# Tipping the Scales: Indiscriminate Use of Interval Scales to Rate Diverse Dysarthric Features

[1]Kaila L. Stipancic, [2]Brooke-Mai Whelan, [3]Lauren Laur, [4]Yunxin Zhao, [5]Andrea Rohl, [6]Inyong Choi, and [6]Mili Kuruvilla-Dugdale

[1]Department of Communicative Sciences and Disorders, University at Buffalo

[2] Division of Speech Pathology, University of Queensland

[3]Department of Speech, Language and Hearing Sciences, University of Missouri

[4]Department of Electrical Engineering and Computer Science, University of Missouri

[5]Department of Neurosurgery, University of Iowa

[6]Department of Communication Sciences and Disorders, University of Iowa

**Author Note**

Running Title:      Tipping the Scales

Correspondence:      Mili Kuruvilla-Dugdale, mkuruvilladugdale@uiowa.edu

Keywords:      Dysarthria, auditory-perceptual scaling, prothetic, metathetic dimensions

**Abstract**

*Purpose:* Error related to incorrect use of rating scales is problematic in the assessment and treatment of dysarthria. The main purpose of this project was to determine scale fit for cardinal speech features of hypokinetic dysarthria. A secondary aim was to determine rater reliability for the two different scales explored.

*Methods:* Forty-three speakers with Parkinson's disease (PD) and 25 neurologically healthy control talkers were recorded reading sentences from the Speech Intelligibility Test. Twenty-two healthy female listeners used both an equal appearing interval (EAI) scale and a direct magnitude estimation (DME) scale to rate five perceptual speech features (i.e., overall speech severity, articulatory imprecision, reduced loudness, short rushes of speech, and monotony) from these recordings. Regression analyses were used to determine the linearity of the relationship between the means of the EAI and DME ratings. Inter- and intrarater reliability was calculated using intraclass correlation coefficients and Spearman's correlation coefficients, respectively, for both EAI and DME ratings.

*Results:* There was a linear relationship between EAI and DME means for monotony, indicating it is a metathetic dimension. Curvilinear relationships were observed between the EAI and DME means for the other four features, indicating prothetic dimensions. Intra- and inter-rater reliability values were similar for EAI and DME ratings.

*Discussion:* Overall, results of this work suggest that DME is the best fit for scaling several hypokinetic dysarthria features, and not the conventionally used EAI scale. Prothetic dimensions best scaled by DME include overall speech severity, articulatory imprecision, reduced loudness, and short rushes of speech. Monotony was the only feature found to be a metathetic dimension and would be best scaled using EAI or DME. Findings call for rethinking the widespread use of EAI scales for rating perceptual features as part of the assessment and treatment of motor speech disorders.

60                                  **Introduction**

61          Auditory-perceptual judgements are the cornerstone of motor speech diagnostics, and

62    are pivotal to tracking speech and voice changes over time (Kent, 1996). Perceptual methods

63    and the speech characteristics they measure vary widely, with a preference toward rating scales

64    for assessing system-level features (e.g., speech intelligibility) over word identification

65    procedures and over subsystem dimensions (e.g., consonant imprecision). For motor speech

66    disorders, subsystem-specific assessment was first introduced through the Mayo Clinic Rating

67    System, which recommended the use of a 5-point scale, to identify clusters of deviant speech

68    features salient to the six dysarthria subtypes (Darley et al., 1969a, 1969b). Over time, the

69    utility and validity of the Mayo Clinic system have been questioned because of its time-

70    intensive nature, limited reliability, and the indiscriminate use of an interval scale applied to

71    all speech features (Ziegler et al., 2017). Notwithstanding the limitations of scale fit (i.e.,

72    construct validity) and reliability, a simplified feature-based system in conjunction with

73    system-level measures provides a comprehensive account of the motor speech impairment and

74    delineates the subsystem impairments contributing to global speech impairment severity and

75    intelligibility loss.

76          In terms of validity, not all perceptual scales are considered equal. The three primary

77    scaling methods include: i) interval scaling, which employs a definitive set of categories or

78    numbers to assign to stimuli (e.g., equal appearing interval [EAI] and visual analog scaling

79    [VAS]); ii) confusion scaling, which requires determination of difference thresholds (e.g.,

80    paired comparison ratings); and iii) magnitude or ratio scaling, which requires the assignment

81    of numerical values proportional to the perceived ratio of the target features in a reference

82    stimulus (e.g., direct magnitude estimation [DME]) (Stevens, 1975; Stevens & Galanter, 1957).

