

Identifying stable speech-language markers of autism in children: Preliminary evidence from a longitudinal telephony-based study

Abstract

This study examined differences in linguistic features produced by autistic and neurotypical (NT) children during brief picture descriptions, and assessed feature stability over time. Weekly speech samples from well-characterized participants were collected using a telephony system designed to improve access for geographically isolated and historically marginalized communities. Results showed stable group differences in certain acoustic features, some of which may potentially serve as key outcome measures in future treatment studies. These results highlight the importance of eliciting semi-structured speech samples in a variety of contexts over time, and adds to a growing body of research showing that fine-grained naturalistic communication features hold promise for intervention research.

1 Introduction

Natural sampling is a rich approach to investigating speech and language in autistic children. Previous studies have shown that language behavior in autism differs from neurotypical (NT) patterns in a number of ways. For example, autistic children who are more severely impacted have been shown to produce less speech,¹ slower speech,^{2,3} and speech with atypical voice quality¹ compared to NT peers. It has also been observed that autistic children's prosody differs from NT children, with pitch descriptions ranging from sing-songy to monotonous.³ In the lexical domain, prior research has shown that autistic children use nouns and cognitive words differently than NT peers when narrating a story from a picture,⁴ use different patterns of filler words during clinical assessments,⁵ and talk less about social topics during get-to-know-you conversations.⁶ Research in this domain

continues to emerge, but samples remain small and results occasionally conflict or fail to replicate.

Prior studies of natural language in autism used a variety of data collection and analysis methods that could critically affect results and may have led to conflicting findings. For example, the presence of an unfamiliar adult during in-person or remote elicitations could adversely impact the behavior of autistic children, thus reducing the quality and informativeness of their language samples.⁷ Also, children's linguistic behavior might differ depending on the specifics of the elicitation task in a given study, i.e., whether natural conversations or semi-structured speech tasks are used, and the characteristics of certain elicitation stimuli.

In order to develop scalable, cost-effective, objective intervention progress monitoring systems of autistic symptoms using speech as a primary target, it is necessary to understand how contextual and testing factors affect children's behavior. Then, it will be possible to identify robust features that reliably index autism symptoms across heterogeneous testing conditions. Toward this goal, we developed a telephony protocol to examine how various factors affect speech performance in autistic children and adolescents. Telephony has particular potential to address service and monitoring gaps for autistic and NT children from historically marginalized and/or low-resource communities, and is a useful alternative to in-person data collection during the COVID-19 pandemic. The final battery of our protocol consisted of seven versions of seven tasks that a parent or legal guardian could independently facilitate. In this preliminary report from an on-going study, we assessed children's speech and language features during one of the seven tasks (picture descriptions) collected in the first and second phone sessions. Our goals were to (1) identify diagnostic group differences in automated speech and language features that are stable over

82 time, and (2) examine potential effects of staff vs.
83 parent administration in each diagnostic group.

84 2 Methods

85 2.1 Participants

86 Study inclusion and exclusion criteria are
87 included in the Appendix. In this report, we
88 analyzed data from 29 children who successfully
89 completed two sessions. Participant groups were
90 matched on age, full-scale IQ, and self-reported
91 race (Table 1). Groups were not matched on sex
92 ($p=0.015$), which is expected due to the prevalence

	Autism (n=13)	NT (n=16)	p-value
Age (years)	9.8 (2.5)	9.6 (2.6)	0.767
Sex (%)	10 boys (76.9%)	6 boys (37.5%)	0.015
Full scale IQ	115.1 (15.4)	119.1 (13.7)	0.469
Race (%)	4 non- whites (30.8%)	5 non- whites (31.3%)	0.69
SCQ (total)	17.0 (6.6)	1.2 (1.1)	<0.001
SRS-2 (total)	70.5 (7)	42.1 (3.5)	<0.001
CCC-2 (speech)	9.2 (2.5)	11.8 (0.8)	<0.001
CCC-2 (non- verbal)	5.5 (2.2)	11.8 (1.3)	<0.001

Table 1: Demographic and clinical characteristics of the participants. SCQ: Social communication questionnaire, ⁸ SRS: Social responsiveness scale, ⁹ CCC: Children’s communication checklist. ¹⁰

93 of ASD in boys, ¹¹ and we are currently addressing
94 with targeted recruitment. One autistic participant
95 identified as non-binary. Autism and NT groups
96 differed in several clinical ratings (Table 1).

