

Mothers Reveal More of Their Vocal Identity When Talking to Infants

T. Kathiresan^{1,9,11}, A. Hervais-Adelman^{2,7,8}, S. Townsend^{3,4,7}, L. Dilley⁵, R. Shi⁶, M. M.

Daum^{2,7,8,10}, V. Dellwo^{1,7,11} *

¹Department of Computational Linguistics, University of Zurich, Switzerland.

²Department of Psychology, University of Zurich, Switzerland.

³Department of Comparative Linguistics, University of Zurich, Switzerland.

⁴Department of Psychology, University of Warwick, United Kingdom.

⁵Department of Communicative Sciences & Disorders, Michigan State University, USA.

⁶Department of Psychology, University of Quebec in Montreal, Canada.

⁷National Competence Centre of Research *Evolving Language*, Switzerland.

⁸Neuroscience Center Zurich, University of Zurich and ETH Zurich, Switzerland.

⁹Telepathy Labs, Zurich.

¹⁰Jacobs Center for Productive Youth Development, University of Zurich, Switzerland.

¹¹Centre for Forensic Phonetics & Acoustics, University of Zurich, Switzerland.

*Correspondence to: volker.dellwo@uzh.ch

1 Summary

2 Voice timbre – the unique acoustic information in a voice by which its speaker can be recognized – is
3 particularly critical in mother-infant interaction. Correct identification of vocal timbre is necessary in
4 order for infants to recognize their mothers as familiar both before and after birth¹⁻⁶, providing a basis
5 for social bonding between infant and mother. The exact mechanisms underlying infant voice
6 recognition remain ambiguous and have predominantly been studied in terms of cognitive voice
7 recognition abilities of the infant. Here, we show – for the first time – that caregivers actively
8 maximize their chances of being correctly recognized by presenting more details of their vocal timbre
9 through adjustments to their voices known as infant-directed speech (IDS) or *baby talk*, a vocal
10 register which is wide-spread through most of the world’s cultures^{4,7,8}. Using acoustic modelling (*k*-
11 means clustering of Mel Frequency Cepstral Coefficients) of IDS in comparison with adult-directed
12 speech (ADS), we found in two cohorts of speakers - US English and Swiss German mothers - that
13 voice timbre clusters of in IDS are significantly larger to comparable clusters in ADS. This effect
14 leads to a more detailed representation of timbre in IDS with subsequent benefits for recognition.
15 Critically, an automatic speaker identification using a Gaussian-mixture model based on Mel
16 Frequency Cepstral Coefficients showed significantly better performance in two experiments when
17 trained with IDS as opposed to ADS. We argue that IDS has evolved as part of an adaptive set of
18 evolutionary strategies that serve to promote indexical signalling by caregivers to their offspring
19 which thereby promote social bonding via voice⁹ and acquiring linguistic systems¹⁰.

20 RESULTS

21 A predominant challenge in voice recognition relates to the immense variability in timbres of voices
22 within individuals, arising from a large spectrum of sources, such as varying registers, vocal effort or
23 different emotions¹¹⁻¹⁹. This variability leads to reduced recognition performance for a voice learned
24 in one register and presented in another. IDS is a register which shows some of the highest within-
25 and between-speaker acoustic variability (see review in²⁰). Recently, Piazza et al.⁴ demonstrated
26 consistent within-speaker shifts in timbre from ADS to IDS, which they interpreted as changes in the
27 unique fingerprint within a caregiver's voice. They argued it helps the infant differentiate between
28 ADS and IDS. However, if adults indeed shifted their unique vocal fingerprint from ADS to IDS, this
29 should correspond to a perceived change in voice identity with detrimental consequences for
30 recognition performance in cross-register recognition (for example, when a voice is learned under IDS
31 and recognition takes place under ADS). There is, however, consent that IDS registers are applied to
32 support attention and thus facilitate social bonding between caregiver and infant²¹⁻²⁴. Social bonding,
33 however, requires a solid recognition of the social interactants, notably by their voices¹⁻⁶. We contend
34 that it is biologically implausible that caregivers make recognition challenging by addressing infants
35 in a register that will not serve to make them more generally recognizable. Additionally, IDS features
36 some of the highest within-register timbre variability of all vocal registers, manifested as tremendous
37 exaggerations of acoustic segmental^{10,25-28} and prosodic characteristics²⁹⁻³⁴. While within-person
38 variability – within or between register – has generally been thought to oppose identity recognition³⁵⁻
39³⁷, results from automatic voice recognition¹⁸ and – more recently – from human voice recognition³⁷⁻
40⁴⁰ revealed that learning identity from voices containing more within-speaker detail, can generalize
41 better to unknown variability in the long-term and thus lead to more robust recognition. It thus seems
42 plausible that the high degree of within-speaker variability observed in IDS might contribute
43 positively to recognition when a system has been trained with it. Here, we tested this hypothesis using
44 acoustic modelling and automated speaker recognition models. To test this, we analysed speech from
45 mothers of two independent cohorts (10 Swiss-German and 27 US-English) reading sentences to their

46 own infants (IDS) and to an adult experimenter (ADS), in recording facilities located in the Upper
47 Midwest, USA and University of Zurich/Switzerland respectively (see STAR methods). Given that
48 infants exhibit a preference for the mother's over the father's voice⁶⁹ and that they show a lack of
49 discrimination between male voices (including the paternal voice⁷⁰), we have – in line with previous
50 studies on IDS – limited our investigation to mothers' voices.

51 **Acoustic modelling: ADS is a timbre subspace of IDS**

52 To better understand the acoustic nature of the timbral shift suggested by Piazza et al. ⁴, we modelled
53 this shift in a multidimensional timbre space. The timbre of a voice that forms its identity consists of a
54 set of multidimensional acoustic variables that change rapidly over time. A typical way of modelling
55 timbre is by using Mel Frequency Cepstral Coefficients (MFCCs, ⁴¹), a series of speech coding
56 techniques arriving at a discrete number of coefficients (typically 13) specifying voice-specific timbre
57 during a specified time unit (typically 25 ms). Clusters of MFCCs in the multi-dimensional space are
58 an efficient way of identifying speaker-specific timbre, henceforth *timbre clusters*. To understand how
59 an ADS timbre space of a speaker compares to the IDS space of the same speaker, we first identified
60 32 timbre clusters in a 13-dimensional MFCC space using *k*-means clustering. We applied the
61 clustering to ADS and IDS combined for each speaker. 32 clusters were chosen because this was the
62 largest number for which each cluster contained sufficient data samples in both IDS and ADS from
63 the same speaker to measure 95% confidence ellipses in dimension-reduced data ⁴². Second, we
64 divided each cluster into ADS and IDS sub-clusters and measured the Euclidian distance between the
65 sub-cluster centroids in the 13-dimensional MFCC space. Fig. 1. (A) illustrates the median distance
66 between IDS and ADS in the timbre space, which was about four times higher than the distance
67 between a random data-split with 5-fold cross-validation at the proportion of the ADS/IDS ratio
68 within a cluster; a linear mixed model showed a significant effect of split type, $F[1,1] = 1425.21, p <$
69 001 , but not of language, $F[1,1] = 0.025, p = .87$; and no interaction: $F[1,1] = 0.33, p = .56$). This

70 shows that, within a timbre cluster, ADS and IDS are centred around different means, consistent with
71 the proposal of a shift in Piazza et al.⁴.

72 To compare the relative extent of the two subclusters in the timbre space, we obtained the variance of
73 the ADS and IDS subgroups within each cluster by applying Principal Component Analysis (PCA) on
74 the MFCC data. Using the first two principal components, we measured the area for IDS and ADS
75 data by fitting a 95% confidence ellipse. Fig. 1. (B) shows the area of ADS and IDS sub-clusters in a
76 representative example of one single *k*-means cluster of a Swiss speaker; areas of IDS and ADS in all
77 clusters for both languages are summarized in Fig. 1. (C). Both languages show a proportionally
78 smaller area for ADS compared to IDS (pairwise t-test with Bonferroni correction for Swiss: $p =$
79 0.0019 ; US: $p < 0.0001$). The difference was significantly larger for US English compared to Swiss
80 German, as revealed by a significant interaction between register and language (stemming from a
81 linear mixed model [register*language] with [speaker] as random factor: $\beta = -0.9641$; ANOVA: $F[1,$
82 $2329] = 9.3711$; $p = 0.00223$).

83 In addition to demonstrating a method by which the timbral shift suggested by Piazza et al.⁴ can be
84 rigorously quantified, the above results crucially identify and confirm a systematic increase in the
85 extent of size of the timbre space occupied by IDS register compared to ADS register. This systematic
86 increase of the timbre space in IDS compared to ADS holds across two distinct languages/dialects –
87 Swiss German and US Midwestern English – involving very disparate geographic regions and
88 populations of individuals. Other than a mere shift between ADS and IDS – as argued by Piazza et al.⁴
89 – this means that the potential vocal space a speaker occupies is better represented in IDS compared to
90 ADS registers, and further, that IDS includes the timbre space of ADS significantly more than the
91 other way around. We thus obtained a rigorous acoustic model that explains the variability advantage
92 demonstrated in prior studies in both computer¹⁹ and human³⁸ voice recognition.