83    Although EAI and VAS are amongst the most widely applied methods to evaluate speech and

84    vocal characteristics (Kreiman et al., 1993), research indicates that EAI scaling is not

85    appropriate for some perceptual dimensions, given evidence of listener bias toward subdividing

86    the lower end of the scale into smaller intervals (Stevens, 1975). EAI scales also offer a limited

87    set of categories to capture the listener's perception and listeners tend to use an equal amount

88    of all the intervals when completing ratings (Zraick & Liss, 2000). For VAS, the level of

89    measurement is unclear. Some researchers have reported that VAS provides only ordinal data

90    (Kersten et al., 2012), while others suggest that VAS behaves like EAI and provides interval

91    data (e.g., Reips & Funke, 2008). Yet other groups deem some types of VASs as ratio or

92    magnitude scales because sensation ratios were in quantitative agreement with VAS ratings of

93    sensory intensity (e.g., Price et al., 1983).

94          Scale validity (i.e., construct validity or how well a scale measures the dysarthric speech

95    feature it is assigned to evaluate) is largely contingent upon the continuum class of a feature

96    (Whitehill et al., 2002), i.e., whether a continuum is prothetic or metathetic. Simply put, a

97    metathetically scaled dimension is a substitutive characteristic, or one that may be altered in

98    terms of quality (e.g., pitch), as opposed to quantity. In contrast, prothetically scaled

99    dimensions refer to additive characteristics that can be measured in relation to incremental

100   changes in quantity or magnitude (e.g., loudness). In terms of the neurobiological difference,

101   prothetic continua involve a quantitative receptor mechanism whereby increasing numbers of

102   sensory receptors respond as the stimulus intensity increases. This differs from the metathetic

103   continuum, which requires a substitutive or qualitative receptor mechanism whereby a different

104   population of receptors is activated as the stimulus intensity increases (Ryan, 1971). A prothetic

105   continuum is also distinct in that the magnitude of the psychological response grows as an

106   exponent of the physical stimulus (e.g., loudness). Contrastively, for metathetic dimensions,

107   the psychological response maintains a constant or uniform distance with the magnitude of the

108   physical stimulus (e.g., pitch). A third difference between these continua is evident when

109   plotting the mean scores of interval (e.g., EAI) and ratio (e.g., DME) scales. A linear

110     relationship suggests a metathetic dimension whereby raters indicate interval equivalence on

111     the measurement scales, and a non-linear relationship points to a prothetic dimension (Stevens,

112     1975). Although the distinction between prothetic and metathetic can be made based on the

113     neural responses associated with each continuum (i.e., substitutive versus additive excitation)

114     and/or the relationship between the physical stimulus and perceptual magnitude (e.g., dB values

115     and loudness perception), it is most easily defined using interval and ratio scales, which is the

116     focus of the current study. EAI scaling was chosen as one of the scaling methods (e.g., as

117     opposed to VAS, for example), as was DME, to follow previous studies in this area.

118        The level of measurement (i.e., ordinal, interval, ratio) from an EAI scale can vary

119     based on the type of continuum, the number of scale intervals (e.g., 5-, 7-, or 9-point), and

120     whether descriptors (e.g., mild, moderate, severe) are used alongside scale values. In general,

121     when listeners' sensitivity to differences is not constant over the scale, the equal interval nature

122     of the data is questioned, and it is recommended that ratings from such a scale be treated as

123     ordinal (Patel et al., 2008).  Only one study has compared three methods for obtaining

124     perceptual judgments based on the premise that equal perceptual distances between the

125     methods suggests equivalent levels of measurement (Patel et al., 2010). The authors found that

126     breathiness ratings obtained from a 7-point scale and a matching task (determined to provide

127     ratio-level measurement) were best fit using a linear function for majority of the talkers, which

128     suggested that the rating scale also has ratio-level measurement properties. This finding

129     conflicts with the common notion that only interval level data is possible with an EAI scale, or

130     that data from an EAI scale should be considered ordinal if listeners cannot apply each interval