97 2.2 Data collection and annotation

98 We developed a telephony platform to support
99 single and dual speaker modes. This platform
100 consisted of a high-availability server, voice over
101 internet protocol (VoIP) service by Vonage,
102 telephony software framework (Asterisk 13.18.3),
103 a relational database, and telephony applications.

104 Prior to the first official data collection call,
105 study staff held an “informational call” with the
106 participating parent to review standard elicitation
107 methods to be utilized across sessions. During the
108 first session with the child, study staff remained on
109 the line and facilitated tasks with the parent and

110 child. During the second session, children and
111 parents independently completed all seven tasks on
112 their own. Children described different pictures
113 during the first and second sessions, and the second
114 session was collected approximately one week
115 after the first session was completed.

116 Recordings were transcribed by trained
117 annotators using a web-based transcription tool
118 with a built-in speech activity detector (SAD)
119 function. For dual speaker mode recordings, SAD
120 ran on each channel separately. Annotators also
121 corrected speech segment boundary errors.

122 2.3 Acoustic and text features

123 Words were automatically tagged for part-of-
124 speech (POS) categories using spaCy. ¹² POS
125 categories, fillers, partial words, repetitions, and
126 “hm” were counted separately and converted to
127 counts per 100 words. Content words were rated
128 for word frequency, ¹³ concreteness, ¹⁴ ambiguity,
129 ¹⁵ age of acquisition (AoA), ¹⁶ and familiarity. ¹⁶ We
130 also ran the Language Inquiry and Word Count
131 program ¹⁷ to calculate additional word-level
132 measures found to be useful in clinical populations.

133 For acoustic processing, stereo recordings were
134 split into single channels for precise audio
135 processing. We extracted low-level descriptors of
136 pitch, jitter, shimmer, harmonic-to-noise ratio
137 (HNR), and four spectral moments (1st order:
138 centroid, 2nd order: standard deviation, 3rd order:
139 skewness, 4th order: kurtosis) from participants’
140 picture descriptions per 10 ms using openSMILE
141 with the ComParE13 configuration file. ¹⁸ Pitch
142 values in hertz were converted to semitones (st)
143 using individuals’ 10th percentiles to normalize
144 physiological differences among participants ($St =$
145 $\log_2(f_0 / 10^{\text{th}} \text{ percentile}) \times 12$). Several durational
146 measures were computed from SAD timestamps.

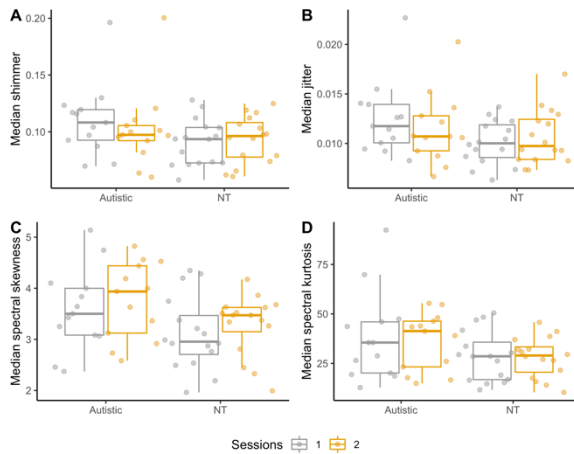
147 2.4 Statistical considerations

148 Preliminary analyses revealed that our variable
149 distributions met the assumptions of parametric
150 tests, so we employed analysis of covariance
151 models. Speech/language features were included as
152 dependent variables, with group, session, and the
153 interaction of group and session as independent
154 variables. Sex was covaried in all models.