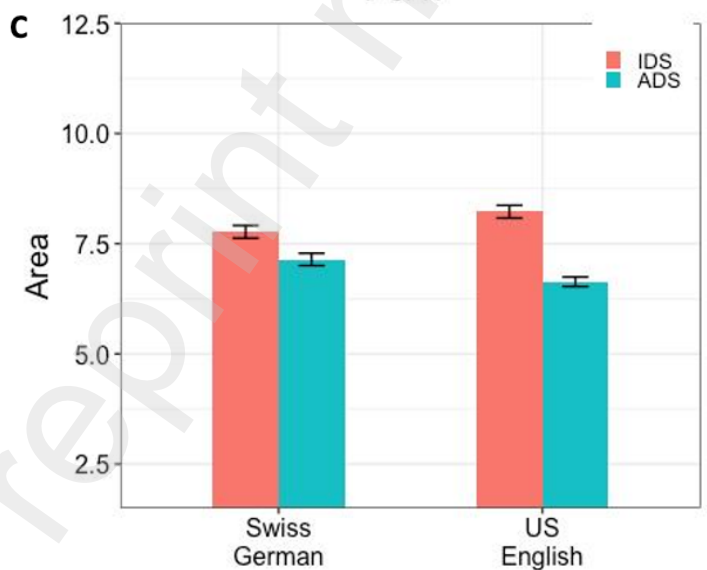
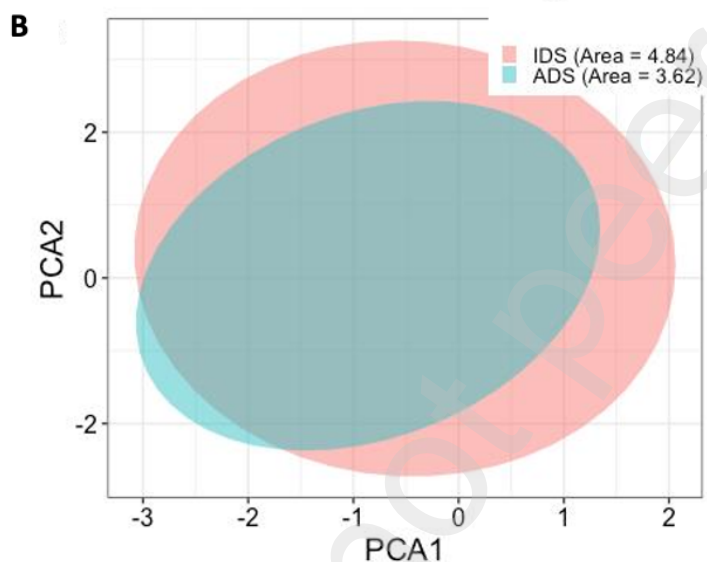
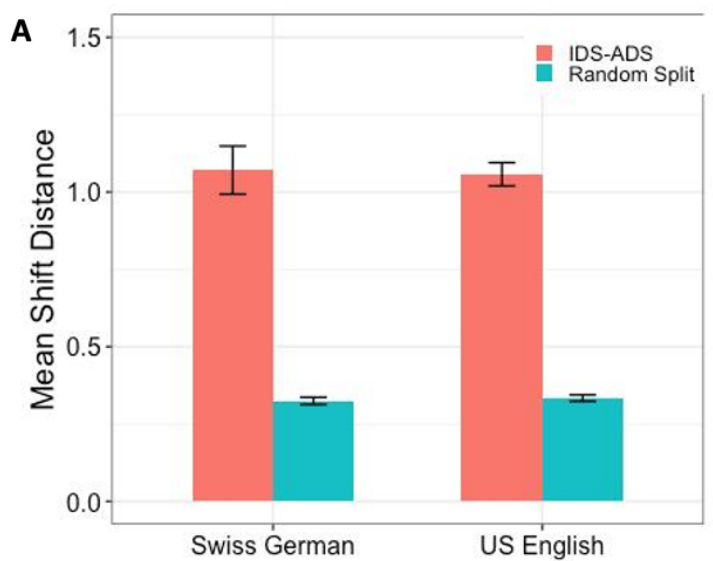


Figure 1. (A) Within-cluster differences between ADS and IDS centres (red) and differences between random within-cluster splits (turquoise). (B) Representative example of IDS (red) and ADS (turquoise) areas in one GMM of one speaker in Swiss German, (C) mean k -means cluster area (averaged \pm 1 stdev over 10 and 27 speakers respectively).

93 **Expanded timbre space leads to a recognition advantage**

94 The finding of an expanded timbre space in IDS compared to ADS has direct implications for voice
95 recognition. Because IDS contains more timbre detail compared to ADS, we hypothesised that
96 learning a voice under the larger IDS timbre space should lead to better recognition performance of
97 voices under ADS compared to learning a voice under ADS and recognizing the speaker under IDS
98 (henceforth: IDS-ADS and ADS-IDS respectively). Recognition was tested based on a Gaussian-
99 Mixture Model (GMM) with 32 mixtures constituting the number of mixtures for which almost no
100 empty clusters in either IDS or ADS were obtained, along with 5-fold cross-validation based on 13-
101 dimensional MFCCs (see Fig. 2). A two-way ANOVA [accuracy by language and speaking-register]
102 had a significant interaction between the effects of language and speaking-register [$F [1,16] = 80.92$,
103 $p_{(\text{Bonferroni})} < 0.001$], indicating effect magnitudes were somewhat larger in US English compared to
104 Swiss German. Testing the simple effect of training-test sequence (IDS-ADS; ADS-IDS) revealed
105 significantly better IDS-ADS recognition in both languages, as shown by a pairwise *t*-test with
106 Bonferroni correction for Swiss: $p = 0.001$; US: $p < 0.001$). Results were replicated (see STAR
107 methods) for different numbers of mixtures (32 to 512) and including temporal features (first- and
108 second-order derivatives; delta and delta-delta). To make sure that the model performs well under
109 within-register conditions, we also tested ADS-ADS and IDS-IDS. Mean within-register recognition
110 led to ceiling effects for ADS (Swiss: 100%, US: 98%) and close to ceiling performance for IDS
111 (Swiss: 99%, US: 88%). This was to be expected as GMMs perform highly on within-register read
112 speech for our group sizes. We relate the decrease in performance within IDS recognition –
113 particularly for US-English speakers – to the generally higher variability in both training and test data.

114 Results suggest that learning vocal identity from acoustically more variable IDS may lead to speaker
115 recognition benefits in other, less variable registers (here ADS), but not vice versa. While using less
116 variable registers in training and testing (e.g., read ADS-ADS) leads to optimal results, this scenario is
117 not very applicable, as listeners would only be prepared to identify speakers under this particular low-

118 variability condition, which is improbable to occur. In real-life situations of kinship recognition,
119 listeners task is to identify speakers under a wide range of within-speaker timbre variability. Thus, by
120 addressing an infant in IDS, the infant may learn critical information about timbre features that serve
121 to identify speakers in ADS and possibly other registers, but the reverse does not hold. It supports the
122 view that the higher amount of attention infants pay to caregivers under IDS²¹⁻²⁴ is accompanied by
123 caregivers revealing more information about their voice timbre variability which in return leads to the
124 acquisition of a more robust representation of that individual's voice. In sum, it suggests that the
125 acoustics of IDS play a crucial role in furnishing the infants with the necessary cues for deriving a
126 highly-generalisable speaker identification system.

127 The fact that the impact of register on both recognition performance (Fig. 2) and timbre differences
128 (Fig. 1) was higher in US English compared to Swiss German may potentially be explained as
129 differences in the circumstances of IDS/ADS production across the recording sessions in different
130 countries/languages (Swiss and US). In comparison to American mothers, who spoke to adults in the
131 absence of their infants, Swiss German mothers carried their baby while they were producing adult
132 directed speech and read a children's story to the adults. Even though Swiss listeners showed a high
133 sensitivity in differentiating IDS from ADS in the Swiss speakers (mean A' across 31 listeners = 0.79;
134 see STAR methods), Swiss ADS had notable IDS influences. In addition, differences between ADS
135 and IDS registers are continuous rather than categorical, and specific acoustic attributes and
136 prevalence differ according to cultural background, geographic region, and individual identity⁴³⁻⁴⁵.
137 We therefore posit that the reason for the smaller effect magnitudes in Swiss speakers may also
138 plausibly reflect the use of a generally less marked IDS register by the Swiss German speaker group
139 compared to the US group.

