

Different in different ways: A network-analysis approach to voice and prosody in Autism Spectrum Disorder

Journal:	<i>Language Learning and Development</i>
Manuscript ID	HLLD-2022-0074.R1
Manuscript Type:	Research Article
Date Submitted by the Author:	03-Jan-2023
Complete List of Authors:	Weed, Ethan; Aarhus University, Fusaroli, Riccardo; Aarhus University Simmons, Elizabeth; Sacred Heart University, Communication Disorders Eigsti, Inge-Marie; University of Connecticut, Psychological Sciences;
Keywords:	Speech production, Computational, Pragmatics

SCHOLARONE™
Manuscripts

1 Introduction

2 Difficulties with speech and language are a defining feature of autism spectrum disorder (ASD)
3 (Association, 2013), and poor communication skills exacerbate the risk being bullied five-fold
4 (Cappadocia, Weiss, & Pepler, 2012). Atypical prosody and vocal presentation impact social
5 integration and employment opportunities. Indeed, some older research suggests that unusual
6 prosody is among the first features that elicit an impression of oddness from others (Mesibov,
7 1992; Shriberg & Widder, 1990; Van Bourgondien & Woods, 1992). Even individuals with ASD
8 whose expressive and receptive language scores are in the typical range have difficulty with
9 prosody (Shriberg, Paul, Black, & Santen, 2011), and children with ASD are rated as more
10 socially awkward on the basis of audio speech samples alone (Grossman, 2015; Redford,
11 Kapatsinski, & Cornell-Fabiano, 2018; Sasson et al., 2017). One mechanism suggested for this is
12 that speakers with ASD are more consistent in their phonetic production than NT speakers, as
13 reported in a study of adults with and without ASD (Kissine, Geelhand, Philippart, Harmegnies
14 & Deliens, 2021). Atypical vocal quality is one of the signals of the “frank” autism presentation,
15 apparent to expert clinicians within just a few seconds (Marchena & Miller, 2017). However,
16 despite their significant clinical impact, the specific acoustic features underlying speech
17 differences in ASD are as yet poorly understood. This is a clinically important gap in knowledge,
18 as a better understanding of what acoustic cues contribute to atypical speech quality would
19 enable us to develop better-informed targets for intervention. The goal of this manuscript is to
20 combine information from acoustic analyses with ratings from naïve and clinical raters, in order
21 to provide some insights into what acoustic qualities make a voice distinctively autistic.

22 Although clinicians since Asperger have described the voices of people with autism as
23 having unusual prosodic (e.g., “robotic”, “flat”, “monotone”), and vocal (e.g., “harsh”, “nasal”

1 Different in different ways

2
3 24 and “hoarse”) qualities, there is little consensus on exactly *which* acoustic features of ASD
4
5 25 speech differ from typical speech (McCann & Peppé, 2003). This is true even though
6
7 26 productions may be reliably reported as “odd” by naïve listeners, as was the case for a study of
8
9 27 utterances produced by 12 youths with Asperger syndrome (Filipe, Castro & Vicente, 1981). A
10
11 28 recent systematic review and meta-analysis (Fusaroli, Lambrechts, Bang, Bowler, & Gaigg,
12
13 29 2017) highlights several potential acoustic contributors to group differences. The meta-analytic
14
15 30 findings point to higher pitch, and increased pitch variability, number of pauses, and pause
16
17 31 duration, although replication attempts show these patterns might not be generalizable across
18
19 32 languages and samples (Parola, 2022; Fusaroli et al, 2021, Rybner et al, 2022). While some
20
21 33 studies have found that acoustic differences map onto expert clinician ratings of autism-specific
22
23 34 symptoms (McCann et al, 2007; Study 1, Diehl et al. 2009), these findings have not always
24
25 35 replicated, even within the same lab (Study 2, *ibid*; Fusaroli et al 2021).

26
27
28
29
30
31 36 Although the speech of autistic people is usually described as having an atypical *prosody*,
32
33 37 the adjectives used to describe the speech often allude to aspects that go beyond prosody. In one
34
35 38 of the earliest descriptions of the voices of children with autism, for instance, Asperger (1991)
36
37 39 described them as “shrill,” “soft,” and “nasal,” qualities that belong more properly to the realm
38
39 40 of *voice quality*. Voice quality has been shown to play an important role in the impressions we
40
41 41 form of speakers. Listeners routinely extract numerous impressions of speaker characteristics,
42
43 42 based on voice quality alone: gender identity, age, mood, socio-economic status, sexual
44
45 43 orientation, social relationships, and even neuropsychiatric conditions (Bryant et al., 2016;
46
47 44 Cummins et al., 2015; Parola, Simonsen, Bliksted, & Fusaroli, 2020; Low, Bentley & Ghosh,
48
49 45 2020, Redford et al., 2018; Weed & Fusaroli, 2020). However, a one-to-one mapping between
50
51 46 acoustic features and perceptual impressions has proven elusive; this likely reflects broader
52
53
54
55
56
57
58
59
60

Different in different ways

1
2
3 47 challenges in mapping acoustic features onto other dimensions of speech, such as phonology,
4
5 48 described as the “many to many” mapping problem (Liberman, Cooper, Shankweiler, &
6
7 49 Studdert-Kennedy, 1967).

10 50 Both expert and non-expert raters are able to reliably distinguish between typical and
11
12 51 atypical prosody in short samples of speech, and are thereby often able to predict which speakers
13
14 52 have an ASD diagnosis (Nadig and Shaw, 2012; Redford et al., 2018). However, consistent with
15
16 53 the complexity noted above, listeners, including trained clinicians, are often not explicitly aware
17
18 54 of which acoustic characteristics seem distinctive in ASD, and thus of which acoustic features or
19
20 55 combination of features they are responding to. Nadig and Shaw (2012) found that speech-
21
22 56 language pathology students could not reliably identify atypical variation in pitch in speakers
23
24 57 with ASD, although the same raters were able to distinguish between broadly typical versus
25
26 58 atypical prosody in the same speakers. In a related study, Redford et al. (2018) found that naïve
27
28 59 raters were able to distinguish between the voices of autistic and neurotypical (NT) children
29
30 60 when asked to score on atypicality, but not when listening to low-pass or high-pass filtered
31
32 61 versions of the same speech, suggesting that “[...] sound pattern differences are subtle enough to
33
34 62 be obscured in degraded speech,” (Redford et al., 2018, p. 289). Dahlgren et al. (2018) found that
35
36 63 a group of three speech language pathologists who specialized in voice were only able to
37
38 64 correctly identify three out of eleven children with ASD on the basis of voice alone, and
39
40 65 concluded that voice and speech characteristics are not unproblematic predictors of diagnosis.

42 66 To further complicate matters, there is no consensus as to what is distinctive about the
43
44 67 voice and prosody of people with ASD. Peppé et al. (2007) point to the wide range of sometimes
45
46 68 contradictory adjectives that have been used to describe the speech of people with ASD. They
47
48 69 note that autistic speech has been simultaneously described as “dull”, “wooden”, “singsong”,
49
50
51
52
53
54
55
56
57
58
59

1 Different in different ways

2
3 70 “robotic”, “stilted”, “over precise” and “bizarre”, but also both “monotonous” and “exaggerated”
4
5 71 as well as “slow” and “fast,” raising the question of whether it possible to draw any consistent
6
7 72 conclusions about expressive prosody in ASD at all.

8
9
10 73 There are several possible explanations for these inconsistent findings. First, it may be
11
12 74 the case that there simply are no consistent patterns to detect. Another possible explanation is
13
14 75 that there are different *profiles* of atypical prosody and voice quality in speakers with ASD. In
15
16 76 other words, the atypical prosody and voice quality of speakers with ASD may fall into a small
17
18 77 (or large) set of distinct patterns, with several different characteristic types of atypicality. This
19
20 78 second explanation is consistent with both the reliable differences in perceptual impressions, and
21
22 79 with the lack of a single clear constellation of acoustic features characteristic of the speech of
23
24 80 people with ASD. While listeners may pick up on these atypicalities, the presence of multiple
25
26 81 different constellations of acoustic features that are all “atypical” may make it difficult to reliably
27
28 82 predict diagnosis on the basis of acoustic features alone. As McCann and Peppé write, “If
29
30 83 findings were consistent, small-scale studies would offer pointers, but as it is these do not inspire
31
32 84 confidence,” (2003, p. 347). This opinion is reiterated a decade later in Fusaroli et al. (2017).

33
34
35 85 Prosody includes changes in pitch, loudness, and rhythm of sounds (Cutler et al., 1997;
36
37 86 Peppé, 2009; Shattuck-Hufnagel & Turk, 1996). *Pitch* describes the perceptual experience of the
38
39 87 fundamental frequency (F0) of a sound (e.g., a high versus low voice), and *variation in pitch* is
40
41 88 most often described by the standard deviation, range, or interquartile range of F0. *Volume* or
42
43 89 loudness refers to the perceptual experience of intensity, or the amount of energy in the sound
44
45 90 waves. Finally, the concept of *rhythm* describes acoustic features related to the duration of
46
47 91 speech and pauses: among them the duration of the utterance, the number and length of pauses,
48
49 92 as well as speech rate and articulation rate. While a high speech rate relates to the amount of
50
51
52
53
54
55
56
57
58
59
60

Different in different ways

93 information (syllables) per unit of time, a high articulation rate indicates more clipped, less
94 drawn-out syllables. An individual could thus have, for example, a low speech rate but a high
95 articulation rate; such a speaker would produce very quick syllables, but have long pauses in
96 their speech. All these measures have been indicated as potential markers of ASD (Fusaroli et al.,
97 2017).

98 Although descriptions of voice quality appear in the literature (e.g., “harsh”, “nasal”, or
99 “hoarse”), voice quality measures are remarkably absent from the study of vocal atypicality in
100 ASD. **Two studies** (Bone et al., 2014; **Kissine & Geelhand, 2019**) have investigated *jitter*, that
101 is, cycle-to-cycle changes in the fundamental period, which is associated with the perceptual
102 qualities of breathiness and hoarseness (Eskenazi, Childers, & Hicks, 1990; Wolfe, Fitch, &
103 Martin, 1997). Other candidate acoustic features are *shimmer* (**Kissine & Geelhand, 2019**), which
104 quantifies cycle-to-cycle fluctuations in the amplitude of the waveform and has been related to
105 both breathiness and hoarseness (Wolfe et al., 1997); *harmonics-to-noise ratio*, which quantifies
106 the relative amount of energy in harmonic portions of the spectrum with other, “noise” energy,
107 and has been related to hoarseness (Yumoto, Gould, & Baer, 1982); and *H1-H2*, the relative
108 amplitudes of the first two harmonics, which is also linked to the perception of breathiness
109 (Hillenbrand & Houde, 1994; Klatt & Klatt, 1990; Kreiman & Gerratt, 2010). Not only do all
110 these features measure voice qualities which have been linked to perceptual qualities noted in
111 ASD, they also tend to be intercorrelated, as they all derive from fluctuations in the glottal source
112 (Murphy, 2000).