155 3 Results

156 3.1 Acoustic measures

157 Median shimmer and jitter values were higher for
 158 autistic children than NT children (shimmer:
 159 $F(1,52)=4.17, p=0.046$; jitter: $F(1,52)=3.96, p=$
 160 0.052 , Figure A-B). Mean, standard deviation
 161 (SD), and interquartile range (IQR) of jitter and
 162 shimmer did not differ by group. Autistic children
 163 also had higher mean (skewness: $F(1,52)=13.46,$
 164 $p<0.001$; kurtosis: $F(1,52)=12.98, p<0.001$),
 165 median (skewness: $F(1,52)=6.17, p=0.016$;
 166 kurtosis: $F(1,52)=4.7, p=0.035$, Figure C-D), SD
 167 (skewness: $F(1,52)=9.89, p=0.003$; kurtosis: $F(1,$
 168 $52)=13.86, p<0.001$), and IQR values (skewness:
 169 $F(1,52)=7, p=0.011$; kurtosis: $F(1,52)=8.26, p=$
 170 0.006) of spectral skewness and kurtosis than NT
 171 children. Groups did not differ in pitch and HNR,
 172 and Session had no significant effect on any
 173 acoustic variables.



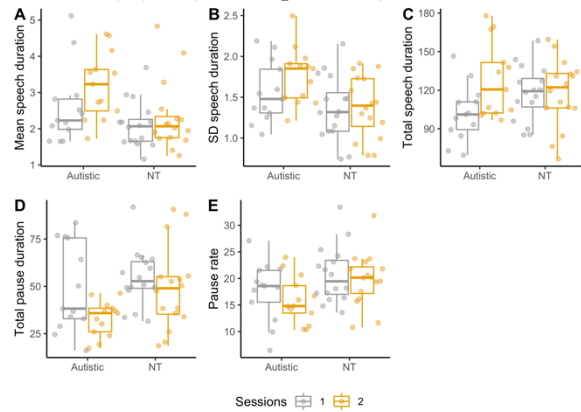
174

175 3.2 Durational measures

176 Autistic children produced longer ($F(1,52)=7.79,$
 177 $p=0.007$) and more variable ($F(1,52)=8.49,$

Figure 1: Acoustic features during picture description tasks.
 178 $p=0.005$) speech segment durations than NT
 179 children (Figure A-B). The difference in total
 180 speech duration between the first and second
 181 sessions was larger for autistic children than NT
 182 children ($F(1,52)=4.34, p=0.042$). Total pause
 183 duration was shorter in autistic participants than
 184 NT children ($F(1,52)=5.14, p=0.028$, Figure C-D),
 185 and children paused longer during the first session
 186 compared to the second ($F(1,52)=4.82, p=0.033$).

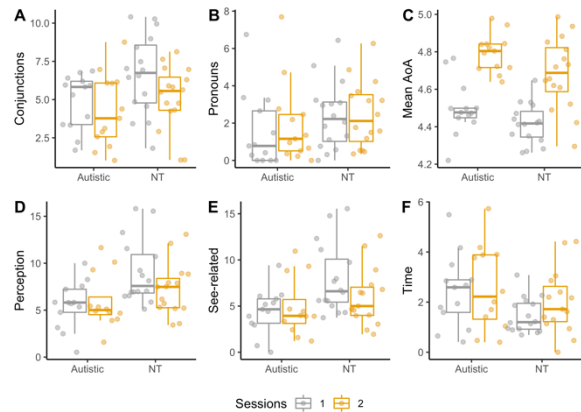
187 Autistic children paused less frequently than NT
 188 children ($F(1,52)=6.33, p=0.015$).



189

Figure 2: Durational measures during picture descriptions. The units of the y-axis are seconds, except the pause rate, where pause rate per minute was plotted.

190 3.3 Textual measures



191

Figure 1: Lexical measures during picture description tasks. All POS counts are per 100 words, and mean AoA was averaged across all content words. LIWC categories were normalized.