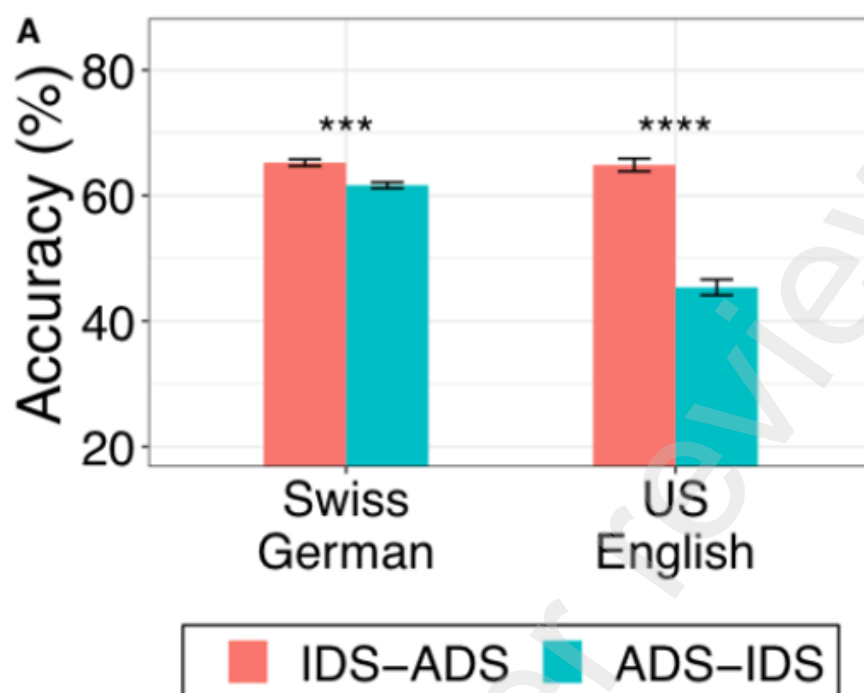


Figure 2. Recognition accuracy for GMM models trained on ADS and tested with IDS (ADS-IDS) and the reverse (IDS-ADS) for both US-English (N=27) and Swiss-German (N=10) speakers (average \pm 1 stdev over 5 folds of cross-validation).

140 **Learning a speaker's voice from IDS results in a recognition advantage in other registers**

141 In a second recognition experiment we used an identical test set to probe recognition performance

142 using the IDS and ADS recognition models from the cross-register condition (previous section) to

143 recognize spontaneous speech elicited in a picture description task and an interview situation. This

144 register was available for the Swiss German speakers only. Both tasks were carried out under IDS and

145 ADS conditions (see STAR Methods). This speech was not subject to the typical constraints of read

146 speech and contained a large variety of different lexical items and grammatical patterns. Fig. 3 (A)

147 shows that the entire timbre space (13-dimensional MFCCs) in a dimension reduced PCA is more

148 extensive than, and mostly occupies a different space to, the read speech of the previous experiment.

149 Recognition performance showed a significant interaction between training set (IDS, ADS) and

150 spontaneous speech type (free, picture-naming) on accuracy [two-way ANOVA accuracy*test set:
 151 $F [1,16] = 25.03, p < 0.0001$]; (Fig. 3. B). Simple effects (Bonferroni corrected p-values) of training
 152 register for each of free and picture-naming speech revealed significantly higher recognition
 153 performances for IDS compared to ADS-trained models (free-speech: $p < 0.043$; picture-naming: $p <$
 154 0.0001). These results provide strong evidence for the view that IDS compared to ADS training
 155 generalizes better to vocal registers that are not part of either of the training registers. Hence, the
 156 timbre variability obtained in IDS is likely to be of advantage to recognize voices under a wide variety
 157 of within-speaker timbre variants. We thus conclude that training a listener with IDS instead of ADS
 158 will prepare them to recognize the speaker under possibly all vocal settings and registers of that
 159 speaker.

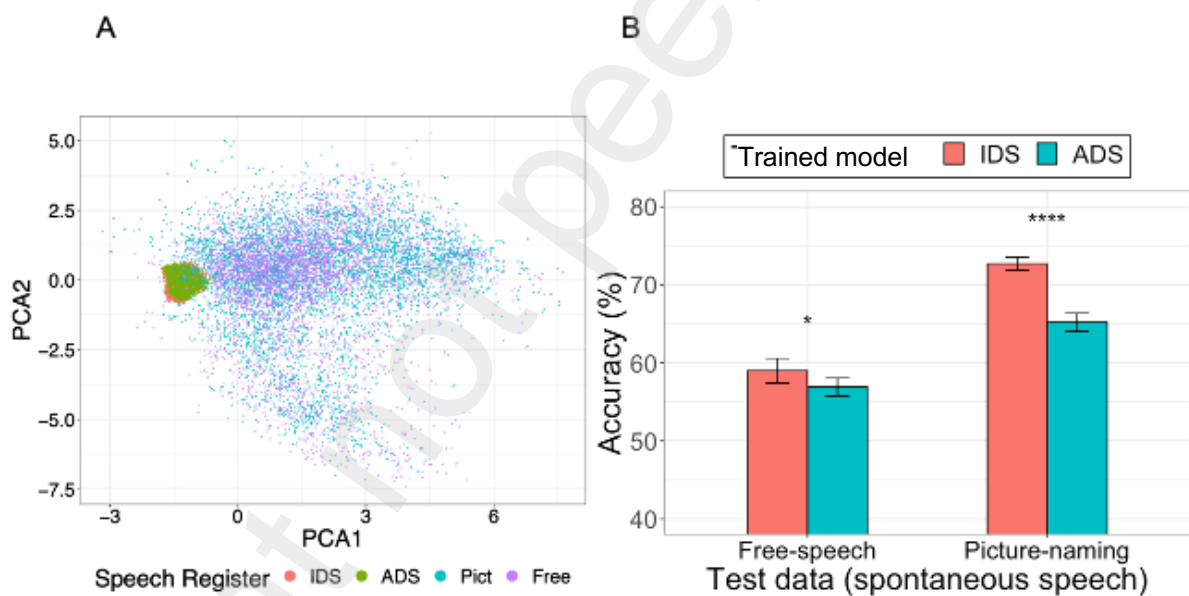


Figure 3. (A) Example of read (ADS, IDS) and spontaneous speech (Picture naming, and free-speech) of one Swiss-German mother in a PCA space. Both spontaneous registers occupy a different timbre space compared to the read registers (B) The trained IDS and ADS voice recognition models of Swiss-German (N=10) mothers separately were tested on two varieties of spontaneous speech 1) free-speech, and 2) picture-naming. The experiment was performed over 5 folds of cross-validation.

160 **DISCUSSION**

161 Voice recognition is a process that is essential for social cohesion of humans and many other voice-
162 recognizing animals, and behaviorally has been shown in infants both before and shortly after birth
163 ^{46,47}. Here, we revealed that the universal phenomenon of IDS typically produced by caregivers to
164 address their infants is a register that supports voice recognition performance when voices are learned
165 under IDS. Using rigorous acoustic modelling we demonstrated – for the first time – that voice-
166 specific timbre clusters capture more of the possible acoustic space in IDS compared to ADS, leading
167 to better speaker recognition when caregivers speak in IDS. As a result, the acoustics of IDS provide
168 optimal training input for the development of a robust voice identification system, one which has
169 adapted to the tremendous situational variability in a caregiver’s voice that the infant will encounter.
170 It is thus likely that IDS is an adaptation applied by parents in infancy to foster acquisition of voice
171 timbre detail by infants to subsequently support robust recognition of the caregiver by the infant, thus
172 providing a basis for close bonding and increased chances of survival.

173 Our results are in line with recent findings from both face^{35–37} and voice recognition [22],^{37–40},^{37–40}
174 revealing that viewing a face from different perspectives or listening to a voice of a speaker under
175 increased variability supports subsequent face or voice recognition, respectively. This is in line with
176 results showing that rich multi-modal information results in better representations of affective
177 categories⁴⁸. Similar evidence also exists from automatic voice recognition, demonstrating that
178 training data from speakers under a variety of speech registers increases correct recognition
179 probability of test data containing high variability¹⁸. This means that maximizing knowledge about
180 within-voice variability supports the receivers’ ability to identify a voice under a wide variety of
181 registers and generalizes to previously unexperienced realizations of particular registers by a speaker.
182 Here, we showed for the first time that the augmentation of within-speaker variability that is present
183 in IDS fosters and supports enhanced speaker recognition as a result of a wide spectrum of within-

184 voice variability of the caregiver. Given a high-variability advantage detected for types of adult
185 directed speech³⁸, it seems plausible that IDS evolved for precisely such scenarios.

186 Within-speaker variability has typically been viewed as detrimental for recognition because the
187 detection of similarity in stimuli with high variability is more challenging^{49,50}. However, our results
188 highlight the importance of variability as advantageous when supplied sufficiently during the training
189 period³⁸. How long then does a familiarisation with IDS need to last to show advantages over speech
190 that consists of low-variability vocalisations? While the answer to this question is subject to future
191 studies, it seems plausible that after familiarization with IDS over the first year of life, the infant will
192 be well prepared to perform caregiver recognition by the time autonomous mobility (i.e., independent
193 crawling) starts and the self-exploration of different social settings consisting of new voices increases.
194 This is probably the point at which robust recognition of the primary caregiver is most crucial. Thus,
195 the point at which IDS changes into so-called child-directed speech (CDS) – which should then
196 contain less indexical signalling – may also be around that time. In fact, there is evidence for the
197 change between IDS and CDS to be indeed around that stage in the infant's life^{51,52, 53,54}.

198 Previous studies have predominantly focused on the hypothesis that the vocal adjustments in IDS
199 support language acquisition, in particular the process by which abstract linguistic forms (e.g.,
200 phonemes, syllables, or words) are acquired from variable acoustic representations¹⁰. However, the
201 type of timbre variability found in IDS introduces highly non-canonical linguistic realisations^{55–57}
202 which have been found to produce ambiguous category information, thus, arguably,
203 counterproductive for the acquisition of linguistic categories⁵⁸ and shown to lead to less discriminable
204 utterances at the phoneme level⁵⁹. A recent machine learning study⁵⁷ found that realisations of IDS
205 phoneme categories were neither more separable nor more robust in comparison to their ADS
206 counterparts. This is entirely consistent with the present results, as increasing voice timbre detail
207 should increase the variability found between segments, rendering them less canonical. Thus, it is
208 plausible that the dominant function of IDS is to elicit infant's attention to voice specific aspects.