131     number as perceptually equidistant from its neighboring intervals (i.e., the perceptual

132     difference between samples rated 4 and 3 should be the same as those rated 4 and 5). Similarly,

133     numerical values combined with descriptors at scale intervals are assumed to provide ordinal-

134     level measurement, where rank order, rather than the magnitude of difference between

135    intervals, is measured. When a scale has fewer, less distinguishable intervals, rank ordering is

136    more likely (Siegel, 1956; Stevens, 1969). To provide further clarification, Snijders and Bosker

137    (2012) stated that variables measured using scales with five to ten intervals and with normal

138    distributions can be considered as continuous rather than categorical variables obtained with

139    an ordinal scale. The current study followed previous studies aimed at establishing

140    psychophysical continua for voice, fluency, and resonance features, by using a 5-point interval

141    scale with descriptors, and as done in these studies, we acknowledge that it cannot be assumed

142    the intervals on an EAI scale represent equally distanced perceptual points.

143        Employing Stevens' (1975) approach, several perceptual speech and voice features

144    have been mapped to metathetic and prothetic scaling dimensions. For example, non-linear

145    relationships and thus, prothetic attribution have been reported for the following perceptual

146    features: overall severity in dysphonic (Eadie & Doyle, 2002b) and tracheoesophageal speech

147    (Eadie & Doyle, 2002a), and hypernasality in speakers with repaired cleft palate (Whitehill et

148    al., 2002; Zraick & Liss, 2000). Prothetic dimensions are best scaled via DME methods; EAI

149    is a poor fit for prothetic dimensions because listeners cannot keep the intervals perceptually

150    equal as they assign stimuli to the various intervals (Stevens, 1971). Conversely, breathiness

151    in dysphonic (Yiu & Ng, 2004) and normal speakers (Sewall et al., 1999), naturalness in

152    stuttered speech (Metz et al., 1990) and tracheoesophageal speech (Eadie & Doyle, 2002a), and

153    pleasantness in dysphonic speakers (Eadie & Doyle, 2002b), have been identified as metathetic.

154    For these dimensions, either DME or EAI scaling is appropriate because listeners' sensitivity

155    to differences is constant over the EAI scale for metathetic continua (Stevens, 1971). Voice

156    disorders have received the most attention in this area and there is a significant paucity of

157    studies on scale fit for measuring dimensions of dysarthria.

158        Error related to incorrect scale use is particularly problematic in the assessment of

159    dysarthrias that progress quickly. As an example, when tracking dysarthria progression, if a

160 change in score from 2 to 3 is more significant than a change from 4 to 5 because listeners tend

161 to divide the lower end of the scale into smaller intervals, it can lead to underestimation of

162 functional loss in the most active phase of the disease when speech changes are accelerated.

163 Moreover, global, system-level measures of dysarthria severity, such as intelligibility and

164 overall speech severity, are very popular in both research and clinical settings and are thought

165 to serve as meaningful indicators of speech loss over time. In speakers with dysphonia, overall

166 severity has been found to be a prothetic dimension best evaluated with DME; however, scale

167 validity for overall severity has not yet been established in speakers with dysarthria. From a

168 treatment standpoint, subsystem features, such as reduced loudness or short rushes of speech

169 often observed in patients with Parkinson's disease (PD), are the targets of therapy aimed at

170 improving speech production. If these are prothetic dimensions, using an EAI scale to quantify

171 intervention outcomes is inappropriate because the change in scores pre- and post-intervention

172 will be difficult to interpret because a change on the lower end is weighted differently than a

173 change on the upper end of the scale. Lastly, relationships between perceptual and acoustic or

174 perceptual and physiologic dimensions can be misinterpreted because invalid scaling options

175 may have been employed. For example, one study showed that the perception of breathiness

176 rated with an EAI scale was not supported by acoustic data (i.e., non-significant noise to

177 harmonic ratio) across both early and late state PD (Holmes et al., 2000) – scale misfit may be

178 one reason, among others, for this discrepancy.

