex acoustic environments,” *ELife* **6** (2017).
- 568 ²³E. Shriberg, “Spontaneous speech: How people really talk and why engineers should care,”
569 in *Ninth European Conference on Speech Communication and Technology* (2005).
- 570 ²⁴J. Barker, S. Watanabe, E. Vincent, and J. Trmal, “The fifth CHiME speech separation
571 and recognition challenge: Dataset, task and baselines,” *CoRR* **abs/1803.10609** (2018)
572 <http://arxiv.org/abs/1803.10609>.

- 573 ²⁵G. Sell and A. McCree, “Multi-speaker conversations, cross-talk, and diarization for
574 speaker recognition,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal Process.* (2017), pp.
575 5425–5429.
- 576 ²⁶O. Novotný, P. Matějka, O. Plchot, O. Glembek, L. Burget, and J. Černocký, “Analysis
577 of speaker recognition systems in realistic scenarios of the SITW 2016 Challenge,” in *Proc.*
578 *INTERSPEECH, ISCA* (2016), pp. 828–832.
- 579 ²⁷X. Huang and K. F. Lee, “On speaker-independent, speaker-dependent, and speaker-
580 adaptive speech recognition,” *IEEE Trans. on Speech and Audio Processing* **1**(2), 150–157
581 (1993).
- 582 ²⁸A. G. Adam, S. S. Kajarekar, and H. Hermansky, “A new speaker change detection method
583 for two-speaker segmentation,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal Process.*,
584 IEEE (2002), Vol. 4, pp. 3908–3911.
- 585 ²⁹J. Ajmera, I. McCowan, and H. Bourlard, “Robust speaker change detection,” *IEEE signal*
586 *processing letters* **11**(8), 649–651 (2004).
- 587 ³⁰N. Dhananjaya and B. Yegnanarayana, “Speaker change detection in casual conversations
588 using excitation source features,” *Speech Communication* **50**(2), 153 – 161 (2008).
- 589 ³¹V. Gupta, “Speaker change point detection using deep neural nets,” in *Proc. IEEE Intl.*
590 *Conf. Acoust. Speech Signal Process.*, IEEE (2015), pp. 4420–4424.
- 591 ³²R. Wang, M. Gu, L. Li, M. Xu, and T. F. Zheng, “Speaker segmentation using deep speaker
592 vectors for fast speaker change scenarios,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal*
593 *Process.*, IEEE (2017), pp. 5420–5424.

- 594 ³³A. Tritschler and R. A. Gopinath, “Improved speaker segmentation and segments clus-
595 tering using the bayesian information criterion,” in *Sixth European Conference on Speech*
596 *Communication and Technology* (1999).
- 597 ³⁴M. Sarma, S. N. Gadre, B. D. Sarma, and S. R. M. Prasanna, “Speaker change detection
598 using excitation source and vocal tract system information,” in *2015 Twenty First National*
599 *Conference on Communications (NCC)*, IEEE (2015), pp. 1–6.
- 600 ³⁵M. Yang, Y. Yang, and Z. Wu, “A pitch-based rapid speech segmentation for speaker
601 indexing,” in *Seventh IEEE International Symposium on Multimedia (ISM’05)* (2005).
- 602 ³⁶B. Abdolali and H. Sameti, “A novel method for speech segmentation based on speakers’
603 characteristics,” CoRR [abs/1205.1794](#) (2012).
- 604 ³⁷W. N. Chan, T. Lee, N. Zheng, and H. Ouyang, “Use of vocal source features in speaker
605 segmentation,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal Process.*, IEEE (2006).
- 606 ³⁸H. Gish, M. H. Siu, and R. Rohlicek, “Segregation of speakers for speech recognition and
607 speaker identification,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal Process.*, IEEE
608 (1991), Vol. 2, pp. 873–876.
- 609 ³⁹S. S. Cheng, H. M. Wang, and H. C. Fu, “BIC-based speaker segmentation using divide-
610 and-conquer strategies with application to speaker diarization,” *IEEE Transactions on Au-*
611 *dio, Speech, and Language Processing* **18**(1), 141–157 (2010).
- 612 ⁴⁰A. S. Malegaonkar, A. M. Ariyaeinia, and P. Sivakumaran, “Efficient speaker change
613 detection using adapted gaussian mixture models,” *IEEE Transactions on Audio, Speech,*
614 *and Language Processing* **15**(6), 1859–1869 (2007).