209 Our results have wide-reaching implications. In evolutionary terms, it seems plausible that IDS has
210 evolved as part of an adaptive set of strategies that serve to promote indexical signalling by an adult to
211 their offspring, thereby promoting adult-infant bonding, increasing the infant's attention, and fostering
212 offspring survival. Given the relatively poor linguistic skills of new-borns and infants in the first few
213 months *post partum*, it seems plausible that the dominant role of voice in communication would be
214 transmitting identity. This would also make offspring-directed vocalisations a plausible phenomenon
215 in voice recognizing animals, where caregiver recognition may be similarly important. In fact,
216 Fernandez and Knörnschild ⁶⁰ recently provided evidence for the existence of an analogue of IDS in
217 bats, which were shown to adopt a more highly variable vocal register when directing their calls to
218 pups than adult conspecifics. Nevertheless, Fernandez and Knörnschild ⁶⁰ automatically classified
219 individual caregiver signals and found low classification rates. They concluded that pup-directed calls
220 “do not seem to encode sufficient interindividual variation to allow for reliable individual
221 discrimination” (p. 5). This, however, is not surprising, given that the advantage of higher variability
222 signals only arises after long-term training ^{49,50}. The pup-directed vocalizations reported in ⁶⁰ from 13
223 female pups – which evidently contain a high timbre variability – must be expected to result in poor
224 classification performance and this result directly replicates findings for human speech ⁶¹ and our
225 results from the within register tests – ADS-ADS and IDS-IDS – in which the higher variability of
226 IDS has disadvantages over ADS recognition (see above). Hence, a lower acoustic variability in adult-
227 directed bat vocalizations would mean that their classification rates should be higher compared to
228 pup-directed vocalizations (classification results not reported in ⁶⁰). Bat pups, however, are evidently
229 good at identifying their caregivers ⁶² and in the light of the results reported here, it is plausible that
230 this performance is supported by the higher variability vocalizations reported in ⁶⁰. We predict that
231 using pup-directed bat vocalizations as training material for a recognition task would show similar
232 training advantages to the human data reported here. It will be interesting to see whether offspring-
233 directed adjustments in vocalisations occur in other vocal recognizing animals (e.g. ⁶³⁻⁶⁶) and whether
234 these vocal adjustments co-occur with individual recognition advantages. Given that both human and

235 non-human animals share this non-linguistic function of vocalisations, it is plausible that vocalising to
236 bootstrap identification may be the ancestral function of infant-directed vocalisations which was later
237 leveraged to aid the development and realisation of linguistic contrasts in human speech ^{67,68}.

238 **ACKNOWLEDGEMENTS**

239 Andrea Fröhlich and Sarah Lim carried out the recordings of the Swiss speakers and Marco Bleiker
240 carried out subject recruitment. Meisam K. Arjmandi assisted with implementing analyses on the US
241 English dataset at Michigan State University. We also acknowledge Tonya Bergeson, Derek Houston,
242 and Maria Kondaurova for assistance with the original US English dataset collection at the Devault
243 Otologic Research Lab. This research was supported by grant #100012_125399 and grant NCCR
244 Evolving Language #51NF40_180888 of the Swiss National Science Foundation.

STAR METHODS

Mothers Reveal More of Their Vocal Identities When Talking to Infants

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Software and Algorithms		
MATLAB R2019b	Mathworks	RRID: SCR_001622
Voicebox Speech processing toolbox	[Brookes, 2011]	http://www.ee.ic.ac.uk/hp/staff/dmb/voicebox/voicebox.html
HTK MFCC toolbox	[Davis and Mermelstein, 1980; Young et al., 2006]	https://ch.mathworks.com/matlabcentral/fileexchange/32849-htk-mfcc-matlab
R Studio	[R Core Team, 2013]	http://www.R-project.org/
ggplot2	[Wickham H, 2016]	https://ggplot2.tidyverse.org .

CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources should be directed to and will be fulfilled by the corresponding author, Volker Dellwo (volker.dellwo@uzh.ch).

EXPERIMENTAL MODEL AND SUBJECT DETAILS

US-English materials: 27 mother–infant dyads were recruited from the local community in Indianapolis, Indiana, USA; all mothers were native US English speakers with self-reported normal hearing who grew up in the Midwestern United States. Mothers participated in the study with their infants. Recordings were randomly selected from among those available such that infants were on average 7.5 months old (range: 4.1 - 9.9 months). The gender of infants was balanced. This research and the recruitment of human subjects were approved by the Indiana University Institutional Review Board. During a lab visit, mothers were asked to read for both IDS and ADS conditions from a storybook. In the IDS speech condition, mothers read the storybook to their infants, and the infants were present in the room with their mothers for the duration of the IDS session. In the ADS condition, mothers were asked to read aloud as if to another adult. The order of IDS and ADS recordings was counterbalanced across mothers and visits. For additional details, see ⁴⁴

263 *Swiss-German materials:* 10 Swiss-German mother–infant dyads participated in the experiment. All
264 the mothers were native Swiss-German speakers with self-reported normal hearing. Infants were 8-9
265 months old, and the gender was balanced. Informed consent was obtained from all participating
266 mothers, and approval of the study was obtained from the University of Zurich. Participants were told
267 that we were broadly interested in “mother-child interaction” and were not aware that we were
268 measuring differences between the acoustic properties of their speech across the conditions. The
269 participants read 12 sentences to their infant, and to an adult experimenter in their own reading
270 register.

271 **METHOD DETAILS**

272 **Equipment**

273 *US-English Recording:* Mothers were digitally recorded reading in a double-walled, copper-shielded
274 sound booth (Industrial Acoustics Company, Bronx, NY) at the DeVault Otologic Research Lab in
275 Indianapolis, Indiana, USA. The speech was recorded in one of two ways: The initial system used an
276 Audio-Technica ES933/H hypercardioid microphone (Audio-Technica, Leeds, UK) powered by a
277 phantom power source and linked to an amplifier (DSC-240; Daqscribe, Centennial, CO) and a Sony
278 DTC-690 digital audio tape recorder (Sony, Tokyo, Japan). The equipment was updated partway
279 through this longitudinal project to an SLX Wireless Microphone System (Shure, Niles, IL). This
280 system included an SLX1 Bodypack transmitter with a built-in microphone and a wireless receiver
281 SLX4, which was connected to a Canon 3CCD Digital Video Camcorder GL2, NTSC
282 (Canon, Melville, NY) and recorded the speech samples directly onto a Mac computer (OSX Version
283 10.4.10; Apple, Inc., Cupertino, CA) via Hack TV (Version 1.11) software. No systematic differences
284 were found across recording sessions or participant groups in terms of recording technology.
285 Recordings were made at a sampling rate of 22050 Hz with 16-bit quantization rate. See ⁴⁴ for more
286 details.

287 *Swiss-German Recording:* Speech data were recorded continuously using Sennheiser (model: MKE 2
288 P) omni-directional prepolarized condenser clip-on microphone. The microphone was attached to

289 each mother's shirt collar. The sound recordings were done in a sound-treated room in the psychology
290 department's baby lab led by author 6 at University of Zurich. Recordings were made at a sampling
291 rate of 48 kHz Hz with 16-bit quantization rate. To obtain the same reading data for mothers and
292 infants we recorded a small passage from a child-book in both the IDS and ADS conditions.

293 **Quantification and Statistical Analysis**

294 Within-speaker vocal variability can occur at different levels (e.g., phonemic, acoustic feature,
295 speaking register; cf. ³⁸), and there is evidence for top-down recognition advantages through multiple
296 levels of linguistic and paralinguistic features contributing to the voice recognition process ⁷¹. Here we
297 limit voice recognition to the process of recognizing a speaker based predominantly on the bottom-up
298 acoustic properties of voice, i.e., independent of lexical choices or other language related properties.
299 This is consistent with infants' limited knowledge of the higher-order linguistic properties of the
300 environmental language, which requires that they must rely predominantly on acoustic characteristics
301 of voice alone to recognize their caregivers.

302 **Feature Extraction**

303 MFCCs were obtained by block through processing the speech segment using a 25 ms Hamming
304 window with an overlap of 10 ms, along with parameters that included pre-emphasis coefficient
305 (0.95), number of filter bank channels (20), and liftering parameter (22), spectral bandwidth (150 -
306 8000 Hz) using MATLAB [Young et al, 2011]. The 0th order coefficient (energy coefficient) were
307 included in the 13-dimension MFCC features. The first- and second-order temporal derivatives were
308 computed from the extracted MFCCs. Finally, 39 dimensional features (13 MFCCs, 13 first-order
309 derivatives, and 13 second-order derivatives) per frame were used to develop a speaker model by
310 training the Gaussian mixture models discussed in the ASR section.

311 **Number of Clusters in the GMM**

312 To compare the data spread of IDS compared to ADS within a cluster, we needed an optimal situation
313 in which all clusters are filled with data from each category. Empty clusters were defined as clusters
314 containing less than 4 data points in either IDS or ADS, corresponding to the minimum number of

315 points to calculate the data spread in 2-dimensional space, below. We performed a k -means cluster
316 analysis with varying cluster numbers (16, 32, 64, 128, 256, 512); we found that from 64 clusters
317 upwards, empty clusters start occurring. Interestingly, the number of empty clusters in ADS was
318 always higher than in IDS. Given this finding, we used the 32 cluster analyses which reflected the
319 point where empty clusters were equal between the two categories; we then measured several cluster
320 properties: (a) distance between cluster means, (b) weight and (c) covariance.