113 We have discussed a handful of the most common features, but many more can be
114 extracted. In studies on dysarthric speech, Borrie et al. (2020) measured over 800 acoustic
115 features, and Al-Qatab & Mustafa (2021) measured over 5000 acoustic features. The ComParE

Different in different ways

1
2
3 116 baseline feature set consists of over 6000 features (Schuller et al., 2016). With so many
4
5 117 potentially important features, the task of feature selection is critical. Unfortunately, although
6
7 118 several different methods are commonly used, there is currently no single accepted method for
8
9 119 reducing the number of features to include in an analysis. Often, features are chosen with the aim
10
11 120 of finding the best combination of features for predicting category membership or some other
12
13 121 measure (e.g., predicting diagnosis, or symptom severity). Borrie et al (2020) used independent
14
15 122 component analysis (ICA) for feature reduction. ICA seeks to identify independent “sources” of
16
17 123 information in a multivariate dataset. Al-Qatab et al (2021) employed seven different algorithmic
18
19 124 feature reduction techniques: Interaction Capping, Conditional Information Feature Extraction,
20
21 125 Conditional Mutual Information Maximization, Double Input Symmetrical Relevance, Joint
22
23 126 Mutual Information, Conditional Redundancy, and Relief. Another common method is Principal
24
25 127 Component Analysis (PCA; e.g. Mittal & Sharma, 2021; Peng et al., 2007). While these
26
27 128 algorithmic methods can be a powerful way to reduce the number of features to a manageable
28
29 129 number, they have the disadvantage of often resulting in feature sets that can be difficult to
30
31 130 interpret in ways that are clinically useful. PCA, for example, reduces individual features to
32
33 131 linear combinations of features, which can be difficult to describe intuitively.

34
35
36
37
38
39
40 132 Given these challenges, the present study adopted an alternative approach. We chose a
41
42 133 small number of acoustic features that were among the most frequently mentioned in the
43
44 134 literature: a measure of pitch variation (standard deviation of fundamental frequency); two
45
46 135 measures of rhythm (speech rate; syllables per second) and articulation rate (syllables per second
47
48 136 after removing pauses, which were defined as absence of voice for more than 200 milliseconds);
49
50 137 and a measure of voice quality (jitter). Note that although variation in fundamental frequency
51
52 138 and jitter are both derived from the glottal source, they are independent: fundamental frequency
53
54
55
56
57
58
59
60

Different in different ways

1
2
3 139 (perceived as pitch) is the lowest of the many harmonics in the voice, while jitter represents the
4
5 140 difference in length between each cycle and the preceding cycle in the sonogram; put broadly,
6
7 141 pitch is estimated by looking for large-scale similarities in the acoustic signal, while jitter is
8
9 142 estimated by measuring small differences in the signal. As measures of rhythm, speech rate and
10
11 143 articulation rate are complementary, capturing two potential sources of the “robotic” quality that
12
13 144 is often ascribed to autistic speech. Furthermore, given that elicited sentences were scripted, and
14
15 145 the number of syllables was therefore held constant, combining speech rate and articulation rate
16
17 146 also provided an implicit measure of pause length. We chose not to include intensity, although it
18
19 147 is among the key elements of prosody, as our participants were at variable distances from the
20
21 148 microphone during recording, which influences intensity. Although algorithmic feature reduction
22
23 149 might well result in a set of features with a higher predictive power, our goal in choosing
24
25 150 common features that have been previously investigated in speakers with ASD, and which cover,
26
27 151 to the extent possible, qualitative reports on the speech of people with ASD, was to select a set of
28
29 152 easily-interpretable features, which would be relevant to clinicians.
30
31
32
33
34

35 153 Certainly, our choice of features presents some disadvantages. Standard deviation of F0 is
36
37 154 a very crude measure of pitch as a contributor to prosody and masks many important nuances of
38
39 155 pitch modulation. Jitter can be affected by factors such as vocal strain (Brockmann-Bauser et al.,
40
41 156 2014; Huang et al., 1995), although jitter may be less affected by room noise and microphone
42
43 157 quality than other measures such as shimmer and harmonics-to-noise ratio (Bottalico et al., 2020;
44
45 158 van der Woerd et al., 2020). In this exploratory study, however, the advantage of working with a
46
47 159 small number of commonly-used features outweighed the disadvantages.
48
49
50

51 160 Rather than trying to identify a single set of acoustic characteristics that describe the
52
53 161 “autistic voice,” the current study aims to identify acoustic profiles that characterize clusters of
54
55
56
57
58
59
60

1 Different in different ways
2

3 162 individuals whose voice and prosody are more similar to each other than to speakers in a
4
5 163 different cluster, and to test how these profiles align to diagnosis and to listener ratings. We
6
7 164 posed the following exploratory hypotheses:

8
9
10 165 H1. Modelling speakers as nodes of acoustic features in a network will allow us to identify
11
12 166 coherent clusters of speakers.

13
14 167 H2. Neurotypical and autistic speakers will tend to cluster differently.

15
16 168 H3. Clusters of speakers will be characterized by distinctive constellations of acoustic
17
18 169 features.

19
20 170 H4. Clusters will reflect the subjective ratings of naïve and clinical raters.

21
22
23
24 171 Network analysis techniques are well suited to this sort of analysis, because they allow
25
26 172 visualization of individuals in relation to each other. To address these hypotheses, we coded
27
28 173 speech samples, and used techniques from network analysis to build clusters of speakers, derived
29
30 174 from the acoustic features of their speech. We then inspected the distribution of autistic and
31
32 175 neurotypical speakers across the clusters, as well as the acoustic profiles that generated these
33
34 176 clusters. Finally, we tested whether naïve and expert raters tended to rate speakers from the
35
36 177 clusters as more or less atypical. We emphasize that our approach here is exploratory, and
37
38 178 acknowledge that our sample size is limited. To our knowledge, the methods we employ here for
39
40 179 identifying speaker clusters on the basis of acoustic measures have not previously been used for
41
42 180 this purpose, although similar graph-based methods are currently being explored in the
43
44 181 development of speaker recognition software (Chen et al., 2021; Wang et al., 2021).
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Different in different ways

182 **Methods**

183 *Participants*

184 Participants were 13 autistic and 13 neurotypical (NT) male adolescents, all native
 185 speakers of American English. Full-scale IQ (FSIQ), verbal IQ (VIQ), and non-verbal IQ
 186 (NVIQ) scores were estimated with the Stanford-Binet Abbreviated IQ (Roid, 2003). Groups did
 187 not differ in chronological age, full-scale IQ, or structural language abilities, measured using the
 188 Clinical Evaluation of Language Fundamentals (CELF; Wiig, Semel, & Secord, 2003). CELF
 189 and IQ scores were in the average range for all participants. ASD diagnoses were confirmed by
 190 trained graduate clinicians using the Autism Diagnostic Observation Schedule (ADOS) Module
 191 3 (Lord et al., 2000), selected sections of the Autism Diagnostic Interview, Revised (ADI- R;
 192 Lord, Rutter, & LeCouteur, 1994), and clinical judgment, based on the DSM-IV-TR criteria
 193 checklist (APA, 2000). Parents reported no co-morbid diagnoses, and completed the Social
 194 Responsiveness Scale (Constantino et al., 2003) to verify group membership of both ASD and
 195 NT participants, and as measure of symptom severity. Participant details are shown in Table 1.

196 *Table 1: Participant Characteristics*

	ASD	NT	Welch Two
	(<i>n</i> = 13)	(<i>n</i> = 13)	Sample t-test
Age (years)	14.24 (1.82)	14.30 (1.41)	$t(22.60) = -0.08, p = 9.36$
FSIQ Standard Score	102.76 (9.75)	104.15 (9.60)	$t(23.99) = -0.36, p = 0.71$

Different in different ways

VIQ Standard Score	9.69 (2.56)	11.53 (2.29)	$t(23.71) = -1.93, p = 0.06$
NVIQ Standard Score	11.30 (2.59)	9.84 (1.95)	$t(22.28) = 1.62, p = 0.11$
CELF Core Language Standard Score	109.41 (10.42)	114.84 (10.08)	$t(22.69) = -1.32, p = 0.19$
Social Responsiveness Scale	74.66 (11.52)	44.84 (7.95)	$t(19.37) = 7.47, p = 0.00004$
ADOS	9.90 (2.84)	N/A	N/A

197

198 *Speech elicitation*

199 The data were originally collected as part of a study investigating how adolescents with
 200 and without ASD used prosody to disambiguate sentences (Mayo, 2015). Although the original
 201 study included a training session, in which participants were provided with suggestions on how
 202 they could modify their prosody to better disambiguate sentences, the utterances used in the
 203 present study were from the pre-training, baseline condition, in which participants each produced
 204 eight sentences with similar sentence structures and high-frequency words. The sentences were
 205 written on printed cards and participants were asked to learn the sentence; once it was
 206 memorized, they were asked to speak the sentence aloud “in a normal voice” without reading
 207 from the card. Sentences were all instructions to interact with an animal, e.g., “Point at the lamb
 208 with the flower.” or “Tap the duck with the lollipop.” Speech was recorded with a Marantz
 209 Professional Model PMD660 audio recorder. A full description of the elicitation procedure, as
 210 well as additional information on recruitment, can be found in Mayo (2015). Although this

Different in different ways

211 method of elicitation lacks the ecological validity of, e.g., recordings of natural conversations,
212 the controlled nature of task was well-suited to the purpose of the current study. Because the
213 speech samples included all identical words, potential prosodic differences due to words with
214 differing numbers of syllables, or pitch fluctuations due to stress patterns or emotional valence
215 associated with different phrases were controlled for, providing a more appropriate framework
216 for our exploratory network analysis than recordings of spontaneous speech.

Feature Extraction

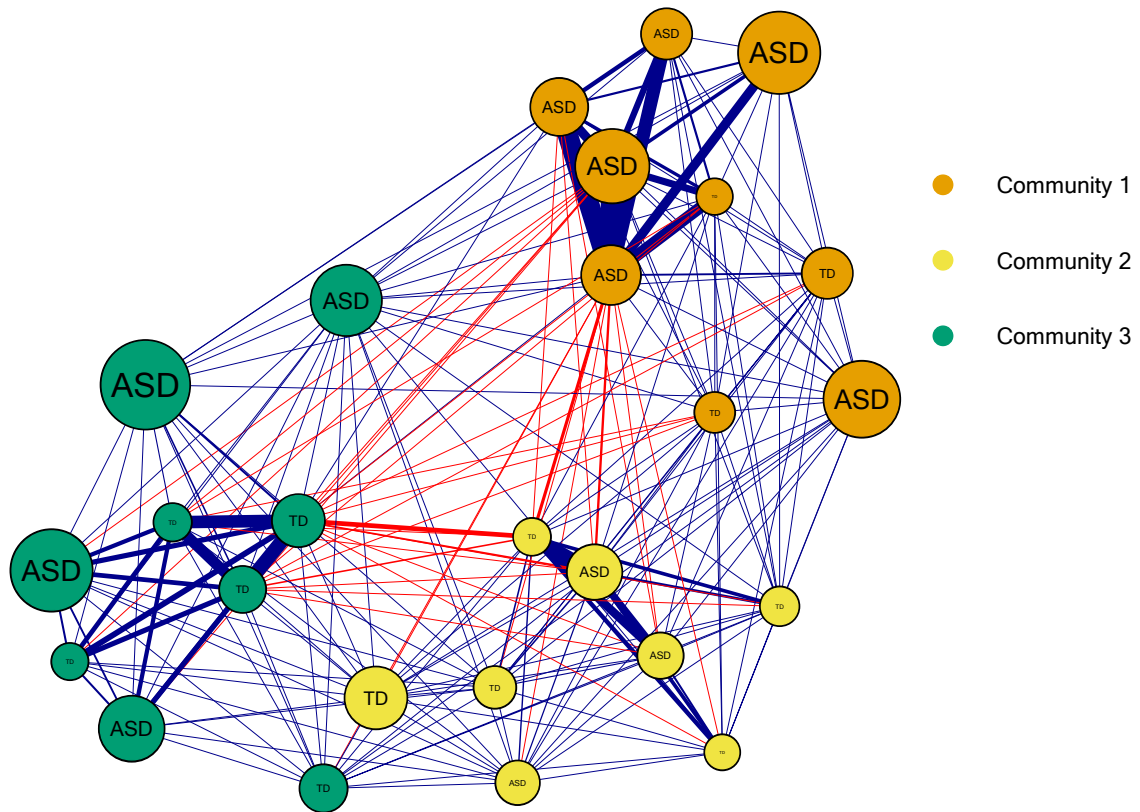
218 F0, jitter, speech rate, and articulation rate were extracted using custom Praat (Boersma
219 & Weenink, 2001) scripts (all scripts and extracted data are available on OSF at
220 https://osf.io/zw67n/?view_only=63441e15ec264184bb7b0b22c4140a22). For F0 extraction, we
221 set a floor value of 50 Hz, and a ceiling of 400 Hz. We inspected individual pitch tracks for
222 evidence of multimodality (period doubling). Fig S1 (available at the OSF repository) displays
223 histograms of the results of the pitch tracking for all participants. Although a few of the speakers
224 voices did show some evidence of multimodality (e.g., participants 9027, 9030, 9033), on closer
225 inspection at least some of these, such as participants 9027 and 9033, could be accounted for by
226 a period of very low, occasionally “creaky” phonation. For the sake of consistency, we applied
227 the 50-400 Hz range to all speakers. F0 was converted to a semitone scale, using the *hqmisc*
228 (Quene, 2022) package for R, using the default value of 50 Hz as the reference frequency.