615 ⁴¹V. Karthik, D. Satish, and C. Sekhar, “Speaker change detection using support vector
616 machine,” in *Proc. 3rd Int. Conf. Non-Linear Speech Process* (2005), pp. 19–22.

617 ⁴²R. Wang, M. Gu, L. Li, M. Xu, and T. F. Zheng, “Speaker segmentation using deep speaker
618 vectors for fast speaker change scenarios,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal
619 Process.* (2017), pp. 5420–5424.

620 ⁴³V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An ASR corpus based
621 on public domain audio books,” in *Proc. IEEE Intl. Conf. Acoust. Speech Signal Process.*
622 (2015), pp. 5206–5210.

623 ⁴⁴H. Kawahara, I. Masuda-Katsuse, and A. de Cheveigné, “Restructuring speech represen-
624 tations using a pitch-adaptive timefrequency smoothing and an instantaneous-frequency-
625 based F0 extraction: Possible role of a repetitive structure in sounds,” *Speech Communi-
626 cation* **27**(34), 187 – 207 (1999).

627 ⁴⁵<https://gorilla.sc>.

628 ⁴⁶See Supplementary material for supplementary experiments and results on change detec-
629 tion.

630 ⁴⁷Sennheiser HD 215 II Closed Over-Ear Back Headphone with High Passive Noise Attenu-
631 ation.

632 ⁴⁸A. Mirzaei, S.-M. Khaligh-Razavi, M. Ghodrati, S. Zabbah, and R. Ebrahimpour, “Pre-
633 dicting the human reaction time based on natural image statistics in a rapid categorization
634 task,” *Vision Research* **81**, 36 – 44 (2013).

635 ⁴⁹R. T. Pramod and S. P. Arun, “Do computational models differ systematically from human
636 object perception?,” in *The IEEE Conference on Computer Vision and Pattern Recognition*
637 *(CVPR)* (2016).

638 ⁵⁰D. Pins and C. Bonnet, “On the relation between stimulus intensity and processing time:
639 Piéron’s law and choice reaction time,” *Perception & Psychophysics* **58**(3), 390–400 (1996).

640 ⁵¹L. R. Rabiner and B.-H. Juang, *Fundamentals of speech recognition*, Vol. 14 (PTR Prentice
641 Hall Englewood Cliffs, 1993).

642 ⁵²G. Peeters, “A large set of audio features for sound description (similarity and classifica-
643 tion) in the CUIDADO project,” Tech. Rep., IRCAM (2004).

644 ⁵³B. Mathieu, S. Essid, T. Fillon, J. Prado, and G. Richard, “Yaafe, an easy to use and
645 efficient audio feature extraction software.,” in *ISMIR* (2010), pp. 441–446.

646 ⁵⁴A. C. Cameron and F. A. G. Windmeijer, “An R-squared measure of goodness of fit for
647 some common nonlinear regression models,” *Journal of Econometrics* **77**(2), 329 – 342
648 (1997).

649 ⁵⁵G. Sell and D. Garcia-Romero, “Speaker diarization with PLDA i-vector scoring and unsu-
650 pervised calibration,” in *Spoken Language Technology Workshop (SLT), 2014 IEEE*, IEEE
651 (2014), pp. 413–417.

652 ⁵⁶<https://github.com/IBM-Bluemix-Docs/speech-to-text> (Last viewed: August 04,
653 2018).

654 ⁵⁷N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, “Front-end factor analysis
655 for speaker verification,” *IEEE/ACM Trans. on Audio, Speech, and Language Processing*

656 **19**(4), 788–798 (2011).

657 ⁵⁸I. Salmun, I. Opher, and I. Lapidot, “On the use of plda i-vector scoring for clustering
658 short segments,” Proc. Odyssey (2016).

659 ⁵⁹D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hanne-
660 mann, P. Motlicek, Y. Qian, P. Schwarz, *et al.*, “The kaldi speech recognition toolkit,”
661 in *IEEE Workshop on Automatic Speech Recognition and Understanding*, EPFL-CONF-
662 192584, IEEE (2011).