321 **Automatic Speaker Recognition (ASR)**

322 We used GMM with diagonal covariance matrices [Reynolds et al., 2000] to design supervised
323 language-independent speaker recognition models. In general, for GMM-based automatic speaker
324 recognition experiments, the optimal number of Gaussian mixtures were estimated by varying the
325 count (2^N ; $N = 1, 2, 3, 4, \dots$) until the best performance was obtained, which turned out to be 32 ($N = 5$)
326 for most of the experiments. The speaker recognition model was designed for every speaker on both
327 speaking registers separately, i.e., IDS and ADS per speaker. The extracted MFCCs with delta and
328 delta-delta (for both, US-English and Swiss German dataset separately) were randomly split into an
329 80-20 ratio for training and testing. The data split and the ASR experiment was repeated using 5-fold
330 cross-validation.

331 **Indexical Feature Expansion**

332 The 13-dimensional MFCC features (without delta and delta-delta) were normalized by speaker (z -
333 score normalization within speakers). The normalized features were clustered using a k -means vector
334 quantization algorithm [Voicebox tool in MATLAB, Brookes, 2011]. The cluster count was varied
335 (shown in Figure S1. and Figure S2) to find the highest number of clusters (k) for which a minimal
336 number of empty clusters (i.e. clusters with only one data sample) occur. To measure the data spread
337 of a cluster, it is necessary to have more than four data samples in a cluster [Stat_Ellipse() function in
338 ggplot2 R package]. At $k = 32$ clusters, there were 1 and 3 empty clusters found in Swiss-German and
339 US-English datasets respectively. Solutions for $k > 32$ had notably more empty clusters, hence our
340 choice was to use 32 clusters for subsequent analyses. After clustering, every data frame was assigned

341 a cluster label (1 - 32) and the speech register label (IDS or ADS). The 13-dimensional MFCC feature
342 space was reduced to 2 dimension using principal component analysis (PCA). The surface area of IDS
343 and ADS points in each cluster was calculated in terms of their ellipse surface area (πab) in the 2-
344 dimensional feature space (normal distribution assumed) and was used as a measure of timbre
345 expansion in each sub-cluster.

346 **Control Analysis: Behavioural Classification of IDS Versus ADS**

347 In the Swiss data retrieval, babies could often not be separated from their mother during the adult
348 recording. This resulted in vocalisations that were audibly less well separated on an IDS-ADS register
349 continuum. To make sure that individual productions of IDS and ADS were categorizable, we tested
350 the sensitivity (A') of 30 native listeners of Swiss German to identify individual productions as either
351 IDS or ADS. The stimulus set consisted of 12 sentences read by 4 speakers in 2 registers ($N = 240$).
352 This set was randomly divided into 3 subsets of $N=80$ each (4 sentences * 10 speakers * 2 registers).
353 Each 10 randomly selected listeners listened to one of the three stimulus sets. The average A' revealed
354 a high sensitivity of 0.79 and no listener bias ($B''D = 0.1$).

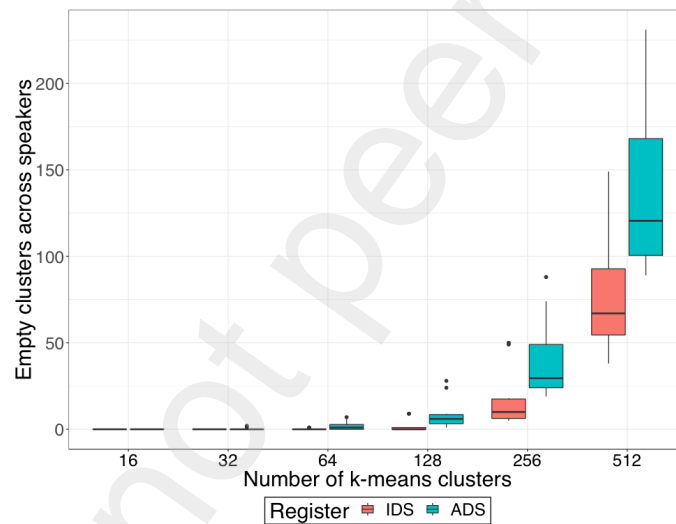
355 **Control Analysis: Automatic Classification of IDS Versus ADS**

356 To ensure that the data is comparable to previous studies, we replicated the results from Piazza et al.
357 (2017) using a support-vector machine classifier (SVM-RBF) [MATLAB SVM classifier
358 Application] for IDS-ADS classification based on 12 MFCCs (frames were not averaged) with a
359 training/test split of 90/10 (5-fold cross-validation) on two datasets. Piazza et al. trained a linear SVM
360 with ADS and found that speaker classification performance was decreased for IDS vs ADS. Based
361 on this finding, they conclude that there is a systematic shift in timbral properties of the two registers.
362 Our data replicates this finding with 70.3% and 75.8% correct classification rates for US English and
363 Swiss-German respectively, consistent with the average classification of $\sim 70\%$ in ⁴ for English and
364 non-English data. Importantly, failure of machine classification alone does not necessarily reflect a
365 systematic acoustic timbre difference (shift) between ADS and IDS.

366 **Supplementary Text**

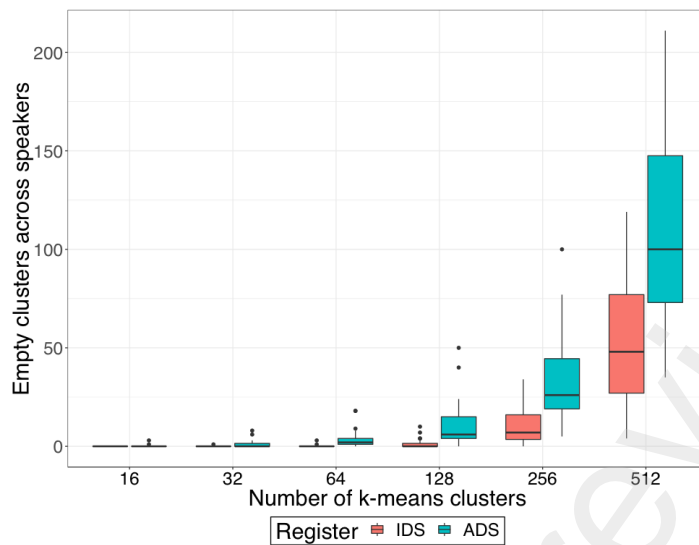
367 All the ASR experiments (Figure S3 and S4) have shown a similar pattern in the results i.e., the
368 recognition accuracy of the model trained on IDS and tested with ADS is always better than the
369 recognition accuracy of the model trained on ADS and tested on IDS independent of the features used,
370 and the number of GMM mixtures varied. Concerning features, MFCCs with delta have shown a
371 better mean accuracy than MFCCs alone in both datasets. The Swiss-German dataset revealed better
372 mean accuracy than the US-English dataset. Overall, the recognition accuracy of the two datasets has
373 declined with increase in GMM mixtures except for the MFCCs with delta in the Swiss-German
374 dataset (Figure S4 b). In Figure S4 b, the peak performance was at 128 GMM mixtures but, for other
375 experiments, the peak performance was at 32 GMM mixtures.

376 **Figures S1-S6**



377

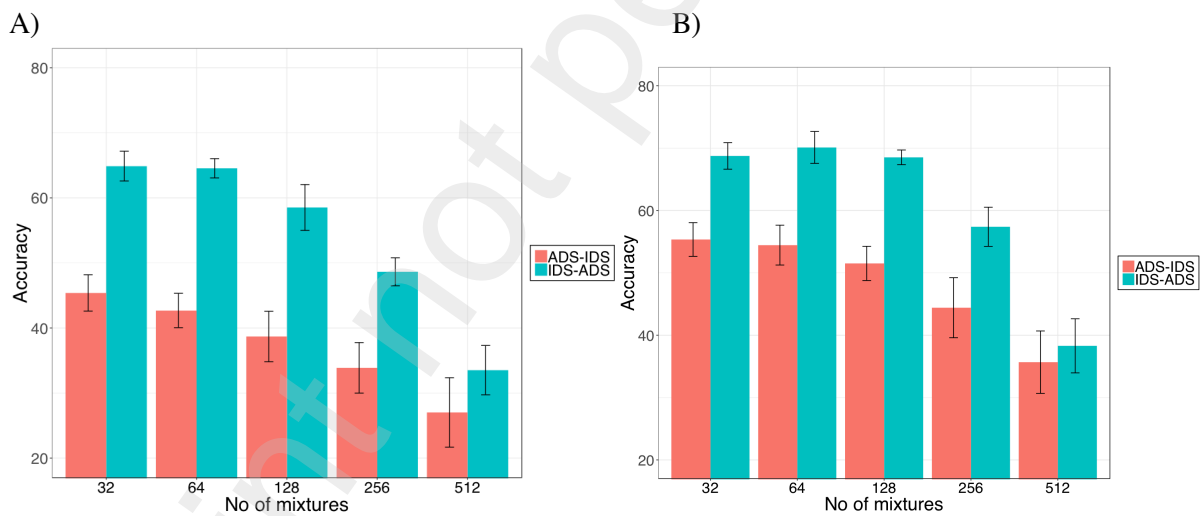
378 **Figure S1.** Empty-clusters in both IDS and ADS registers of 10 Swiss-German mothers in k -means
379 clustering. The x -axis shows the number of k -means clusters used in the analysis, and the y -axis shows
380 the variation of empty-clusters across 10 speakers. The vertical bar in the graph shows the standard
381 deviation.