Network Models and Community Detection

230 To identify groups of speakers with similar acoustic patterns, we first constructed a
231 network of speakers. We used the *qgraph* (Epskamp, Cramer, Waldorp, Schmittmann, &
232 Borsboom, 2012) package for R Statistical Software (v4.2.1; R Core Team, 2022) to create a
233 partial correlation matrix of speakers, with each speaker represented by the four selected prosody

Different in different ways

234 and voice features. We plotted this matrix as a network graph, in which each speaker is
 235 represented by a circle (node); the lines (edges) connecting the nodes represent the pairwise
 236 correlations (similarity) between each pair of speakers, when the effects of all other speakers
 237 have been removed; see Figure 1.



238

239 Figure 1: Network of acoustically-defined communities. Each node represents a single speaker.
 240 Blue connections between nodes represent positive correlations, red lines represent negative
 241 correlations. Line thickness indicates the absolute strength of the correlation.

242

243 We used the edges between nodes to identify clusters or “communities” of speakers with
 244 similar voice/prosody profiles. To do this, we used a spin glass community detection algorithm,

Different in different ways

245 as proposed in Reichardt & Bornholdt (2006) and Traag & Bruggeman (2008), and implemented
246 in the *iGraph* package (Csardi & Nepusz, 2006) for R. The spin glass approach attempts to find
247 communities of nodes by maximizing the number of *positive edges* (positive correlations)
248 between nodes within the community and minimizing the number of positive edges between
249 members of the community and members outside the community. At the same time, the
250 algorithm maximizes the number of *negative edges* (negative correlations) with nodes outside the
251 community, while minimizing the number of negative edges within the community. Although the
252 user sets an upper limit to the number of communities that the algorithm will detect, the optimal
253 number of communities may be lower than this upper limit (Csardi & Nepusz, 2006). We set the
254 upper limit of possible communities at ten. Because the algorithm is non-deterministic, it may
255 not always settle on the same number of communities each time it is run. For this reason, we
256 followed the procedure outlined in Djelantik et al. (2020) and seeded the random number
257 generator with a seed which would obtain the median number of communities identified in 1000
258 runs of the algorithm. Using this method, the algorithm always settled on a three-community
259 solution, and we therefore chose a seed which produced a 3-community solution, and used these
260 three communities for further analysis.

Listener judgments

262 Ten clinicians with ASD expertise and 15 undergraduates without ASD experience, all
263 naïve to the study hypotheses, listened to unaltered speech samples, and rated each sample as
264 “atypical or unusual” on a 1-3 scale (1 = typical; 2 = somewhat unusual; 3 = definitely atypical;
265 undergraduates) or as autistic or non-autistic (0 = NT-like, 1 = ASD-like) on a binary scale
266 (clinicians). The clinical and naïve raters’ scores were calculated as the mean of all the raters’
267 scores for speech samples from that individual, normalized to a 0 to 1 scale. Interrater reliability

Different in different ways

268 was calculated as Average Score Intraclass Correlation, using the *irr* package (Gamer et al.,
269 2012) for R. Estimates were based on a mean rating (naïve: $k = 17$, expert: $k = 5$) agreement, 2-
270 way random-effects model (Koo & Li, 2016).

271 **Results**

272 *Community Detection.* Community membership was very consistent. Out of 1000 runs,
273 there were five speakers who were occasionally placed in different communities; see Table 2.
274 The most inconsistently-grouped was an NT speaker, who was placed in Community 1 in 67.9%
275 of runs and in Community 2 in 32.1% of runs. Another NT speaker was placed in Community 3
276 in 70.3% runs and in Community 2 in 29.7% runs. The other three were placed in other
277 communities in less than 2% of runs. Each of these speakers were placed in their “preferred”
278 community in the final network.

279 Table 2: Community assignments across 1000 runs (inconsistently assigned participants only).
280 All participants not listed in the table were assigned to the same community across all runs.
281 Values indicate percentage of time the participant was assigned to the community.

ID	Dx	Com1	Com2	Com3
9004	NT	67.9	32.1	0
9009	ASD	1.2	0	98.8
9021	ASD	99.9	0.1	0
9025	NT	0	98.6	1.4
9029	NT	0	29.7	70.3

282

283

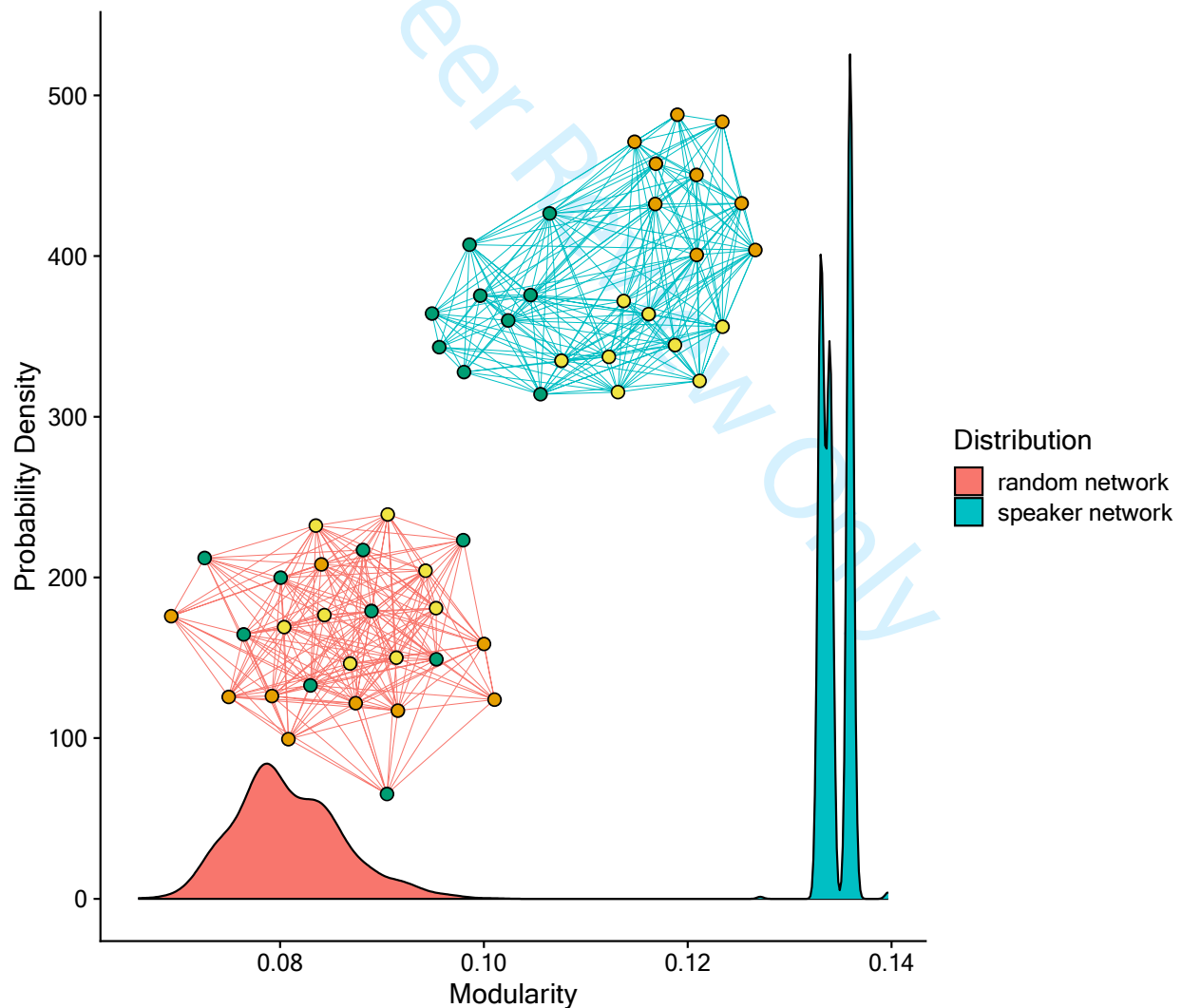
Different in different ways

284 *Network modularity*. Because the community-detection algorithm can identify
285 “communities” even in randomly-generated networks, it is important to assess the strength of the
286 communities detected. This can be done by calculating a “modularity” score, that is, an estimate
287 of the degree of “clumpiness” in the network, and then comparing this score with the distribution
288 of modularity scores from randomly-generated networks. To test whether the placement of
289 speakers into communities in our network was at greater than chance level, we first computed the
290 level of modularity of the network, where modularity is defined as the number of edges that fall
291 within a community divided by the number of edges that would be expected to fall within a
292 community in a randomized network, in which each node has the same number of connections as
293 the actual network, but in which these connections are randomized (Clauset et al, 2004). This
294 resulted in a modularity score of 0.1345. We then compared the modularity of our network with
295 the distribution of modularity scores from a population of random networks built with the same
296 constraints of numbers of nodes and number of edges incident upon each node (Csárdi et al,
297 2016; Maslov & Sneppen, 2002) by generating 1000 random networks with the same constraints
298 as our network of speakers, detecting communities for each of them using the same spin glass
299 algorithm, and calculating the modularity score for each random network. The resulting mean
300 modularity score was 0.0808, with a standard deviation of 0.0052. Thus, the modularity score of
301 our network lay approximately eleven standard deviations above the mean; see Figure 2. By
302 Chebyshev's inequality (Knuth, 1997 pg. 98), no more than 0.82 percent of values can lie 11
303 standard deviations away from the mean of the probability distribution of modularities calculated
304 on random networks, implying that it is highly unlikely that our network was drawn from the
305 same distribution as the random networks. To further quantify the modularity of the speaker
306 network, we compared the distribution of modularities from 1000 runs of the community

Different in different ways

307 detection algorithm on the actual data with the modularities from 1000 random networks with the
 308 same constraints as the speaker network with a Welch Two Sample t-test, $t(1315.8) = -337.42$, p
 309 $= 2.2 \cdot 10^{-16}$)

310 *Clustering of Diagnosis and Acoustic Features.* Community 1 (Fig. 1, orange group; $n =$
 311 9) consisted of six autistic speakers and three NT speakers, Community 2 (Fig. 1, yellow group;
 312 $n = 8$) consisted of five NT speakers and three autistic speakers, and Community 3 (Fig. 1, green
 313 group; $n = 9$) consisted of four autistic speakers and five NT speakers.