192 Autistic participants produced fewer conjunctions
 193 ($F(1,52)=5.06, p=0.029$) and pronouns ($F(1,52)=$
 194 $4.75, p=0.034$) than NT children, and their content
 195 words had a higher AoA than those of NT children
 196 ($F(1,52)=6.35, p=0.015$, Figure 1A-C). Also,
 197 autistic children produced fewer perception ($F(1,$
 198 $52)=9.17, p=0.004$) and see-related words ($F(1,52)$
 199 $=7.1, p=0.01$) and more time-related words ($F(1,$
 200 $52)=4.79, p=0.033$) than NT children (Figure 1).

201 Regardless of diagnostic status, children
 202 produced more adverbs ($F(1,52)=9.08, p=0.003$)
 203 and prepositions ($F(1,52)=6.47, p=0.014$) during
 204 the second session than the first (not shown in the
 205 figure). Children also produced content words that
 206 were more ambiguous ($F(1,52)=10.82, p=0.002$),

207 later acquired ($F(1,52)=54.9$, $p<0.001$), and
208 familiar ($F(1,52)=14.85$, $p<0.001$) during the
209 second session than the first session. Finally,
210 several LIWC categories, including anger ($F(1,52)$
211 $=4.69$, $p=0.035$), difference ($F(1,52)=5.55$, $p=$
212 0.023), feeling ($F(1,52)=4.06$, $p=0.049$), bio ($F(1,$
213 $52)=4.99$, $p=0.03$), and ingestion ($F(1,52)=19$,
214 $p<0.001$), showed significant effects of Session.

215 **4 Discussion**

216 In this study, we elicited picture descriptions from
217 autistic and NT children using a telephony
218 platform, and tested for the presence of diagnostic
219 group differences in a variety of acoustic and
220 lexical features over two sessions. Results showed
221 that autistic children produced greater local jitter,
222 shimmer and the third and fourth orders of spectral
223 moments, as well as shorter and less frequent
224 pauses compared to NT children, across two
225 sessions and with different stimuli. Autistic
226 children produced more speech during the second
227 session when parents administered the task without
228 study staff, compared to the first session, while NT
229 children's speech duration did not differ by session.
230 Lexically, autistic children produced fewer
231 conjunctions and pronouns than NT children, and
232 used later-acquired content words compared to NT
233 peers. Our results also showed that autistic children
234 used fewer see- or perception-related words and
235 more time-related words than NT children.
236 However, many other lexical features differed by
237 session without group differences, suggesting that
238 the picture stimuli may have had more influence
239 than diagnostic group on lexical production.

240 Given that the acoustic features described here
241 remained stable from the first to the second
242 telephony session, and also distinguished the
243 groups, they might hold potential as reliable speech
244 markers of autism. High jitter (variability in pitch)
245 and shimmer (variability in intensity) are perceived
246 as harsh, hoarse, or breathy voice.¹⁹ The
247 observation that autistic children's jitter and jitter
248 variability were higher than NT peers is consistent
249 with prior research that also showed positive
250 correlations between jitter and autism
251 symptomology.¹ However, prior research also
252 found lower HNR values for autistic children
253 compared to NT peers, with no significant
254 differences in shimmer; this differs from our
255 pattern of results. Spectral moments in autism have
256 rarely been studied, even though these measures
257 are known to characterize individuals' voice

258 timbre.²⁰ We plan to study these features further in
259 a larger sample, to explore whether they could
260 serve as validated speech markers of autism.

261 Autistic participants spoke longer and paused
262 less frequently during the second session than the
263 first session, whereas NT children's duration
264 measures did not differ by session. This might be
265 because autistic individuals experience social-
266 communicative challenges which might have
267 hindered their willingness to speak freely in the
268 presence of unfamiliar study staff. In this case, they
269 may have spoken longer in the second session
270 because their parent administered the task. Thus, it
271 is important to consider the presence of study staff
272 when interpreting studies of speech in autism.