179       Reliability of raters has also been an area of concern in the perceptual evaluation of

180 dysarthric speech. A study by Bunton et al. (2007) examined both intra- and inter-rater

181 reliability of EAI for the 38 perceptual dimensions in the Mayo Clinic Rating System across

182 inexperienced and experienced listeners. In their study, speech samples from 47 speakers with

183 dysarthria were rated by 20 listeners (i.e., 10 inexperienced listeners and 10 clinicians) on the

184 38 perceptual features using a 7-point EAI scale. Percentage of exact score agreements across

185 the features ranged from 32.78% (for pitch level) to 100% (for grunt at the end of expiration).

186 The authors also found no difference in agreement between inexperienced and experienced

187 listeners. The wide range of agreement suggests that although listeners are highly reliable for

188 some features, listeners may have more difficulty (or are more variable) with other perceptual

189 features. Even rating scales that have reliability that is upwards of 50-70% agreement may be

190 insufficient for clinical use, although there is not an agreed-upon standard for this. The authors

191 also noted that more research is needed to determine how to improve reliability of these ratings.

192 Eadie and Doyle (2002b) reported good reliability (i.e., ICCs ranging from 0.692-0.984) for

193 *both* DME and EAI judgments of voice pleasantness and overall severity. These ratings were

194 completed by 12 speech-language pathology (SLP) graduate students who listened to 24

195 speakers with dysphonia and 24 control speakers. Although the authors did not comment on a

196 comparison of these reliability statistics in the discussion, it appears that reliability was grossly

197 similar for the DME ratings and the EAI ratings for the features rated in their study. From these

198 previous studies, it is clear that the reliability of ratings must be considered when attempting

199 to characterize the psychometric properties of rating scales.

200 **Current Study Aims**

201 The first aim of the study was to determine scale fit for the cardinal speech features of

202 hypokinetic dysarthria, namely consonant imprecision, reduced loudness, short rushes of

203 speech, monotony, and overall severity with non-expert listeners as the judges. The focus was

204 on PD because unlike more rapidly progressing dysarthria types (e.g., spastic-flaccid dysarthria

205 due to amyotrophic lateral sclerosis) where there is overreliance on system-level dimensions,

206 both system and subsystem-specific features are frequently used to assess and manage

207 hypokinetic dysarthria. The salient features were selected to cover the speech subsystems

208 predominantly involved in PD, i.e., reduced loudness (phonatory/respiratory subsystems),

209 consonant imprecision (articulatory subsystem), short rushes of speech (articulatory/prosodic

210 subsystems), and monotony (phonatory/respiratory/prosodic subsystems), in addition to overall

211 speech impairment severity (several speech subsystems contribute to this dimension).

212       Based on the existing voice literature, we expected overall speech impairment severity

213 to be prothetic. Because there are no existing studies on the remaining dysarthria features, we

214 could not generate hypotheses for articulatory imprecision, reduced loudness, short rushes of

215 speech, and monotony. Although the psychophysical scaling literature shows that both

216 loudness and duration are prothetic continua (Stevens & Galanter, 1957), it is unclear if

217 derivatives of these features that are considered pathological like reduced loudness and short

218 rushes of speech are also prothetic in nature. Additionally, as mentioned earlier, the continuum

219 class of various features can only be determined in these specific ways: (i) the relationship

220 between scaled interval and ratio scores, (ii) the physical stimulus-perceptual response

221 relationship, and (iii) the neural mechanisms that subserve each continuum. Therefore, it is

222 challenging to hypothesize which dysarthria dimensions are prothetic versus metathetic and

223 which features would be better suited to EAI versus DME measurement solely from perceptual

224 impressions or knowledge about the nature of each feature.

225       Only inexperienced listeners were included in this study because they are more

226 representative of the general population and tend to weight perceptual characteristics

227 differently than highly trained experts. Further, untrained listeners play a role in evaluating

228 dysarthric speech, both in clinical and research settings, when familiarity bias is likely to

229 impede speech evaluation, and when the real-world significance of communication limitations

230 needs to be determined (Lehner, 2021). It is not uncommon for nonexperts to complete fine-

231 grained auditory analysis, for example, master's students in communication sciences and

232 disorders, who are largely untrained when they enter research and clinical settings yet perform

233 these types of ratings. Additionally, from a research perspective, inexperienced listeners are

234 commonly employed to provide ratings of dysarthric speech. Examining how inexperienced

235 listeners use rating scales is a necessary endeavor to support future research efforts that will

236 continue to recruit nonexpert listeners. In psychophysics, most of the seminal studies that

237 established the principles of optimal scaling used nonexpert listeners; more recent studies have

238 also followed previous work in this area by using similar listeners.