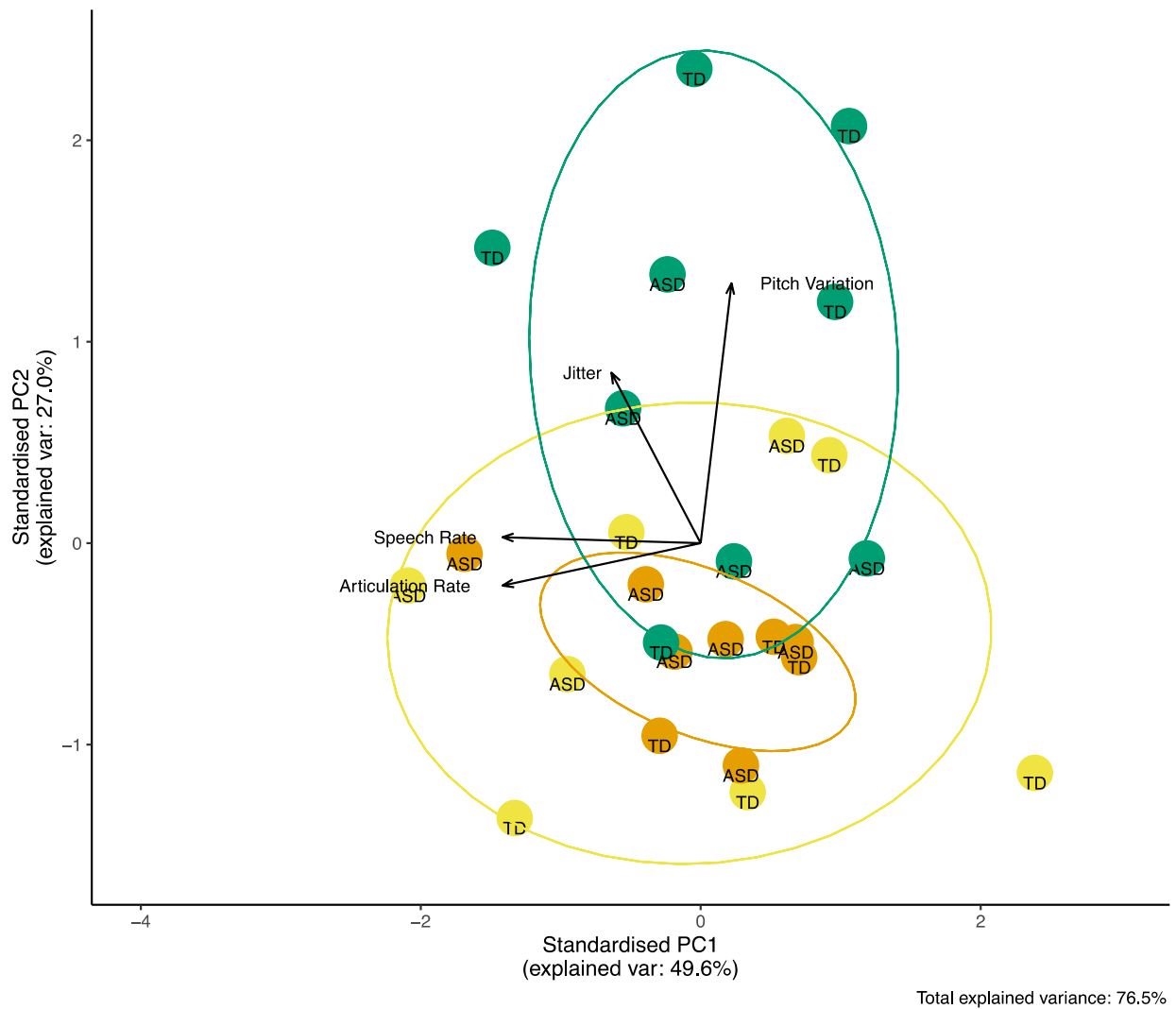
315

Different in different ways

1
2
3 316 Figure 2: Distributions of modularity scores for 1000 random networks built with the same
4 317 constraints as the actual network of speakers, and the distribution of modularity scores for 1000
5 318 runs of community detection on the speaker network. Also displayed are a sample random
6 319 network and the actual speaker network. The colors of the connections in the sample networks
7 320 indicate which distribution they are drawn from. The colors of the nodes indicate community
8 321 membership.

322

323



324

325 Figure 3: Speakers and communities visualized along the first two dimensions of a principal
326 components decomposition. Ellipses colors and point colors indicate communities detected by the
327 spin glass algorithm (orange = Community 1, yellow = Community 2, green = Community 3).
328 Arrows indicate the influence of each of the features on the components.

1 Different in different ways
2
3

4 329

5 330 To better visualize the patterns of acoustic variables that characterized the three
6
7 331 communities, we plotted these variables in a principal components analysis, using the *pca*
8
9 332 function from the AMR package for R (Berends et al., 2021). The first component accounted for
10
11 333 49.6 percent of the variance, and was heavily influenced by the two rhythm features, while the
12
13 334 second component accounted for 27 percent of the variance, and was dominated by pitch
14
15 335 variation and jitter, both of which are derived from the glottal source. Communities 1 and 2
16
17 336 (orange and yellow) occupied a similar area in the PCA space, although speakers in Community
18
19 337 2 covered a broader range than speakers in Community 1. Community 3 (green), on the other
20
21 338 hand, was more heavily influenced by the glottal measures, in particular pitch variation.
22
23
24

25 339 *Subjective ratings.* Clinical and naïve ratings were very highly correlated, $r(24) = 0.84, p$
26
27 340 < 0.0001 , with high interrater reliability (intraclass correlation coefficient): naïve raters: ICC =
28
29 341 $.957, F(26,192) = 29.9, p < .0001, CI(95\%) 0.929 < ICC < 0.978$; expert raters: ICC = $.882,$
30
31 342 $F(26,60.6) = 10.1, p < .0001, CI(95\%) = 0.785 < ICC < 0.941$. Separate 2-way ANOVAs for
32
33 343 naïve and expert raters showed a significant main effect of diagnosis: naïve: $F(1,22) = 22.18, p <$
34
35 344 $.0001$; expert: $F(1,22) = 51.38, p < .0001$, indicating that both sets of raters were more likely to
36
37 345 rate a talker as atypical or autistic, if that talker was from the ASD group. For the naïve raters
38
39 346 there was a significant main effect of Community, but not for the expert: naïve: $F(2,22) = 4.012,$
40
41 347 $p < 0.05$, expert: $F(2,22) = 2.070, p < .14$). The greatest divergence between naïve and expert
42
43 348 raters was in their assessment of the three autistic speakers and one NT speaker who were placed
44
45 349 in Community 2. Two of these three autistic speakers were rated at a similar level as the NT
46
47 350 speakers by the naïve raters, and a single NT speaker from this community was rated similarly to
48
49 351 the autistic speakers from the other communities. The expert raters, on the other hand, while
50
51
52
53
54
55
56
57
58
59
60

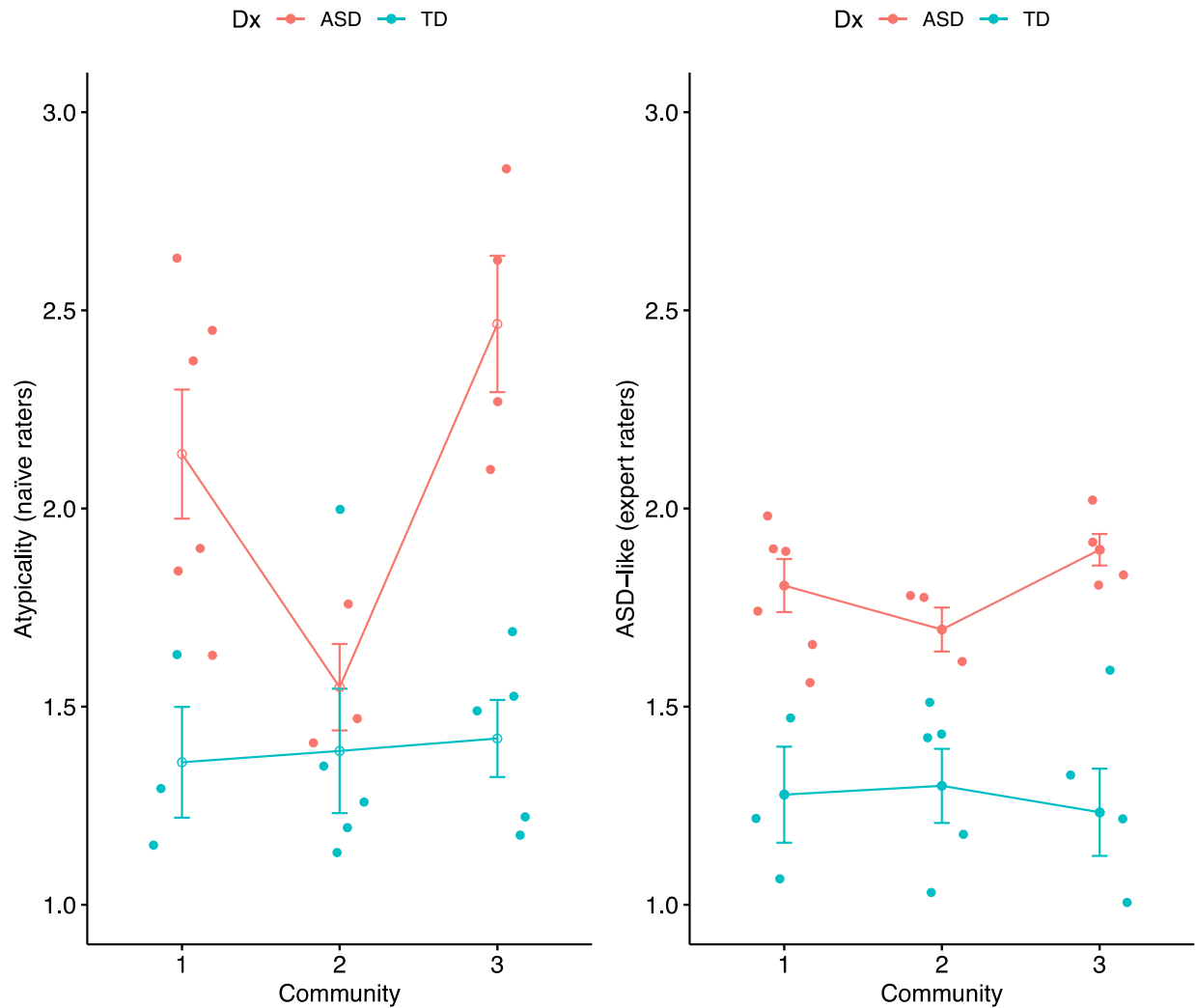
Different in different ways

352 rating the speakers from Community 2 as slightly less likely to have an autism diagnosis than
353 those in Communities 1 and 3, were much more accurate in their assessment.

354 *Symptom severity.* Scores on the Social Responsiveness Scale differed as a function of
355 group (Table 1). However, separate ANOVAs for NT and ASD speakers showed no relation
356 between symptom severity (SRS) and community (NT: $F(1,12) = 0.553, p = 0.471$; ASD:
357 $F(1,10) = 0.004, p = 0.952$). SRS data were missing for one ASD speaker.