273 Finally, our study also found that autistic
274 children produced fewer conjunctions, pronouns,
275 see- and perception-related words with high AoA
276 than NT children. We also observed that many
277 word-level features differed by session in both the
278 autistic and NT groups, suggesting that picture
279 selection has an outsized effect on lexical features.
280 In this study, we selected seven different pictures to
281 prevent boredom and practice effects across
282 multiple sessions. However, since different
283 pictures include unique objects that children are
284 likely to list in their descriptions, this will result in
285 significant session-based differences in word-level
286 features. As data collection continues in the current
287 study, we will investigate whether group
288 differences in more abstract lexical features (e.g.,
289 pronoun use) might remain stable across sessions.

290 **5 Conclusion**

291 Telephony carries great potential as a low-cost and
292 scalable platform for monitoring intervention
293 responses from afar, as well as measuring
294 longitudinal developmental changes in individual
295 children. Acoustic features extracted from data
296 collected using a telephony system, which
297 delivered consistent, high-quality recordings,
298 could be important tools for identifying speech
299 markers of autism.

300 **Acknowledgments**

301 We thank the children and parents who participated
302 in this study. This study was funded by Roche, Ltd
303 (PI: Parish-Morris), and R01DC018289 (PI:
304 Parish-Morris).
305

306 **References**

- 307 Daniel Bone, Chi-Chun Lee, Matthew P. Black, Marian
308 E. Williams, Sungbok Lee, Pat Levitt, and Shrikanth
309 Narayanan. 2014. The Psychologist as an
310 interlocutor in Autism Spectrum Disorder
311 assessment: Insights from a study of spontaneous
312 prosody. *Journal of Speech, Language, and Hearing
313 Research*, 57(4):1162–1177.
- 314 Julia Parish-Morris, Mark Liberman, Neville Ryant,
315 Christopher Cieri, Leila Bateman, Emily Ferguson,
316 and Robert T. Schultz. 2016. Exploring autism
317 spectrum disorders using HLT. In *Proceedings of
318 North American Association of Computational
319 Linguistics, Comp Ling and Clin Psych*, pages 74–
320 84.
- 321 Bonne Yoram, Levanon Yoram, Dean-Pardo Omrit,
322 Lossos Lan, and Adini Yael. 2011. Abnormal
323 Speech Spectrum and Increased Pitch Variability in
324 Young Autistic Children. *Frontiers in Human
325 Neuroscience*, 4.
326 [https://www.frontiersin.org/article/10.3389/fnhum.
327 2010.00237](https://www.frontiersin.org/article/10.3389/fnhum.2010.00237)
- 328 Jaclin Boorse, Meredith Cola, Samantha Plate, Lisa
329 Yankowitz, Juhi Pandey, Robert T. Schultz and Julia
330 Parish-Morris. 2019. Linguistic markers of autism
331 in girls: Evidence of a “blended phenotype” during
332 storytelling. *Molecular Autism*, 10:14.
333 <https://doi.org/10.1186/s13229-019-0268-2>
- 334 Julia Parish-Morris, Mark Liberman, Christopher
335 Cieri, John D. Herrington, Benjamin E. Yerys, Leila
336 Bateman, Joseph Donaher, Emily Ferguson, Juhi
337 Pandey, and Robert T. Schultz. 2017. Linguistic
338 camouflage in girls with autism spectrum disorder.
339 *Molecular Autism*, 8:48.
340 <https://doi.org/10.1186/s13229-017-0164-6>
- 341 Amber Song, Meredith Cola, Samantha Plate, Victoria
342 Petrulla, Lisa Yankowitz, Juhi Pandey, Robert T
343 Schultz, and Julia Parish-Morris. 2021. Natural
344 language markers of social phenotype in girls with
345 autism. *Journal of Child Psychology and
346 Psychiatry*, 62(8): 949-960.
- 347 Mihaela Barokova and Helen Tager-Flusberg. 2020.
348 Commentary: Measuring language change through
349 natural language samples. *Journal of autism and
350 developmental disorders*, 50(7): 2287-2306.
- 351 Michael Rutter, Anthony Bailey, and Catherine Lord.