239 The second aim of the study was to determine inter- and intrarater reliability for each

240 scale and feature. Based on prior voice findings, for prothetic features, EAI scaling was

241 expected to show lower inter- and intrarater reliability compared to DME. For metathetic

242 features, the inter- and intra-rater reliability was expected to be similar for EAI and DME.

## Methods

244 The study was conducted and received approval from the Institutional Review Board at

245 the University of Missouri. A consent script was read to all the listeners, and written informed

246 consent was received from the speakers. Both listeners and speakers were compensated for

247 participating in the study.

**Participants**

*Speakers*

250 Forty-three individuals with Parkinson's disease (18 females, 25 males) and 25

251 neurotypical controls (11 females, 14 males) provided speech samples for the study. The mean

252 age of the PD group was 68.14 years ($SD = 7.65$), and of the control group was 70.32 years

253 ($SD = 8.43$). All participants met the following criteria: (i) negative history of speech, language,

254 and hearing disorders (except those related to the diagnosis for the PD group); (ii) absence of

255 a co-existing neurological diagnosis for the participants with PD, and any neurological

256 diagnosis for the controls; (iii) negative history of head and neck surgery; (iv) no hearing aids

257 or a prescription for hearing aids; and (v) monolingual, native speakers of American English.

258 The participants with PD displayed a range of dysarthria severities with the majority showing

259 mild dysarthria ($n = 14$), followed by moderate dysarthria ($n = 13$), severe dysarthria ($n = 4$),

260    and lastly, profound dysarthria ($n$ = 3), as determined by the clinical impressions of two

261    experienced SLPs. Nine PD participants displayed typical speech.

262    *Listeners*

263        Twenty non-expert neurologically healthy females rated the speech samples provided

264    by the PD and control groups using EAI and DME scales. The mean age of the listeners was

265    23.91 years (*SD* = 4.39). All listeners passed a bilateral hearing screening at 25 dB HL at

266    500Hz, 1 kHz, 2 kHz, and 4 kHz. They also met all the same inclusionary criteria as the

267    speakers and had minimal exposure to communication disorders (i.e., they had not worked with

268    clinical populations or received any formal instruction in communication disorders).

269    **Experimental Task**

270    *Speakers*

271        The speakers were asked to read aloud sentences from the Speech Intelligibility Test

272    (SIT; Yorkston et al., 2007) at their typical rate and loudness. A print version of the SIT

273    sentences was presented to each speaker. Font size was checked, and readability was confirmed

274    before commencing the recording. The 11 sentences presented to each speaker were randomly

275    generated by the SIT software such that majority (if not all) of the sentences differed across

276    speakers. The speakers read each sentence aloud in their typical rate and loudness after the

277    experimenter announced the sentence number. The speech samples were recorded in a quiet

278    laboratory setting as part of other studies conducted in the last author's lab. A condenser

279    microphone (Shure, Model PG42, Niles, IL) placed 20cm away from the mouth was used to

280    record audio at a sampling rate of 22kHz and the samples were stored on a digital recorder

281    (Marantz, Model PDM670, Eindhoven, Netherlands). Of the 11 SIT sentences, only one

282    sentence per speaker was included in the listening task. The length of the selected sentence

283    varied from 12 to 15 words across the speakers. For each speaker, we selected the sentence

284    with the highest number of hypokinetic speech features represented (i.e., reduced loudness,

285     consonant imprecision, short rushes of speech, and monotony). To determine which features

286     were present in each sample, two trained research assistants, who had completed the graduate

287     Motor Speech Disorders course, used the Dysarthria Rating Scale (Darley et al., 1969 a & b)

288     to make independent judgments. They arrived at a consensus for features with divergent

289     ratings, and the consensus ratings were used for sample selection.

290     ***Listeners***

291        The listeners performed the perceptual ratings over two sessions that were about one

292     week apart and lasted approximately an hour each. Either a 5-point EAI scale or modulus DME

293     scale was used in each session. Across listeners, the order of the scales and speech samples was

294     randomized. Listeners rated 82 samples (68 samples plus 14 re-rated for calculation of intra-

295     rater reliability) in each session. Among the 68 samples, none of the sentences were repeated.