358 *Verbal IQ.* Although there was no significant difference between the NT and ASD groups
359 on VIQ (Table 1), because the p-value for this comparison was close to the 0.05 alpha level
360 (0.06), we also performed separate ANOVAs to see whether VIQ variably predicted community
361 membership for these two groups. VIQ was not a significant predictor of community
362 membership for either the NT ($F(1, 12) = 0.01, p = 0.917$) or the ASD ($F(1, 10) = 0.502, p =$
363 0.495) groups.

Different in different ways



364

365 Figure 4: Relationship between ratings by naïve (left) and clinically-trained (right) raters and
 366 community membership for ASD (red) and NT (blue) participants

367

368

Discussion

369

This exploratory study aimed to characterize the acoustic features that underlie

370

perceptually distinctive qualities of autistic speech. Drawing on the scant prior literature

371

examining these features, we expected that variation in fundamental frequency, jitter, speech

372

rate, and articulation rate were likely to be important in differentiating the speech of individuals

Different in different ways

1
2
3 373 with ASD. We compared the acoustic qualities of age- and NVIQ-matched groups of speakers
4
5 374 with a history of autism or typical development as they spoke a series of eight identical
6
7 375 utterances. Given the failure of prior studies to find generalizable group-level differences in
8
9 376 individual acoustic cues, we hypothesized that a network approach might be more successful,
10
11 377 given that it permits the assessment of a constellation of features. In contrast to supervised
12
13 378 predictive approaches, which attempt to build models of the voice and prosody of individuals
14
15 379 with ASD for the purpose of making diagnostic predictions, our approach was exploratory and
16
17 380 unsupervised, and thereby made no *a priori* assumptions about the role of diagnosis. This
18
19 381 method opens the possibility that multiple prosodic and voice quality profiles can be identified in
20
21 382 a data-driven fashion, with the hope that the identification of these profiles can inform not only
22
23 383 clinical practice but also future predictive modelling attempts.

24
25
26 384 Based on the literature, we proposed four hypotheses relevant to identifying a prosodic
27
28 385 profile for speakers with ASD.

29
30
31 386 *H1. Modelling speakers as nodes of acoustic features in a network will allow us to*
32
33 387 *identify coherent clusters of speakers.*

34
35
36 388 To our knowledge, this approach of modelling a group of speakers as nodes in a network
37
38 389 defined by the similarity or dissimilarity of acoustic properties of the voice has not been
39
40 390 previously reported, so our results are necessarily both tentative and exploratory. At the same
41
42 391 time, although the number of participants in the study was small, the separation of speakers into
43
44 392 groups was quite robust. Given four distinct acoustic measures of prosody and voice as measured
45
46 393 in a set of eight utterances (identical across speakers), the algorithm reliably detected three
47
48 394 groups of speakers over many runs with different random seeds. Although this method needs
49
50 395 further exploration and validation, including experimenting with the inclusion of a wider variety
51
52
53
54
55
56
57
58
59
60

1 Different in different ways
2

3 396 of acoustic features, and different types of utterances, the current results suggest that it is
4
5 397 possible to reliably identify clusters of individual speakers on the basis of acoustic variables.
6

7
8 398 *H2. NT speakers and speakers with ASD will tend to cluster differently.*
9

10 399 The speakers were divided into three clusters; ASD and NT speakers were not evenly
11
12 400 distributed over the three groups. Approximately 69% of the speakers in Community 1 but only
13
14 401 37.5% of the speakers in Community 2 were autistic. This result suggests that the acoustic
15
16 402 measures used by the community-detection algorithm to cluster speakers did to some degree
17
18 403 correlate with the presence or absence of ASD. However, Community 3 adds complexity to this
19
20 404 picture. While 17 of the speakers were grouped into two communities which might otherwise be
21
22 405 called the “autistic” community and the “NT” community, over a third of the participants were
23
24 406 clustered in Community 3, which was composed of a nearly equal number of NT and ASD
25
26 407 speakers. Clearly, sorting the speakers by diagnostic group cannot rely on acoustic features
27
28 408 alone, at least not using the features that we chose. At the same time, both naïve and expert raters
29
30 409 were quite successful in their categorization of these speakers.
31
32

33 410 *H3. Clusters of speakers will be characterized by distinctive constellations of acoustic*
34
35 411 *features.*
36
37

38
39 412 Although our results are preliminary, the combined patterns of acoustic features did form
40
41 413 qualitatively different clusters of speech samples. Variation in Communities 1 and 2 were
42
43 414 accounted for primarily by differences in the rhythmic features speech rate and articulation rate,
44
45 415 while Community 3 was mostly characterized by variation in F0 and jitter.
46
47
48

49 416 *H4. Clusters will reflect the subjective ratings of naïve and clinical raters.*
50

51 417 Both naïve and expert raters were very successful in distinguishing autistic and neurotypical
52
53 418 speakers. Interestingly, this was the case even for speakers from Community 3, which consisted
54
55
56
57
58
59
60

Different in different ways

1
2
3 419 of a nearly even mixture of ASD and NT speakers whose voice and speech patterns closely
4
5 420 resembled each other, at least as measured with our four acoustic variables. Since the content of
6
7 421 the sentences uttered by the speakers was the same, these raters likely relied on subtle
8
9
10 422 characteristics of either intonation, voice quality, or both, which at least partially eluded our
11
12 423 algorithmic analysis.

13
14 424 In addition to the limited sample size, there are at least three reasons for the failure of the
15
16 425 algorithmic approach to group the speakers into distinct diagnostic group clusters. **First**, the
17
18 426 network may have included too few acoustic features, or included some features which are not
19
20 427 distinct across groups. **A second possibility** is that the network used the “correct” set of features,
21
22 428 but these features were not accurately or appropriately measured or processed. **Finally**, of
23
24 429 course, the common perception that there is an “autistic voice” may be inaccurate; perhaps there
25
26 430 is in fact no single acoustically-definable profile. We discuss each possibility in turn.

27
28
29 431 Regarding the first possibility, it is certainly plausible that adding other acoustic features
30
31 432 such as shimmer, H1-H2, MFCC’s (Mel Frequency Cepstral Coefficients), or raw spectrograms
32
33 433 to the mix might improve our ability to algorithmically cluster speakers into groups of primarily
34
35 434 ASD or NT speakers. However, this raises new challenges. First, adding more features increases
36
37 435 the risk of overfitting. This could be alleviated by a subsequent feature-reduction step, such as
38
39 436 using principal components analysis (PCA) to identify composite variables with the greatest
40
41 437 predictive power. In addition to overfitting, adding additional features can further contribute to
42
43 438 difficulties in interpreting the communities. The most useful outcome for training new clinicians,
44
45 439 for example, would be to identify a set of easily-identified acoustic features that have been
46
47 440 shown to be characteristic of the speech of autistic people. Neither a long list of relatively
48
49 441 complex acoustic features nor an inscrutable composite component from a PCA achieves this

1 Different in different ways
2

3 442 goal (see supporting information S4 of Fusaroli et al (2021) for a discussion of the interpretation
4
5 443 of PCA components in acoustic analysis). While likely incomplete, a list of approximately three
6
7 444 to five intuitively understandable features may be of the greatest utility for intervention. The
8
9 445 prior literature, together with the current results, indicates that pitch variation and rhythm
10
11 446 features are likely to be important.
12
13

14 447 Regarding the second possibility, one could retain the current list of features, with some
15
16 448 small adjustments. Pause behavior (as indirectly measured by relationship of speech rate to
17
18 449 articulation rate in the present study) appears to be an important distinguishing factor in our data,
19
20 450 but the window for defining a pause could be adjusted in either direction; in this study, it was
21
22 451 defined as silences of 0.2 seconds or longer. Further, a simple measure like speech rate or
23
24 452 articulation rate may not be refined enough to capture the essence of the atypical speech rhythms
25
26 453 that our naïve and expert raters picked up on. For example, there is some research suggesting that
27
28 454 analysis tools, such as recurrence quantification analysis (RQA), that embrace the temporal
29
30 455 aspect of speech and identify recurring patterns over time, can be combined with features such as
31
32 456 pause length to identify speakers with ASD (Fusaroli, Bang, & Weed, 2013; Fusaroli, Grossman,
33
34 457 Cantio, Bilenberg, & Weed, 2015). However, like feature reduction with PCA, this approach
35
36 458 results in a set of predictive features that are less easy to interpret intuitively.
37
38
39

40 459 This brings us to the third possibility: that there is no unique or defining “autistic voice.”
41
42 460 This argument is highly appealing on several grounds. First, the empirical, clinical experience
43
44 461 and especially self-reports from autistic people, all emphasize the tremendous heterogeneity of
45
46 462 ASD; see Mottron & Bzdok (2020) and Waterhouse (2013) for discussion. As an example,
47
48 463 although reduced pitch variation may be typical of many, though not all, autistic people, the
49
50 464 speaker in our sample most consistently identified by both naïve and expert raters as either
51
52
53
54
55
56
57
58
59
60

1 Different in different ways
2

3 465 “atypical” or “likely to be diagnosed with ASD” had a relatively *high* degree of pitch variation.
4
5 466 Clearly, an algorithm that builds predictions based alone or in part on pitch variation would be
6
7 467 likely to misclassify this speaker as NT, but the human raters were not in doubt. A qualitative
8
9 468 assessment of this speaker’s speech offers some clues to this apparent contradiction, as discussed
10
11 469 further below.

12
13
14 470 *Different in different ways?*

15
16
17 471 The network community-detection algorithm split the speakers into clusters that were
18
19 472 related to diagnosis, but did not separate them evenly into two clusters defined by diagnosis. This
20
21 473 lends support to the idea that while acoustic features of prosody and voice are related to
22
23 474 diagnosis, there is no single acoustic feature or even constellation of features which is
24
25 475 consistently associated with ASD or with the perception of atypicality. This tracks with the
26
27 476 initial observation that radically different qualitative descriptions have been used of the prosody
28
29 477 of people with autism: monotone versus singsong, slow versus fast. There may be a similar
30
31 478 underlying cause that results in these seemingly opposite outcomes; for example, an awareness
32
33 479 that NT speakers use pitch to communicate intention or attitude, but a lack of intuitive
34
35 480 understanding for how this is done might lead speakers with ASD to either ignore this aspect of
36
37 481 speech, or to use it in an atypical fashion. Of course, the degree to which speakers are more or
38
39 482 less monotone or dynamic varies in neurotypical speakers as well.

40
41
42 483 As mentioned above, the speaker most consistently rated by both expert and naïve raters
43
44 484 as atypical was one of the three speakers with ASD in Community 2. This speaker’s speech is
45
46 485 indeed quite dynamic in terms of pitch variation, and **the peaks and valleys in** pitch subjectively
47
48 486 **seem to fall in about the same places as in** neurotypical speakers’ productions of the same
49
50 487 sentences (consistent with findings in Geelhand, Papastamou & Kissine, 2021). However, there
51
52
53
54
55
56
57
58
59
60

1 Different in different ways

2
3 488 is an exaggerated and distinctive quality to this speaker's use of pitch variation. Furthermore,
4
5 489 another salient aspect of this speaker's distinctive speech was *not* among the acoustic features
6
7 490 measured in this study; several vowels were produced with an unusual quality, sounding like
8
9
10 491 phonemes produced by a speaker of English as a second language. **Indeed, broad measures such**
11
12 492 **as pitch variation and articulation rate cannot capture the many small but important differences in**
13
14 493 **sentence intonation that can convey critical information to the hearer.**

15
16
17 494 How should this sort of subjective assessment of participants' speech be incorporated into
18
19 495 our models? If vowel quality is of interest, then formant measures should be included into the
20
21 496 acoustic profile. Alternatively, the atypicality of this speaker's speech could reflect an interaction
22
23 497 of vowel quality with pitch variation and timing, whereas a different speaker's speech might be
24
25 498 rated as similarly atypical, reflecting an entirely different combination of interacting features.
26
27
28 499 The proposition of building a generalizable predictive model of the autistic voice reflects the
29
30 500 enormous challenges in understanding the many-to-many mapping problem in the acoustics of
31
32 501 speech (Lieberman et al., 1967). We propose that all three possibilities described above (the
33
34 502 inclusion of too few acoustic features, the misquantification of features, and the possibility that
35
36 503 there is no single acoustically-definable profile of the autistic voice) present important challenges
37
38 504 to building predictive models of atypical speech. They are not mutually exclusive, and all require
39
40 505 further investigation. How researchers choose to address them will be in part based on the goals
41
42 506 guiding the research.