352 2003. *SCQ: The Social Communication
353 Questionnaire*. Los Angeles: Western Psychological
354 Services.
- 355 John N. Constantino. 2011. *Social Responsiveness
356 Scale, Second Edition*. Los Angeles, CA: Western
357 Psychological Services.
- 358 Dorothy Bishop. 2006. *Children’s Communication
359 Checklist-2 U.S. Edition*. San Antonio, TX:
360 Psychological Corporation.
- 361 Baio J. Prevalence of Autism Spectrum Disorder
362 Among Children Aged 8 Years — Autism and
363 Developmental Disabilities Monitoring Network,
364 11 Sites, United States, 2014. *MMWR Surveill
365 Summ*. 2018;67 Available from: [https://
366 www.cdc.gov/mmwr/volumes/67/ss/ss6706a1.htm](https://www.cdc.gov/mmwr/volumes/67/ss/ss6706a1.htm).
- 367 Matthew Honnibal and Mark Johnson. 2015. An
368 improved non-monotonic transition system for
369 dependency parsing. In *EMNLP 2015: Conference
370 on empirical methods in natural language
371 processing*, pages 1373-1378.
- 372 Mark Brysbaert and Boris New. 2009. Moving beyond
373 Ku cera and Francis: A critical evaluation of current
374 word frequency norms and the introduction of a new
375 and improved word frequency measure for
376 American English. *Behavior Research Methods*,
377 41(4): 977-990.
- 378 Mark Brysbaert, Amy B. Warriner, and Victor
379 Kuperman. 2014. Concreteness ratings for 40
380 thousand generally known English word lemmas.
381 *Behavior Research Methods*, 46(3): 904-911.
- 382 Paul Hoffman, Matthew A. Lambon Ralph, and
383 Timothy Rogers. 2013. Semantic diversity: A
384 measure of semantic ambiguity based on variability
385 in the contextual usage of words. *Behavior Research
386 Methods*, 45(3): 718-730.
- 387 Mark Brysbaert, Paweł Mandera, Samntha F.
388 McCormick, and Emmanuel Keuleers. 2018. Word
389 prevalence norms for 62,000 English lemmas.
390 *Behavior Research Methods*, 51(2): 467-479.
- 391 Florian Eyben, Felix Weninger, Florian Gross, and
392 Björn Schuller. 2013. Recent developments in
393 openSMILE, the Munich Open-Source Multimedia
394 Feature Extractor. In *Proceedings of ACM
395 Multimedia*, pages 835–838.
- 396 Athanasios Tsanas, Max A. Little, Patrick E. McSharry,
397 and Lorraine O. Ramig. 2011. Nonlinear speech
398 analysis algorithms mapped to a standard metric
399 achieve clinically useful quantification of average
400 Parkinson’s disease symptom severity. *Journal of
401 the Royal Society*, 8: 842–855.
- 402 Alexander Lerch. 2012. *An introduction to audio
403 content analysis: Applications in signal processing
404 and music informatics*. John Wiley & Sons.

405 **A Appendix: Inclusion and Exclusion
406 Criteria**

407 Inclusion criteria for participants were the
408 following:

- 409 • Subjects age 6 – 17.99
- 410 • English is participant's first language
- 411 • Verbally fluent – language on grade
412 level/consistent with chronological age
- 413 • Strongly suspected/confirmed diagnosis
414 of autism or typical development
- 415 • Full-scale and verbal IQ > 75
- 416 • For autistic children, current SCQ score
417 ≥ 11
- 418 • For the NT group, current SCQ scores <
419 11

420 Exclusion criteria for participants were the
421 following:

- 422 • Known genetic condition that impacts
423 neurodevelopment or vocal
424 production/language
- 425 • History of persistent language deficits
426 that are currently affecting child's
427 language abilities such that it impacts
428 their ability to have a conversation
- 429 • Extreme prematurity (<32 weeks)
- 430 • History of severe neurological injury
431 likely to affect expressive language and
432 communication behavior
- 433 • If NT, no first-degree family members
434 with autism
- 435 • Plan to begin or change medication
436 during study duration
- 437 • Plan to begin or change an intervention
438 during study duration.
- 439 • Diagnosis of hearing impairment or
440 cochlear implant