382

383 **Figure S2.** The occurrence of empty-clusters in both IDS and ADS registers of 27 US-English mothers
 384 in k -means clustering. The x -axis shows the number of k -means clusters used in the analysis, and the y -
 385 axis shows the variation of empty-clusters across 27 speakers. The vertical bar in the graph shows the
 386 standard deviation.

387



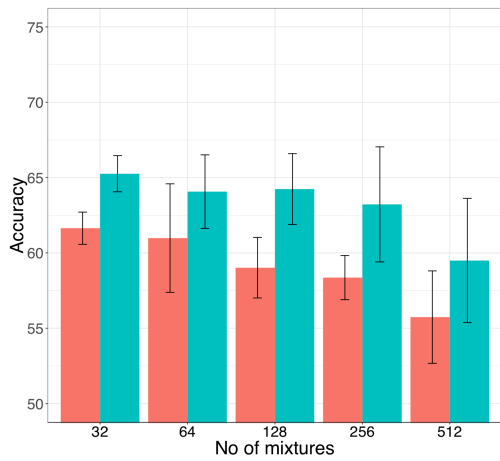
388

389 **Figure S3.** Cross-register voice recognition accuracy (mean and standard deviation) for US-English
 390 dataset using a) MFCCs, and b) MFCCs with delta. The x -axis shows the number of mixtures in the
 391 GMM classifier, and the y -axis shows the recognition accuracy in %. The vertical bar in the graph shows
 392 the standard deviation of 5-fold cross validation.

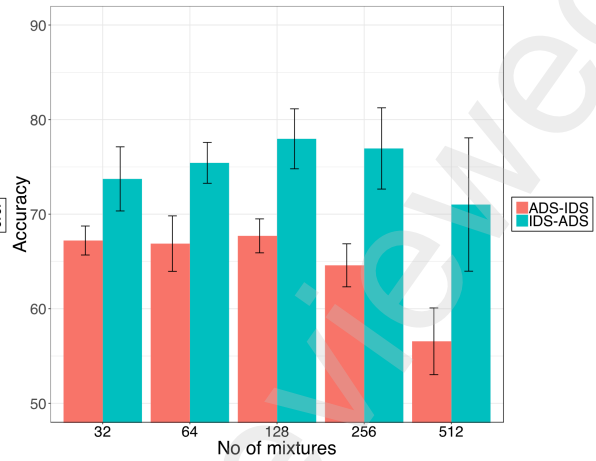
393

394

A)



B)



395

396

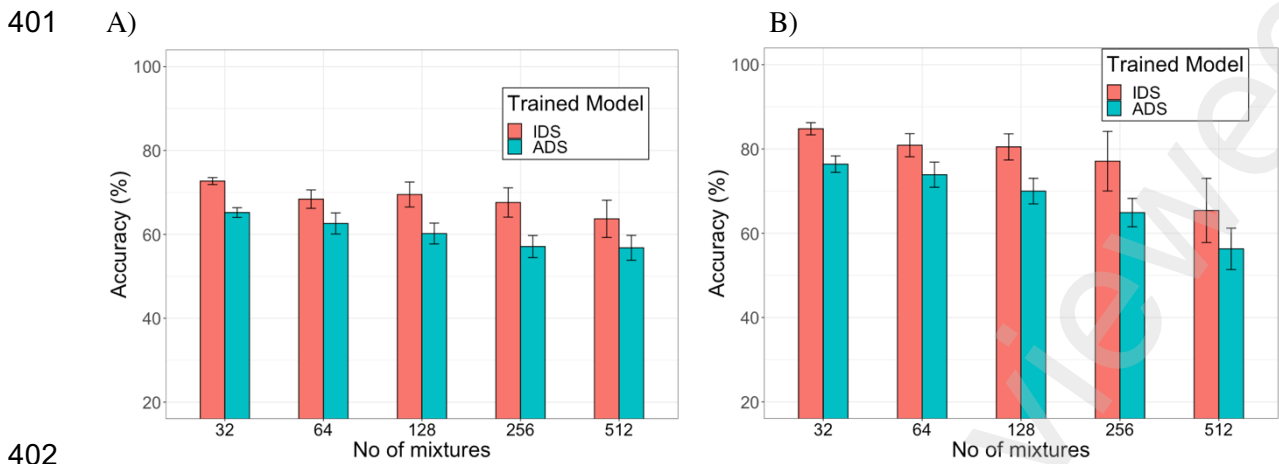
397

398

399

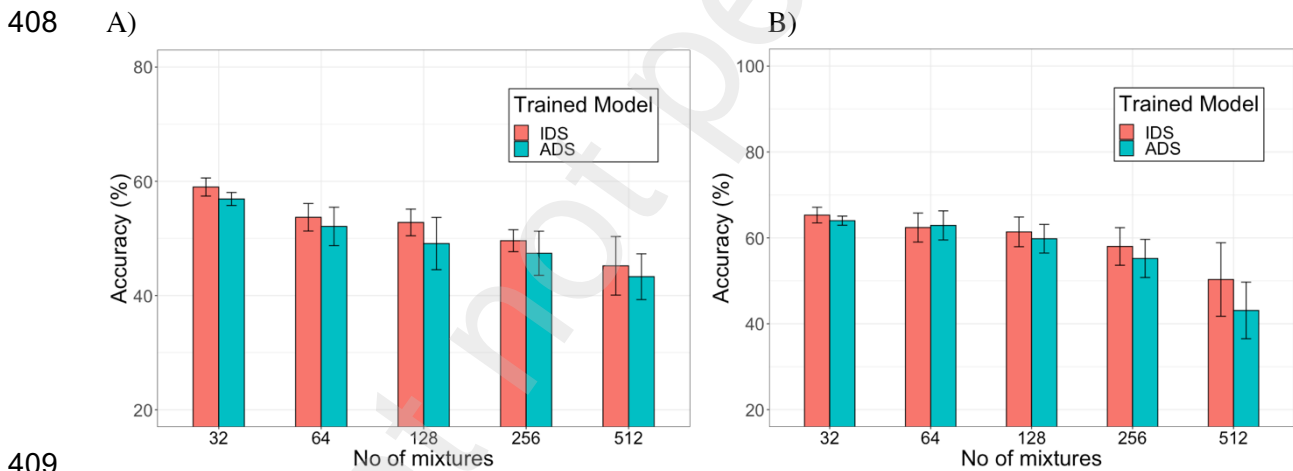
400

Figure S4. Cross-register voice recognition accuracy on Swiss-German dataset using a) MFCCs, and b) MFCCs with delta. The number of mixtures in the GMM classifier varied from 32 to 512, in the power of 2. The X-axis shows the number of mixtures in the GMM classifier, and the Y-axis shows the recognition accuracy in %, the vertical bar in the graph shows the standard-deviation of 5-fold cross validation.



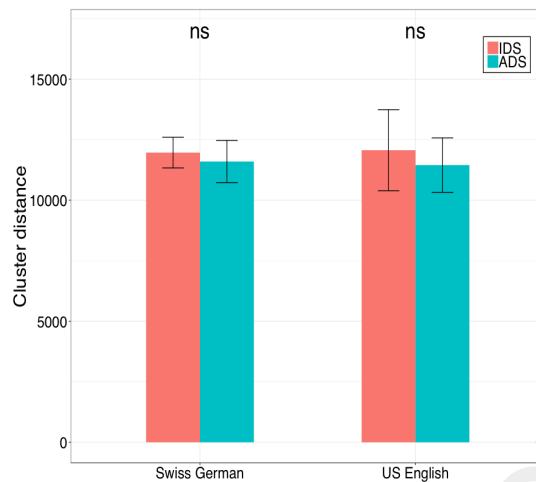
402

403 **Figure S5.** Speaker recognition accuracy on spontaneous speech (picture naming task in Swiss-German
 404 dataset) using a) MFCCs, and b) MFCCs with delta. The number of mixtures in the GMM classifier
 405 varied from 32 to 512, in the power of 2. The x -axis shows the number of mixtures in the GMM
 406 classifier, and the y -axis shows the recognition accuracy in %. The vertical bar in the graph shows the
 407 standard deviation of 5-fold cross validation.



409

410 **Figure S6.** Voice recognition accuracy on spontaneous speech (free speech in Swiss-German dataset)
 411 using a) MFCCs, and b) MFCCs with delta. The number of mixtures in the GMM classifier varied from
 412 32 to 512, in the power of 2. The x -axis shows the number of mixtures in the GMM classifier, and the
 413 y -axis shows the recognition accuracy in %. The vertical bar in the graph shows the standard-deviation
 414 of 5-fold cross validation.