43
44
45
46
47 507 An important aspect of ensuring progress on these complicated questions is the open
48
49 508 sharing of data and code. While sharing raw audio files presents important data privacy
50
51 509 problems, sharing anonymized acoustic features, as well as the code used for extracting and
52
53 510 analyzing those features will be crucial for advancing research in this domain. To this end, we

Different in different ways

511 provide full data and code at the online repository OSF:

512 https://osf.io/zw67n/?view_only=63441e15ec264184bb7b0b22c4140a22.

513 *Limitations and future research*

514 The present study is exploratory in nature, and as such, there are several obvious avenues
515 for future research. Among the most important next steps will be expanding the number and
516 variety of features used to model the speakers' productions. As mentioned in the introduction,
517 sentence-level pitch variation alone is a poor measure of intonation. Emphasizing a word in
518 conversation, for example, can be achieved with a change of pitch, but also with a change in
519 intensity, the addition of a pause, the lengthening or shortening of a vowel, or a combination of
520 all of these. In addition, while scripted productions like the ones used in our study are useful
521 because they are easily comparable, they are clearly less than ideal as a means of eliciting
522 conversational prosody. Another obvious avenue for future research is thus the extension of the
523 methods suggested here to more naturalistic data, including conversational data. Not only is
524 prosody more vital for conducting successful conversations than it is for producing utterances in
525 the lab, but the artificial setting of our elicitation may have impacted the speakers in a variety of
526 ways, including the stress of having to "perform". A third opportunity for improvement will be
527 obtaining higher-quality recordings, reducing background noise, and being more discerning in
528 setting individualized windows for pitch extraction. Finally, it will be important to expand this
529 research to larger, more diverse samples (including a wider range of language abilities, and
530 speakers more diverse in sociodemographic factors) and to include other languages.
531 Unfortunately, we expect that it will be difficult to achieve all of these within the context of a
532 single study. For example, more naturalistic data will likely result in noisier recordings, larger
533 samples will make careful, individualized feature extraction more difficult, and different

1 Different in different ways
2

3 534 languages may make different use of prosodic features, making generalization difficult. Despite
4
5 535 these difficulties, we find that the present results suggest the need to think in a more nuanced
6
7 536 way about when, how, and whether the speech of people with ASD diverges from that of TD
8
9 537 speakers, especially in conversation.
10
11

12 538 When de Marchena et al (2017) asked clinicians to identify the behavioral features they
13
14 539 most associated with ‘frank’ ASD, atypical prosody was among the most common features
15
16 540 mentioned, and they propose the pursuit of a ‘narrow frank feature (e.g. avoidance of eye
17
18 541 contact, or unusual prosody)’ as an effective way to study and understand frank presentation
19
20 542 autism (pg. 660). While our data suggest a viable means to explore and quantify these prosodic
21
22 543 differences, they also point to the fact that even the relatively ‘narrow’ feature of prosody may
23
24 544 contain multitudes.
25
26
27

28 545 **Conclusion**

29
30
31 546 Our main goal in this investigation was to test the use of a network-based community-
32
33 547 detection algorithm in investigating atypical prosody and voice quality in ASD. Using four
34
35 548 acoustic features, we built a network of speakers which reliably consisted of three communities:
36
37 549 speakers that were more acoustically alike within their community and more dislike speakers
38
39 550 outside their community. While two of the communities tended toward either **more ASD or more**
40
41 551 **NT speakers**, the third community was populated **by a nearly even split** of NT and ASD
42
43 552 speakers. Both expert and naïve raters were highly successful in identifying autistic speakers,
44
45 553 regardless of which acoustic community the speakers had been placed in by the algorithm.
46
47 554 Although we regard this exploratory study as a first proof of concept, we suggest that using
48
49 555 network analysis to identify clusters of acoustically-similar speakers may lead to important
50
51
52
53 556 insights into how atypical prosody, while always atypical, may be different in different ways.
54
55
56
57
58
59
60

Different in different ways

557 **References**

558 Al-Qatab, B. A., & Mustafa, M. B. (2021). Classification of dysarthric speech according
559 to the severity of impairment: An analysis of acoustic features. *IEEE Access*, 9, 18183-18194.

560 Asperger, H. (1991). "Autistic psychopathy" in childhood. In U. Frith (Ed.), *Autism and*
561 *Asperger syndrome* (pp. 37–91). Cambridge: Cambridge University Press. (Original work
562 published 1944)

563 American Psychological Association (2013). *Diagnostic and statistical manual of mental*
564 *disorders (DSM-5)*. American Psychiatric Pub.

565 Berends, M. S., Luz, C. F., Friedrich, A. W., Sinha, B. N., Albers, C. J., & Glasner, C.
566 (2021). AMR-An R Package for working with antimicrobial resistance data. *BioRxiv*, 810622.

567 Boersma, P., & Weenink, D. (2001). *Praat, a system for doing phonetics by computer*.

568 Bone, D., Lee, C.-C., Black, M. P., Williams, M. E., Lee, S., Levitt, P., & Narayanan, S.
569 (2014). The psychologist as an interlocutor in autism spectrum disorder assessment: Insights
570 from a study of spontaneous prosody. *Journal of Speech, Language, and Hearing Research*,
571 57(4), 1162–1177.

572 Bottalico, P., Codino, J., Cantor-Cutiva, L. C., Marks, K., Nudelman, C. J., Skeffington,
573 J., Shrivastav, R., Jackson-Menaldi, M. C., Hunter, E. J., & Rubin, A. D. (2020). Reproducibility
574 of Voice Parameters: The Effect of Room Acoustics and Microphones. *Journal of Voice*, 34(3),
575 320–334. <https://doi.org/10.1016/j.jvoice.2018.10.016>

576 Borrie, S. A., Barrett, T. S., Liss, J. M., & Berisha, V. (2020). Sync pending:
577 Characterizing conversational entrainment in dysarthria using a multidimensional, clinically
578 informed approach. *Journal of Speech, Language, and Hearing Research*, 63(1), 83-94.

Different in different ways

579 Brockmann-Bauser, M., Beyer, D., & Bohlender, J. E. (2014). Clinical relevance of
580 speaking voice intensity effects on acoustic jitter and shimmer in children between 5;0 and 9;11
581 years. *International Journal of Pediatric Otorhinolaryngology*, 78(12), 2121–2126.

582 <https://doi.org/10.1016/j.ijporl.2014.09.020>

583 Bryant, G. A., Fessler, D. M., Fusaroli, R., Clint, E., Aarøe, L., Apicella, C. L., ... others.
584 (2016). Detecting affiliation in co-laughter across 24 societies. *Proceedings of the National
585 Academy of Sciences*, 113(17), 4682–4687.

586 Cappadocia, M. C., Weiss, J. A., & Pepler, D. (2012). Bullying experiences among
587 children and youth with autism spectrum disorders. *Journal of Autism and Developmental
588 Disorders*, 42(2), 266–277.

589 Chen, L., Ravichandran, V., & Stolcke, A. (2021). Graph-based label propagation for
590 semi-supervised speaker identification. arXiv preprint arXiv:2106.08207.

591 Clauset, A., Newman, M. E., & Moore, C. (2004). Finding community structure in very
592 large networks. *Physical review E*, 70(6), 066111.

593 Constantino, J. N., Davis, S. A., Todd, R. D., Schindler, M. K., Gross, M. M., Brophy, S.
594 L., ... Reich, W. (2003). Validation of a brief quantitative measure of autistic traits: Comparison
595 of the social responsiveness scale with the autism diagnostic interview-revised. *Journal of
596 Autism and Developmental Disorders*, 33(4), 427–433.

597 Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network
598 research. *InterJournal, Complex Systems*, 1695(5), 1–9.

599 Cummins, N., Scherer, S., Krajewski, J., Schnieder, S., Epps, J., & Quatieri, T. F. (2015).
600 A review of depression and suicide risk assessment using speech analysis. *Speech
601 Communication*, 71, 10–49.

Different in different ways

602 Cutler, A., Dahan, D., & Van Donselaar, W. (1997). Prosody in the comprehension of
603 spoken language: A literature review. *Language and speech*, 40(2), 141-201.

604 Dahlgren, S., Sandberg, A. D., Strömbergsson, S., Wenhov, L., Råstam, M., &
605 Nettelblatt, U. (2018). Prosodic traits in speech produced by children with autism spectrum
606 disorders—Perceptual and acoustic measurements. *Autism and Developmental Language*
607 *Impairments*, 3, 2396941518764527.

608 Degottex, G., Kane, J., Drugman, T., Raitio, T., & Scherer, S. (2014). *Covarep - a*
609 *collaborative voice analysis repository for speech technologies* (pp. 960–964). IEEE.

610 Diehl, D. B., Joshua J. Watson. (2009). An acoustic analysis of prosody in high-
611 functioning autism. *Applied Psycholinguistics*, 30(3), 385–404.

612 Diehl, J. J., & Paul, R. (2013). Acoustic and perceptual measurements of prosody
613 production on the profiling elements of prosodic systems in children by children with autism
614 spectrum disorders. *Applied Psycholinguistics*, 34(1), 135–161.

615 Djelantik, A. M. J., Robinaugh, D. J., Kleber, R. J., Smid, G. E., & Boelen, P. A. (2020).
616 Symptomatology following loss and trauma: Latent class and network analyses of prolonged
617 grief disorder, posttraumatic stress disorder, and depression in a treatment-seeking trauma-
618 exposed sample. *Depression and Anxiety*, 37(1), 26–34.

619 Drugman, T., & Alwan, A. (2019). Joint robust voicing detection and pitch estimation
620 based on residual harmonics. arXiv preprint *arXiv:2001.00459*.

621 Epskamp, S., Cramer, A. O. J., Waldorp, L. J., Schmittmann, V. D., & Borsboom, D.
622 (2012). qgraph: Network visualizations of relationships in psychometric data. *Journal of*
623 *Statistical Software*, 48(4), 1–18. Retrieved from <http://www.jstatsoft.org/v48/i04/>

1 Different in different ways

2
3 624 Eskenazi, L., Childers, D. G., & Hicks, D. M. (1990). Acoustic correlates of vocal
4
5 625 quality. *Journal of Speech, Language, and Hearing Research*, 33(2), 298–306.

6
7 626 Filipe, M. G., Frota, S., Castro, S. L., & Vicente, S. G. (2014). Atypical prosody in
8
9 627 Asperger syndrome: Perceptual and acoustic measurements. *Journal of Autism and*
10
11 628 *Developmental Disorders*, 44(8), 1972-1981.

12
13
14 629 Fusaroli, R., Bang, D., & Weed, E. (2013). Non-linear analyses of speech and prosody in
15
16 630 Asperger's syndrome. *International meeting for autism research*.