296     Before commencing each task, the experimenter (last author) provided definitions for each of

297     the five features to be rated, followed by instructions on appropriate scale use. For example,

298     for overall severity, listeners were asked to rate the samples based on their general impression

299     of speech impairment severity and not on understandability. For reduced loudness, they were

300     instructed to rate the extent that the voice was insufficiently loud. For consonant imprecision,

301     listeners were asked to determine if some or most sounds were produced crisply and sharply,

302     and for short rushes they were asked to listen for rapid bursts of speech separated by pauses.

303     For monotony, listeners were told to rate the extent that the sample sounded flat in terms of

304     pitch, loudness, or duration. Regarding scale use, listeners were instructed to use the full range

305     of the scale. Listeners were strongly encouraged to rate the first three features after listening to

306     the sample once and rate the next two features after listening to the sample a second time to

307     avoid confusing the features to be rated. The features were rated in the following order across

308     all samples for both EAI and DME: overall speech impairment severity, articulatory

309     imprecision, reduced loudness, short rushes of speech, and monotony. Listeners were only

310    allowed to listen to each sample twice. The 5-point EAI scale had the following intervals:

311    1=typical, 2=mild, 3=moderate, 4=severe, and 5=profound. A 5-point scale was used to be

312    consistent with the procedures of the Mayo Clinic Rating System from where the features were

313    taken. The severity descriptors used at the intervals of the scale also follow the Mayo Clinic

314    Rating System.

315         For DME, the listeners were instructed to first listen to the moduli (or references) for

316    the first three features. Each modulus represented a moderate level of severity and was given a

317    score of 100; each of the features had a different modulus. After listening to the moduli, the

318    listeners played the sample and provided a comparative score between the modulus and the

319    sample. Specifically, listeners were told that if they perceived the sample to be twice as severe

320    as the modulus, a score of 200 would be appropriate. Similarly, if they perceived the sample to

321    be half as severe as the modulus, a score of 50 would be appropriate. The lower limit was fixed

322    at 1 to carry out geometric mean calculations; no upper limit was specified. After entering their

323    scores for the first three features, listeners proceeded to listen to the last two moduli before

324    listening to the sample a second time and scoring the last two features. Listeners were required

325    to relisten to the moduli after every 6[th] sample to maintain referential value of the modulus and

326    to minimize shifting of the listeners' internal standards when performing the ratings (Eadie &

327    Doyle, 2002; Kreiman et al., 1993).

328         The DME moduli were initially selected by an experienced SLP who listened to

329    dysarthria samples from the *Audio Seminar Series* (Darley et al., 1975) and chose samples for

330    each feature representing mild, moderate, severe, and profound severity. Once these samples

331    were identified for each feature, the last author rated them independently for severity, as well

332    as the feature(s) represented in each sample. Consensus was sought when the two experts

333    disagreed about feature representation or severity level; both experts confirmed that the

334    modulus for each feature represented moderate severity.

335 **Data Acquisition**

336          All audio samples were presented through headphones (Panasonic, RP-HC200-K and

337     RP-DJS150) from a desktop computer (Dell, OptiPlex 7010) in a quiet laboratory setting. The

338     sound files recorded on different dates had different intensity levels. Therefore, the intensity

339     levels were normalized across all the sound files using the Yoon script (Yoon, 2022) in Praat

340     (Boersma & Weenink, 2021). When the script is run, the sound files are scanned to identify the

341     maximum and minimum intensity in dB. The sound files are then normalized to 65 dB. Because

342     two different computers and headsets were used, the sound pressure level (SPL) from each

343     headphone was calibrated with a half-inch precision condenser microphone (PCB Piezotronics,

344     377C13) and a sound level meter (Larson Davis, SoundTrack LxT®), by connecting the

345     headphone to an ear simulator (Larson Davis, AEC201). A one-minute speech sample devoid

346     of long pauses or silences was played through the headphone connected to the artificial ear,

347     and the SPL, via the sound card interface, was adjusted until the sound level meter showed

348     75dB SPL. The volume level corresponding to 75 dB SPL was noted for each computer-headset

349     pair and checked before each session to ensure that all stimuli were played at the same loudness

350     level across listener participants.