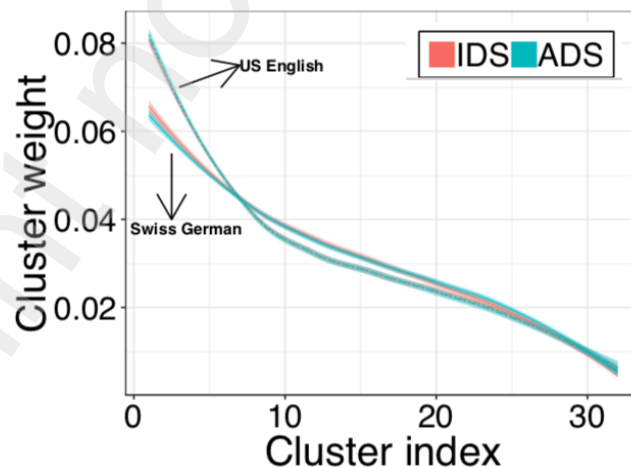


415

416 **Figure S7.** Sum of distance between centroids of the 32 IDS and ADS registers of Swiss German and
 417 US English dataset.

418

419 We also obtained the sum of the distances between cluster centroids in the IDS and the ADS feature
 420 space. Fig. S7 shows that the overall distance between IDS sub-clusters was higher compared to ADS
 421 sub-clusters (nonsignificant interaction between register and language: $F[1,35] = 0.3612$; $p = 0.5517$;
 422 significant main effect of registers: $F[1,35] = 5.6817$, $p < 0.0227$, but not of language: $F[1,35] = 0.0030$,
 423 $p < 0.9566$), meaning that IDS clusters are on average more widely dispersed compared to ADS clusters.



424

425 **Figure S8.** Cluster weight distribution of 32 clusters in IDS and ADS registers of the Swiss German
 426 and US English datasets.

427 **REFERENCES**

- 428 1. Kisilevsky, B.S., Hains, S.M.J., Lee, K., Xie, X., Huang, H., Ye, H.H., Zhang, K., and
 429 Wang, Z. (2003). Effects of experience on fetal voice recognition. *Psychological*
 430 *Science* 14, 220–224.
- 431 2. Beauchemin, M., González-Frankenberger, B., Tremblay, J., Vannasing, P., Martínez-
 432 Montes, E., Belin, P., Béland, R., Francoeur, D., Carceller, A.M., Wallois, F., et al.
 433 (2011). Mother and stranger: An electrophysiological study of voice processing in
 434 newborns. *Cerebral Cortex* 21, 1705–1711.
- 435 3. Johnson, E.K., Westrek, E., Nazzi, T., and Cutler, A. (2011). Infant ability to tell voices
 436 apart rests on language experience. *Developmental Science* 14, 1002–1011.
- 437 4. Piazza, E.A., Iordan, M.C., and Lew-Williams, C. (2017). Mothers Consistently Alter
 438 Their Unique Vocal Fingerprints When Communicating with Infants. *Current Biology*
 439 27, 3162-3167.e3.
- 440 5. Webb, A.R., Heller, H.T., Benson, C.B., and Lahav, A. (2015). Mother’s voice and
 441 heartbeat sounds elicit auditory plasticity in the human brain before full gestation.
 442 *Proc Natl Acad Sci U S A* 112, 3152–3157.
- 443 6. Blasi, A., Mercure, E., Lloyd-Fox, S., Thomson, A., Brammer, M., Sauter, D., Deeley,
 444 Q., Barker, G.J., Renvall, V., Deoni, S., et al. (2011). Early specialization for voice and
 445 emotion processing in the infant brain. *Current Biology* 21, 1220–1224.
- 446 7. Moser, C.J., Lee-Rubin, H., Bainbridge, C.M., Atwood, S., Simson, J., Knox, D.,
 447 Glowacki, L., Galbarczyk, A., Jasienska, G., Ross, C.T., et al. (2020). Acoustic
 448 regularities in infant-directed vocalizations across cultures. *bioRxiv*,
 449 2020.04.09.032995.
- 450 8. Broesch, T.L., and Bryant, G.A. (2015). Prosody in Infant-Directed Speech Is Similar
 451 Across Western and Traditional Cultures. *Journal of Cognition and Development* 16,
 452 31–43.
- 453 9. Kendrick, K.M. (2006). Introduction. The neurobiology of social recognition, attraction
 454 and bonding. *Philos Trans R Soc Lond B Biol Sci* 361, 2057–2059.
- 455 10. Kuhl, P.K., Andruski, J.E., Chistovich, I.A., Chistovich, L.A., Kozhevnikova, E. v.,
 456 Ryskina, V.L., Stolyarova, E.I., Sundberg, U., and Lacerda, F. (1997). Cross-language
 457 analysis of phonetic units in language addressed to infants. *Science* (1979) 277, 684–
 458 686.
- 459 11. Latinus, M., McAleer, P., Bestelmeyer, P.E.G., and Belin, P. (2013). Norm-based
 460 coding of voice identity in human auditory cortex. *Current Biology* 23, 1075–1080.
- 461 12. Lavan, N., Burston, L.F.K., and Garrido, L. (2019). How many voices did you hear?
 462 Natural variability disrupts identity perception from unfamiliar voices. *British Journal of*
 463 *Psychology* 110, 576–593.
- 464 13. Lee, Y., Keating, P., and Kreiman, J. (2019). Acoustic voice variation within and
 465 between speakers. *J Acoust Soc Am* 146, 1568–1579.
- 466 14. Hansen, J.H.L.L., and Hasan, T. (2015). Speaker recognition by machines and
 467 humans: A tutorial review. *IEEE Signal Processing Magazine* 32, 74–99.
- 468 15. Smith, H.M.J., Baguley, T.S., Robson, J., Dunn, A.K., and Stacey, P.C. (2019).
 469 Forensic voice discrimination by lay listeners: The effect of speech type and
 470 background noise on performance. *Applied Cognitive Psychology* 33, 272–287.

- 471 16. Ramírez López, A., Saeidi, R., Juvela, L., and Alku, P. (2017). Normal-to-shouted
472 speech spectral mapping for speaker recognition under vocal effort mismatch. In
473 ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing
474 - Proceedings, pp. 4940–4944.
- 475 17. Jokinen, E., Saeidi, R., Kinnunen, T., and Alku, P. (2019). Vocal effort compensation
476 for MFCC feature extraction in a shouted versus normal speaker recognition task.
477 *Computer Speech and Language* 53, 1–11.
- 478 18. Shriberg, E., Graciarena, M., Bratt, H., Kathol, A., Kajarekar, S., Jameel, H., Richey,
479 C., and Goodman, F. (2008). Effects of vocal effort and speaking style on text-
480 independent speaker verification. *Proceedings of the Annual Conference of the*
481 *International Speech Communication Association, INTERSPEECH*, 609–612.
- 482 19. Shriberg, E., Graciarena, M., Bratt, H., Kathol, A., International, S.R.I., and Park, M.
483 (2008). Effects of Vocal Effort and Speaking Style on Text-Independent Speaker
484 Verification. In *INTER_SPEECH 2008, 9th Annual Conference of the International*
485 *Speech Communication Association*, pp. 609–612.
- 486 20. Kuhl, P.K. (2004). Early language acquisition: cracking the speech code. *Nature*
487 *Reviews Neuroscience* 2004 5:11 5, 831–843.
- 488 21. Kitamura, C., and Burnham, D. (2003). Pitch and Communicative Intent in Mother's
489 Speech: Adjustments for Age and Sex in the First Year. *Infancy* 4, 85–110.
- 490 22. Singh, L., Morgan, J.L., and Best, C.T. (2002). Infants' Listening Preferences: Baby
491 Talk or Happy Talk? *Infancy* 3, 365–394.
- 492 23. Golinkoff, R.M., Can, D.D., Soderstrom, M., and Hirsh-Pasek, K. (2015). (Baby)Talk to
493 Me: The Social Context of Infant-Directed Speech and Its Effects on Early Language
494 Acquisition. *Current Directions in Psychological Science* 24, 339–344.
- 495 24. Gauthier, B., and Shi, R. (2011). A connectionist study on the role of pitch in infant-
496 directed speech. *J Acoust Soc Am* 130, EL380–EL386.
- 497 25. Kalashnikova, M., Carignan, C., and Burnham, D. (2017). The origins of babytalk :
498 smiling , teaching or social convergence ? Subject Category : Subject Areas : Author
499 for correspondence :
- 500 26. Burnham, D., Kitamura, C., and Vollmer-Conna, U. (2002). What's New, Pussycat?
501 On Talking to Babies and Animals. *Science* (1979) 296, 1435 LP – 1435.
- 502 27. Miyazawa, K., Shinya, T., Martin, A., Kikuchi, H., and Mazuka, R. (2017). Vowels in
503 infant-directed speech: More breathy and more variable, but not clearer. *Cognition*
504 166, 84–93.
- 505 28. van der Feest, S.V.H., Blanco, C.P., and Smiljanic, R. (2019). Influence of speaking
506 style adaptations and semantic context on the time course of word recognition in quiet
507 and in noise. *Journal of Phonetics* 73, 158–177.
- 508 29. Cooper, R.P., and Aslin, R.N. (1989). The language environment of the young infant:
509 Implications for early perceptual development. *Canadian Journal of*
510 *Psychology/Revue canadienne de psychologie* 43, 247–265.
- 511 30. Fernald, A., Taeschner, T., Dunn, J., Papousek, M., de Boysson-Bardies, B., and
512 Fukui, I. (1989). A cross-language study of prosodic modifications in mothers' and
513 fathers' speech to preverbal infants. *Journal of Child Language* 16, 477–501.
- 514 31. Fernald, A., and Simon, T. (1984). Expanded intonation contours in mothers' speech
515 to newborns. *Developmental Psychology* 20, 104–113.