17
18
19 631 Fusaroli, R., Grossman, R., Bilenberg, N., Cantio, C., Jepsen, J. R. M., & Weed, E.
20
21 632 (2021). Toward a cumulative science of vocal markers of autism: A cross-linguistic
22
23 633 meta-analysis-based investigation of acoustic markers in American and Danish autistic children.
24
25 634 *Autism Research*, 1-12. doi: 10.1002/aur.2661

26
27
28 635 Fusaroli, R., Grossman, R., Cantio, C., Bilenberg, N., & Weed, E. (2015). *The temporal*
29
30 636 *structure of the autistic voice: A cross-linguistic investigation*. Poster session presented at the
31
32 637 *International Meeting for Autism Research*.

33
34
35 638 Fusaroli, R., Lambrechts, A., Bang, D., Bowler, D., & Gaigg, S. (2017). Is voice a marker
36
37 639 for autism spectrum disorder? A systematic review and meta-analysis". *Autism Research*, 10(3),
38
39 640 384–407.

40
41
42 641 Gamer, M., Lemon, J., Gamer, M. M., Robinson, A., & Kendall's, W. (2012). irr:
43
44 642 Various Coefficients of Interrater Reliability and Agreement. R package version 0.84.1.

45
46 643 <https://CRAN.R-project.org/package=irr>

47
48
49 644 Geelhand, P., Papastamou, F., & Kissine, M. (2021). How do autistic adults use syntactic
50
51 645 and prosodic cues to manage spoken discourse?. *Clinical Linguistics and Phonetics*, 35(12),
52
53 646 1184-1209.

Different in different ways

647 Grossman, R. (2015). Judgments of social awkwardness from brief exposure to children
648 with and without high-functioning autism. *Autism, 19*(5), 580–587.

649 Hillenbrand, J., & Houde, R. A. (1994). Acoustic correlates of breathy vocal quality:
650 Dysphonic voices and continuous speech. *Journal of Speech, Language, and Hearing Research,*
651 *39*(2), 311–321.

652 Huang, D. Z., Minifie, F. D., Kasuya, H., & Lin, S. X. (1995). Measures of vocal
653 function during changes in vocal effort level. *Journal of Voice, 9*(4), 429–438.
654 [https://doi.org/10.1016/S0892-1997\(05\)80206-1](https://doi.org/10.1016/S0892-1997(05)80206-1)

655 Kissine, M., Geelhand, P., Philippart De Foy, M., Harmegnies, B., & Deliens, G. (2021).
656 Phonetic inflexibility in autistic adults. *Autism Research, 14*(6), 1186–1196.

657 Kissine, M. & Geelhand, P. (2019). Brief Report: Acoustic Evidence for Increased
658 Articulatory Stability in the Speech of Adults with Autism Spectrum Disorder. *J. Autism Dev.*
659 *Disord. 49*, 2572–2580

660 Klatt, D. H., & Klatt, L. C. (1990). Analysis, synthesis, and perception of voice quality
661 variations among female and male talkers. *The Journal of the Acoustical Society of America,*
662 *87*(2), 820–857.

663 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass
664 correlation coefficients for reliability research. *Journal of Chiropractic Medicine, 15*(2), 155-
665 163.

666 Kreiman, J., & Gerratt, B. (2010). Perceptual sensitivity to first harmonic amplitude in
667 the voice source. *The Journal of the Acoustical Society of America, 128*(4), 2085–2089.

668 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., & Studdert-Kennedy, M. (1967).
669 Perception of the speech code. *Psychological Review, 74*(6), 431.

Different in different ways

670 Lord, C., Risi, S., Lambrecht, L., Cook, E. H., Leventhal, B. L., DiLavore, P. C., ...

671 Rutter, M. (2000). The Autism Diagnostic Observation Schedule—Generic: A standard measure
672 of social and communication deficits associated with the spectrum of autism. *Journal of Autism
673 and Developmental Disorders, 30*(3), 205–223.

674 Low, D. M., Bentley, K. H., & Ghosh, S. S. (2020). Automated assessment of psychiatric
675 disorders using speech: A systematic review. *Laryngoscope Investigative Otolaryngology, 5*(1),
676 96-116.

677 Marchena, A. de, & Miller, J. (2017). “Frank” presentations as a novel research construct
678 and element of diagnostic decision-making in autism spectrum disorder. *Autism Research, 10*(4),
679 653–662.

680 Maslov, S., & Sneppen, K. (2002). Specificity and stability in topology of protein
681 networks. *Science, 296*(5569), 910-913.

682 Mayo, J. (2015). Prosodic Phrasing in Adolescents with High Functioning Autism:
683 Production Following Intervention and Under Dual Load Conditions. Ph.D. Thesis, University of
684 Connecticut.

685 McCann, J., & Peppé, S. (2003). Prosody in autism spectrum disorders: A critical review.
686 *Language and Communication Disorders, 38*(4), 325–350.

687 Mittal, V., & Sharma, R. K. (2021). Machine learning approach for classification of
688 Parkinson disease using acoustic features. *Journal of Reliable Intelligent Environments, 7*(3),
689 233-239.

690 Mesibov, G. B. (1992). Treatment issues with high-functioning adolescents and adults
691 with autism. In *High-functioning individuals with autism* (pp. 143–155). Springer.

Different in different ways

692 Mottron, L., & Bzdok, D. (2020). Autism spectrum heterogeneity: Fact or artifact?

693 *Molecular Psychiatry*, 25(12), 3178-3185.

694 Murphy, P. J. (2000). Spectral characterization of jitter, shimmer, and additive noise in
695 synthetically generated voice signals. *The Journal of the Acoustical Society of America*, 107(2),
696 978-988.

697 Nadig, A., & Shaw, H. (2012). Acoustic and perceptual measurement of expressive
698 prosody in high-functioning autism: Increased pitch range and what it means to listeners. *Journal*
699 *of Autism and Developmental Disorders*, 42(4), 499–511.

700 Nakai, Y., Takashima, R., Takiguchi, T., & Takada, S. (2014). Speech intonation in
701 children with autism spectrum disorder. *Brain and Development*, 36(6), 516–522.

702 Parola, A., Simonsen, A., Bliksted, V., & Fusaroli, R. (2020). Voice patterns in
703 schizophrenia: A systematic review and Bayesian meta-analysis. *Schizophrenia Research*, 216,
704 24–40.

705 Parola, A., Simonsen, A., Lin, J. M., Zhou, Y., Huiling, W., Ubukata, S., ... & Fusaroli,
706 R. (2022). Voice patterns as markers of schizophrenia: building a cumulative generalizable
707 approach via cross-linguistic and meta-analysis based investigation. medRxiv.

708 Peng, C., Chen, W., Zhu, X., Wan, B., & Wei, D. (2007). Pathological voice
709 classification based on a single Vowel's acoustic features. In 7th IEEE International Conference
710 on Computer and Information Technology (CIT 2007) (pp. 1106-1110). IEEE.

711 Peppé, S. J. E. (2009). Why is prosody in speech-language pathology so difficult.
712 *International Journal of Speech and Language Pathology*, 11(4), 258–271.

1 Different in different ways

- 2
3 713 Peppe, S., McCann, J., Gibbon, F., O'Hare, A., & Rutherford, M. (2007). Receptive and
4
5 714 expressive prosodic ability in children with high-functioning autism. *Journal of Speech,*
6
7 715 *Language and Hearing Research, 50*(4), 1015-1028.
- 8
9
10 716 Quene H (2022). *_hqmisc: Miscellaneous Convenience Functions and Dataset*. R package
11
12 717 *version 0.2-1*, <https://CRAN.R-project.org/package=hqmisc>
13
14
15 718 R Core Team. (2022). *R: A Language and Environment for Statistical Computing*, R
16
17 719 *Foundation for Statistical Computing, Vienna, Austria*. <https://www.R-project.org/>.
18
- 19 720 Redford, M. A., Kapatsinski, V., & Cornell-Fabiano, J. (2018). Lay listener classification
20
21 721 and evaluation of typical and atypical children's speech. *Language and Speech, 61*(2), 277–302.
- 22
23
24 722 Reichardt, J., & Bornholdt, S. (2006). Statistical mechanics of community detection.
25
26 723 *arXivPhys. Rev. E 74 (2006) 016110*, 0603718v1.
- 27
28 724 Roid, G. H. (2003) *Stanford-Binet Intelligence Scales: Fifth Edition*. Itasca, IL: Riverside
- 29
30
31 725 Sasson, N. J., Faso, D. J., Nugent, J., Lovell, S., Kennedy, D. P., & Grossman, R. B.
32
33 726 (2017). Neurotypical peers are less willing to interact with those with autism based on thin slice
34
35 727 judgments. *Scientific Reports, 7*(1).
- 36
37 728 Shattuck-Hufnagel, S., & Turk, A. E. (1996). A prosody tutorial for investigators of
38
39 729 auditory sentence processing. *Journal of psycholinguistic research, 25*(2), 193-247.
- 40
41
42 730 Shriberg, L. D., Paul, R., Black, L. M., & Santen, J. P. van. (2011). The hypothesis of
43
44 731 apraxia of speech in children with autism spectrum disorder. *Journal of Autism and*
45
46 732 *Developmental Disorders, 41*(4), 405–426.
- 47
48
49 733 Shriberg, L. D., Paul, R., McSweeny, J. L., Klin, A., Cohen, D. J., & Volkmar, F. R.
50
51 734 (2001). Speech and prosody characteristics of adolescents and adults with high-functioning
52
53 735 autism and asperger syndrome. *Journal of Speech, Language, and Hearing Research*.

Different in different ways

- 736 Shriberg, L. D., & Widder, C. J. (1990). Speech and prosody characteristics of adults
737 with mental retardation. *Journal of Speech, Language, and Hearing Research*, 33(4), 627–653.
- 738 Schuller, B., Steidl, S., Batliner, A., Hirschberg, J., Burgoon, J. K., Baird, A., ... &
739 Evanini, K. (2016). The interspeech 2016 computational paralinguistics challenge: Deception,
740 sincerity & native language. In 17TH Annual Conference of the International Speech
741 Communication Association (Interspeech 2016), Vols 1-5 (pp. 2001-2005).
- 742 Team, R. C. (2018). *R: A language and environment for statistical computing*. Vienna: R
743 Foundation for Statistical Computing.
- 744 Traag, V. A., & Bruggeman, J. (2008). Community detection in networks with positive
745 and negative links. *arXivPhys. Rev. E* 80, 036115, (2009), 0811.2329v3.
- 746 Van Bourgondien, M. E., & Woods, A. V. (1992). Vocational possibilities for high-
747 functioning adults with autism. In *High-functioning individuals with autism* (pp. 227–239).
748 Springer.
- 749 Wang, J., Xiao, X., Wu, J., Ramamurthy, R., Rudzicz, F., & Brudno, M. (2021). Speaker
750 attribution with voice profiles by graph-based semi-supervised learning. arXiv preprint
751 arXiv:2102.03634.
- 752 Waterhouse, L. (2013). *Rethinking autism: Variation and complexity*. Academic Press.
- 753 Weed, E., & Fusaroli, R. (2020). Acoustic measures of prosody in right-hemisphere
754 damage: A systematic review and meta-analysis. *Journal of Speech, Language, and Hearing
755 Research*, 63(6), 1762–1775.
- 756 Wiig, E. H., Semel, E. M., & Secord, W. (2003). *CELF 5: Clinical Evaluation of
757 Language Fundamentals*. Pearson/PsychCorp.

Different in different ways

1
2
3 758 Wolfe, V., Fitch, J., & Martin, D. (1997). Acoustic measures of dysphonic severity across
4
5 759 and within voice types. *Folia Phoniatica Et Logopaedica*, 49(6), 292–299.