663 ⁶⁰D. Dimitriadis and P. Fousek, “Developing on-line speaker diarization system,” in *INTER-*
664 *SPEECH* (2017).

665 ⁶¹Z. Meng, L. Mou, and Z. Jin, “Hierarchical rnn with static sentence-level attention for
666 text-based speaker change detection,” in *Proceedings of the 2017 ACM on Conference on*
667 *Information and Knowledge Management*, ACM (2017), pp. 2203–2206.

668 ⁶²I. V. Serban and J. Pineau, “Text-based speaker identification for multi-participant open-
669 domain dialogue systems,” in *NIPS Workshop on Machine Learning for Spoken Language*
670 *Understanding*, Montreal, Canada (2015).

671 ⁶³R. Řehůřek and P. Sojka, “Software Framework for Topic Modelling with Large Corpora,”
672 in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*,
673 ELRA, Valletta, Malta (2010), pp. 45–50.

674 ⁶⁴F. Chollet *et al.*, “Keras,” <https://keras.io> (2015).

675 **LIST OF TABLES**

676	I	Acoustic features used in the regression model analysis.	17
-----	---	---	----

677 **LIST OF FIGURES**

678	1	(color online) Illustration of the proposed talker change detection (TCD)	
679		paradigm used in the present listening test study.	4
680	2	A comparison of talker attributes with respect to (a) fundamental frequency	
681		variation, and (b) word speaking rate. The vertical bars indicate one standard	
682		deviation spread around the mean value.	8
683	3	An illustration of a listening test trial.	10
684	4	(color online) (a) Illustration of human reaction time (RT) versus the ground-	
685		truth talker change instant (t_r vs t_c) across a total of 2720 trials (with $T_x \neq T_y$)	
686		over 17 subjects. The three inclined gray lines from bottom to top correspond	
687		to $t_r = t_c$, $t_c + 225$, $t_c + 2000$, respectively. NP stands for no button press.	
688		(b) Subject-wise summary using a boxplot of RTs in trials with hits. The	
689		black dots correspond to means. (c) Subject-wise miss and false alarm rates,	
690		and d-prime obtained from 200 trials for each subject.	12
691	5	(color online) Illustration of the distribution (obtained as a histogram) of the	
692		RT data for trials on which there was a hit for (a) raw RT data and (b)	
693		log-transformed RT data to improve the fit to a normal distribution.	14

694	6	Dependence of average RT on speech duration before the change instant. The	
695		black line is the linear regression fit.....	15
696	7	(color online) Proposed approach to model RT using acoustic features before	
697		and after change instant.	16
698	8	(color online) (a) Linear regression slope indicating the performance of in-	
699		dividual feature sets in accounting for the variance in human TCD reaction	
700		times. (b) The percentage of explained variance (r-square) for the linear	
701		regression model, averaged across subjects. (c) Illustration of significance (p-	
702		value) of different feature sets in multiple linear regression. (d) Scatter plot	
703		(polling all subjects' trials) of true RT versus estimated RT obtained from a	
704		multiple linear regression model computed across all feature sets.	20
705	9	(color online) Dependence of: (a) average RT on talker pairs (T_x-T_y), (b)	
706		average miss rate on talker pairs (T_x-T_y), and (c) average false alarm rate	
707		across talkers.	22
708	10	Proposed sequential diarization system as an on-line system for TCD.	23
709	11	Illustration of text-based context change classification.	26
710	12	(color online) Detection Error Trade-off (DET) curve for all the TCD systems	
711		evaluated on the stimuli set used in the listening experiments. The IBM	
712		Watson system is a single operating point. The human score is also marked	
713		in this figure for reference.	28

714 **REFERENCES**

715 ¹<https://gorilla.sc>.

716 ²See Supplementary material for supplementary experiments and results on change detec-
717 tion.

718 ³Sennheiser HD 215 II Closed Over-Ear Back Headphone with High Passive Noise Attenu-
719 ation.

720 ⁴<https://github.com/IBM-Bluemix-Docs/speech-to-text> (Last viewed: August 04,
721 2018).