- 516 32. Lee, C.S., Kitamura, C., Burnham, D., and McAngus Todd, N.P. (2014). On the
517 rhythm of infant- versus adult-directed speech in Australian English. *J Acoust Soc Am*
518 *136*, 357–365.
- 519 33. Payne, E.M., Post, B., Astruc, L., Prieto, P., and Vanrell, M. (2015). Rhythmic
520 modification in child directed speech. In *Supplementi alla biblioteca di linguistica.*, M.
521 Russo, ed. (Aracne), pp. 147–183.
- 522 34. Leong, V., Kalashnikova, M., Burnham, D., and Goswami, U. (2017). The Temporal
523 Modulation Structure of Infant-Directed Speech. *Open Mind: Discoveries in Cognitive*
524 *Science*, 1–13.
- 525 35. Kramer, R.S.S., Young, A.W., and Burton, A.M. (2018). Understanding face
526 familiarity. *Cognition* *172*, 46–58.
- 527 36. Kramer, R.S.S., Jenkins, R., Young, A.W., and Burton, A.M. (2017). Natural variability
528 is essential to learning new faces. *Visual Cognition* *25*, 470–476.
- 529 37. Burton, A.M., Kramer, R.S.S., Ritchie, K.L., and Jenkins, R. (2016). Identity From
530 Variation: Representations of Faces Derived From Multiple Instances. *Cognitive*
531 *Science* *40*, 202–223.
- 532 38. Lavan, N., Knight, S., Hazan, V., and McGettigan, C. (2019). The effects of high
533 variability training on voice identity learning. *Cognition* *193*, 104026.
- 534 39. Karlsson, I. (1999). Within-speaker variability in the VeriVox database. *Fonetik '99:*
535 *Proceedings from the Twelfth Swedish Phonetics Conference* *81*, 93–96.
- 536 40. Karlsson, I., Banziger, T., Dankovicová, J., Johnstone, T., Lindberg, J., Melin, H.,
537 Nolan, F., and Scherer, K. (2000). Speaker verification with elicited speaking styles in
538 the VeriVox project. *Speech Communication* *31*, 121–129.
- 539 41. Davis, S.B., and Mermelstein, P. (1980). Comparison of Parametric Representations
540 for. *Transactions on Acoustics, Speech, and Signal Processing* *28*, 357–366.
- 541 42. Friendly, M., Monette, G., and Fox, J. (2013). Elliptical insights: Understanding
542 statistical methods through elliptical geometry. *Statistical Science* *28*, 1–39.
- 543 43. Wieland, E.A., Burnham, E.B., Kondaurova, M., Bergeson, T.R., and Dilley, L.C.
544 (2015). Vowel space characteristics of speech directed to children with and without
545 hearing loss. *Journal of Speech, Language, and Hearing Research* *58*, 254–267.
- 546 44. Dilley, L., Lehet, M., Wieland, E.A., Arjmandi, M.K., Kondaurova, M., Wang, Y., Reed,
547 J., Svirsky, M., Houston, D., and Bergeson, T. (2020). Individual Differences in
548 Mothers' Spontaneous Infant-Directed Speech Predict Language Attainment in
549 Children With Cochlear Implants. *Journal of Speech, Language, and Hearing*
550 *Research* *63*, 2453–2467.
- 551 45. Cristia, A. (2022). A systematic review suggests marked differences in the prevalence
552 of infant-directed vocalization across groups of populations. *Developmental Science*.
- 553 46. DeCasper, A.J., and Fifer, W.P. (1980). Of Human Bonding - Newborns Prefer Their
554 Mothers Voices. *Science* (1979) *208*, 1174–1176.
- 555 47. Fecher, N., and Johnson, E.K. (2019). By 4.5 Months, Linguistic Experience Already
556 Affects Infants' Talker Processing Abilities. *Child Development* *90*, 1535–1543.
- 557 48. Flom, R., and Bahrick, L.E. (2007). The development of infant discrimination of affect
558 in multimodal and unimodal stimulation: The role of intersensory redundancy.
559 *Developmental Psychology* *43*, 238–252.

- 560 49. Lavan, N., Knight, S., Hazan, V., and McGettigan, C. (2016). No clear advantage for
561 high variability training during voice identity learning.
- 562 50. Lavan, N., Burton, A.M., Scott, S.K., and McGettigan, C. (2019). Flexible voices:
563 Identity perception from variable vocal signals. *Psychonomic Bulletin and Review* 26,
564 90–102.
- 565 51. Vosoughi, S., and Roy, D. (2012). A Longitudinal Study of Prosodic Exaggeration in
566 Child-directed Speech. *Proceedings of the 6th International Conference on Speech*
567 *Prosody (SP2012)*, 194–197.
- 568 52. Liu, H.-M., Tsao, F.-M., and Kuhl, P.K. (2009). Age-related changes in acoustic
569 modifications of Mandarin maternal speech to preverbal infants and five-year-old
570 children: a longitudinal study. *J. Child Lang.* 36, 909–922.
- 571 53. Ratner, N.B. (1984). Patterns of vowel modification in mother–child speech. *Journal of*
572 *Child Language* 11, 557–578.
- 573 54. Ko, E.-S. (2012). Nonlinear development of speaking rate in child-directed speech.
574 *Lingua* 122, 841–857.
- 575 55. Friedrichs, D., Maurer, D., Rosen, S., and Dellwo, V. (2017). Vowel recognition at
576 fundamental frequencies up to 1 kHz reveals point vowels as acoustic landmarks. *J*
577 *Acoust Soc Am* 142, 1025–1033.
- 578 56. Kathiresan, T., Maurer, D., and Dellwo, V. (2019). Highly spectrally undersampled
579 vowels can be classified by machines without supervision. *J Acoust Soc Am* 146,
580 EL1–EL7.
- 581 57. Ludusan, B., Mazuka, R., and Dupoux, E. (2021). Does Infant-Directed Speech Help
582 Phonetic Learning? A Machine Learning Investigation. *Cognitive Science* 45, 12946.
- 583 58. McMurray, B., Kovack-Lesh, K.A., Goodwin, D., and McEchron, W. (2013). Infant
584 directed speech and the development of speech perception: Enhancing development
585 or an unintended consequence? *Cognition* 129, 362–378.
- 586 59. Guevara-Rukoz, A., Cristia, A., Ludusan, B., Thiollière, R., Martin, A., Mazuka, R.,
587 and Dupoux, E. (2018). Are Words Easier to Learn From Infant- Than Adult-Directed
588 Speech? A Quantitative Corpus-Based Investigation. *Cognitive Science* 42, 1586–
589 1617.
- 590 60. Fernandez, A.A., and Knörnschild, M. (2020). Pup Directed Vocalizations of Adult
591 Females and Males in a Vocal Learning Bat. *Frontiers in Ecology and Evolution* 8.
- 592 61. Lavan, N., Knight, S., Hazan, V., and McGettigan, C. (2016). No clear advantage for
593 high variability training during voice identity learning.
- 594 62. Balcombe, J.P., and McCracken, G.F. (1992). Vocal recognition in Mexican free-tailed
595 bats: Do pups recognize mothers? *Animal Behaviour* 43, 79–87.
- 596 63. Jouventin, P., and Aubin, T. (2002). Acoustic systems are adapted to breeding
597 ecologies: individual recognition in nesting penguins. *Animal Behaviour* 64, 747–757.
- 598 64. Insley, S.J. (2001). Mother-Offspring vocal recognition in northern fur seals is mutual
599 but asymmetrical. *Anim Behav* 61, 129–137.
- 600 65. de Fanis, E., and Jones, G. (1995). Post-natal growth, mother-infant interactions and
601 development of vocalizations in the vespertilionid bat *Plecotus auritus*. *Journal of*
602 *Zoology* 235, 85–97.

- 603 66. Torriani, M.V.G., Vannoni, E., and McElligott, A.G. (2006). Mother-Young Recognition
604 in an Ungulate Hider Species: A Unidirectional Process. *The American Naturalist* 168,
605 412–420.
- 606 67. Creel, S.C., and Bregman, M.R. (2011). How Talker Identity Relates to Language
607 Processing. *Linguistics and Language Compass* 5, 190–204.
- 608 68. Belin, P. (2006). Voice processing in human and non-human primates. *Philosophical
609 Transactions of the Royal Society B: Biological Sciences* 361, 2091–2107.
- 610 69. Lee, G.Y., and Kisilevsky, B.S. (2014). Fetuses respond to father’s voice but prefer
611 mother’s voice after birth. *Developmental Psychobiology* 56, 1–11.
- 612 70. Ward, C.D., and Cooper, R.P. (1999). A lack of evidence in 4-month-old human
613 infants for paternal voice preference. *Developmental Psychobiology* 35, 49–59.
- 614 71. Zarate, J.M., Tian, X., Woods, K.J.P., and Poeppel, D. (2015). Multiple levels of
615 linguistic and paralinguistic features contribute to voice recognition. *Scientific Reports*
616 5, 11475.
617