6
7 760 Yumoto, E., Gould, W. J., & Baer, T. (1982). Harmonics-to-noise ratio as an index of the
8
9 761 degree of hoarseness. *The Journal of the Acoustical Society of America*, 71(6), 1544–1550.

10
11
12 762
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

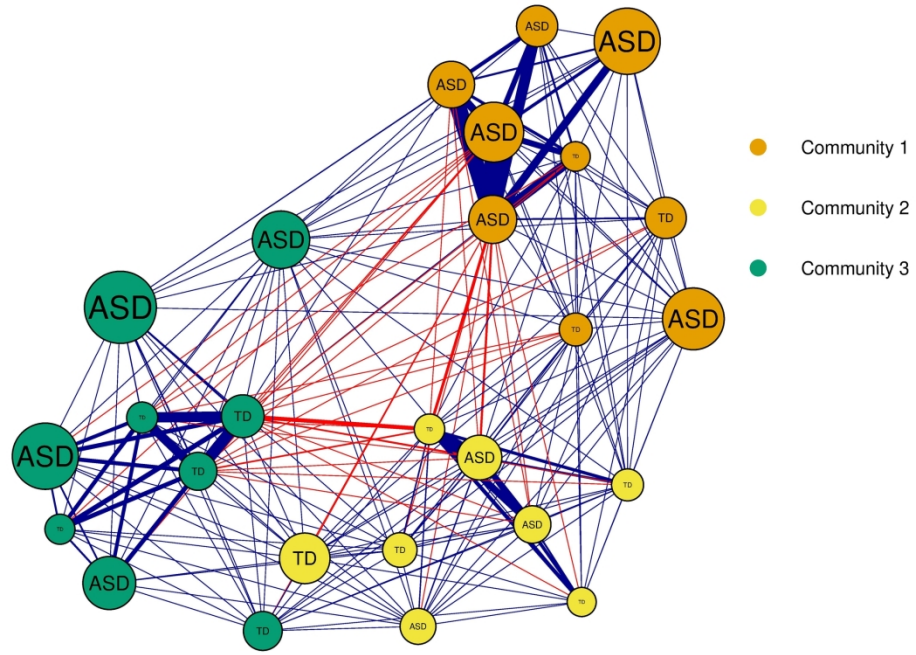


Figure 1: Network of acoustically-defined communities. Each node represents a single speaker. Blue connections between nodes represent positive correlations, red lines represent negative correlations. Line thickness indicates the absolute strength of the correlation.

210x181mm (300 x 300 DPI)

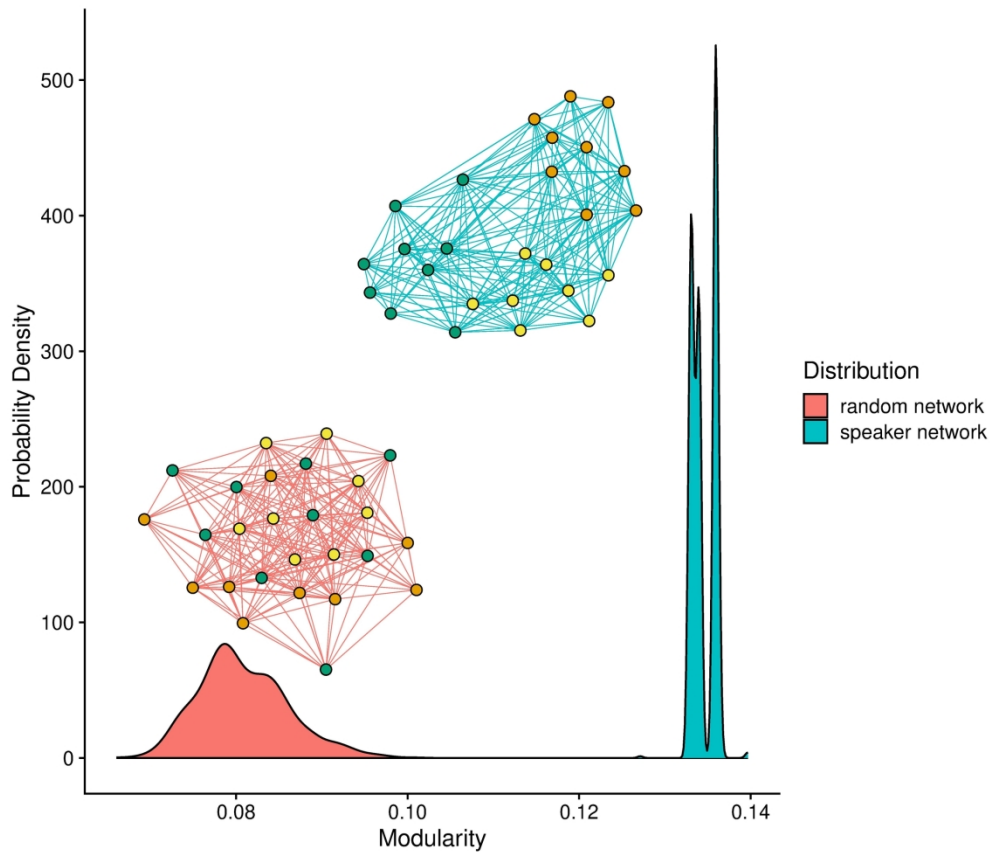


Figure 2: Distributions of modularity scores for 1000 random networks built with the same constraints as the actual network of speakers, and the distribution of modularity scores for 1000 runs of community detection on the speaker network. Also displayed are a sample random network and the actual speaker network. The colors of the connections in the sample networks indicate which distribution they are drawn from. The colors of the nodes indicate community membership.

210x181mm (300 x 300 DPI)

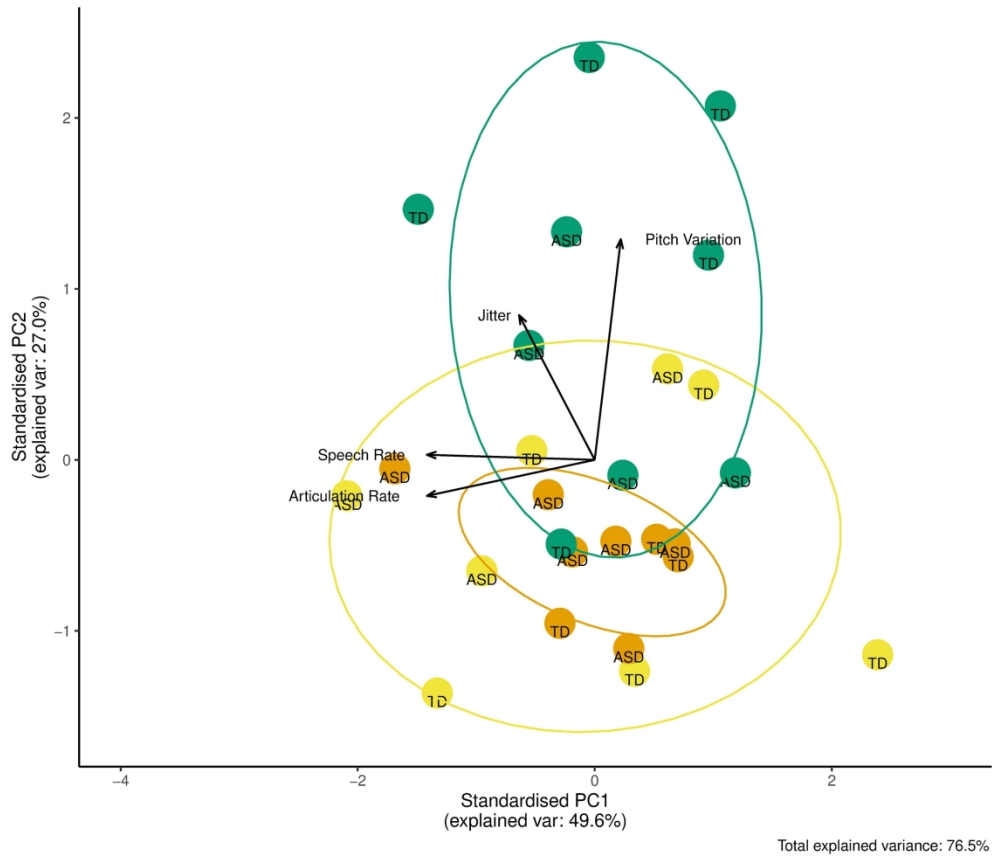


Figure 3: Speakers and communities visualized along the first two dimensions of a principal components decomposition. Ellipses colors and point colors indicate communities detected by the spin glass algorithm (orange = Community 1, yellow = Community 2, green = Community 3). Arrows indicate the influence of each of the features on the components.

210x181mm (300 x 300 DPI)

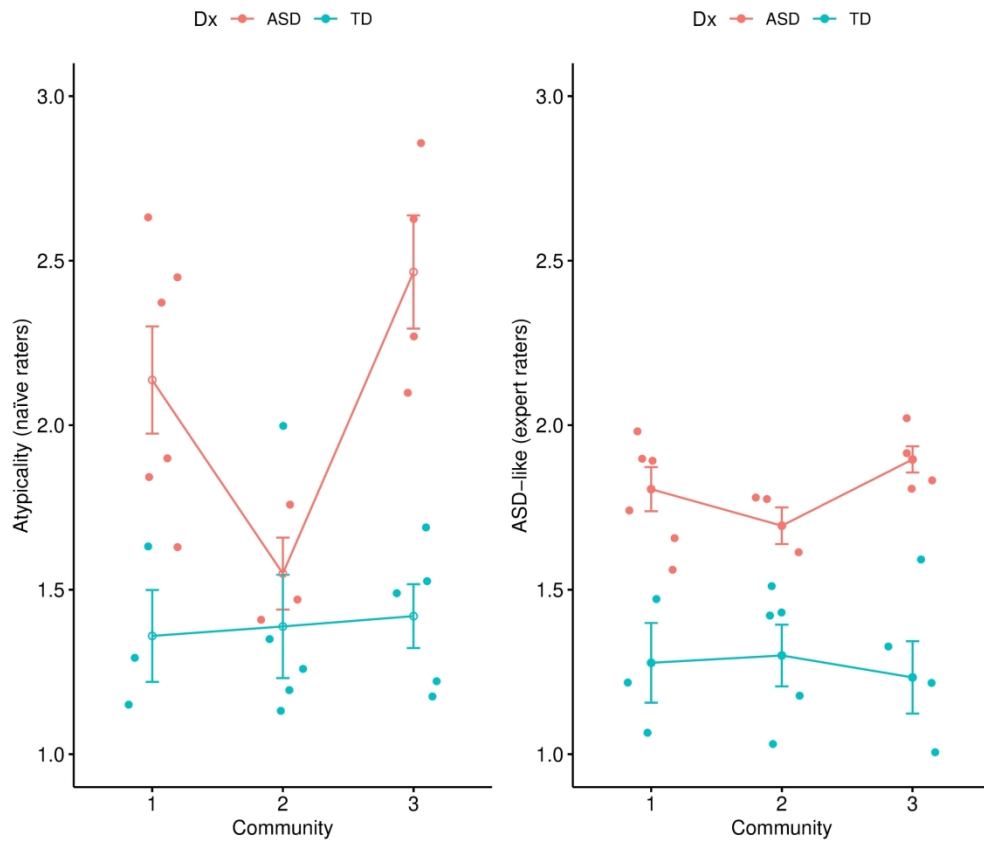


Figure 4: Relationship between ratings by naïve (left) and clinically-trained (right) raters and community membership for ASD (red) and NT (blue) participants

210x181mm (300 x 300 DPI)