351          Five EAI scales clearly demarcated for each feature were available to the listeners via

352     a custom-built MATLAB GUI. Listeners were instructed to use the correct scale to judge each

353     dimension from 1 (i.e., no impairment) to 5 (i.e., profound impairment). A second custom-built

354     MATLAB GUI was used for the DME ratings. The moduli were embedded at the top of the

355     screen and five textboxes were available at the bottom of the screen to enter the comparative

356     scores for the features. Renderings of the GUIs provided to listener participants are displayed

357     in Figure 1. For both scales, listeners could only advance to the next sample after all five scores

358     were entered. Moreover, for DME, listeners could not proceed if they missed a modulus or

359     attempted to listen to a modulus more than once or when the modulus was inactive.

360 **Figure 1.** Renderings of the MATLAB GUIs used by listeners for (A) Equal Appearing Interval
361 (EAI) Scales; (B) Direct Magnitude Estimation (DME).
362



363
364

365 **Data Analysis**

366      EAI and DME scores were generated automatically as part of the MATLAB output file.

367 DME and EAI ratings were compared using Stevens' (1975) method. Geometric rather than

368 arithmetic means are required for DME because the relationship between scores is

369 multiplicative or exponential. Contrastively, using arithmetic means is appropriate for EAI

370     because the relationship between scores is additive. To compute the arithmetic means from

371     EAI scores, the scores were averaged across all 20 listeners for each feature and speaker. For

372     DME, the geometric mean was calculated across all the listeners for each feature and speaker

373     using the formula $\sqrt[n]{X_1 X_2 \ldots X_n}$ (where $X$ = values of DME ratings, $n$ = number of

374     speakers). The EAI arithmetic means were plotted against the DME geometric means for all

375     samples (excluding the samples repeated for intrarater reliability). Following Stevens' (1975)

376     interpretation, a linear relationship between the two sets of means indicated a metathetic

377     continuum, whereas a downward bowed curvilinear function indicated a prothetic continuum.

378     **Statistical Analysis**

379          To determine scale fit for each of the five features the following steps were taken in

380     SPSS version 28 (IBM SPSS, Armonk, NY):

381     ***Linear Regression***

382          A linear regression was performed for each feature to determine the linearity of the

383     relationship between the means of the EAI and DME ratings. The geometric mean was placed

384     in block 1, the computed square of the geometric mean in block 2, and the computed cube of

385     the geometric mean in block 3. The arithmetic mean was the dependent variable in each model.

386     The model summary, specifically the adjusted $R^2$ values and significant $F$ change were checked

387     for each feature to determine fit. Significance was set at $p < .05$. Cook's distance was examined

388     to identify outliers; for all features, data were removed from one participant with PD due to

389     Cook's distance being greater than 1.0. For the same reason, data from a second PD participant

390     was also removed for all features except monotony. For reduced loudness, visual inspection of

391     the predicted values versus residual plots revealed heteroscedasticity; therefore, the arithmetic

392     means of the EAI scores were log transformed. Linear regressions were performed again after

393     removing the outliers and correcting heteroscedasticity. Visual inspection of the residuals

394     showed no signs of violating the normality assumption. Additionally, the Shapiro Wilk test

15

395   was performed and showed normal distribution of residuals for overall severity ($W(61) = 0.97$,

396   $p > 0.05$), articulatory imprecision ($W(61) = 0.96$, $p > 0.05$), reduced loudness ($W(61) = 0.98$,

397   $p > 0.05$), short rushes ($W(61) = 0.97$, $p > 0.05$), and monotony ($W(61) = 0.99$, $p > 0.05$).

398   ***Collinearity Diagnostics***

399           This step helped determine whether collinearity affected the fit of the linear, quadratic,

400   and cubic models. When the variance inflation factor was above 10, the condition index was

401   inspected for values greater than 30, and variance proportions greater than 0.90 among two of

402   the predictors. For all five features, only the cubic model met these three criteria; therefore, the

403   cubic model was excluded for all features.

404           Linear regression assumes interval or ratio-level data, and taking an average of scaled

405   scores is regarded as requiring an interval scale. Because none of the final models were found

406   to violate regression assumptions of normality, collinearity, homoscedasticity, and outliers, the

407   regression results are considered valid.

408   ***Rater Reliability***

409           Reliability between raters (i.e., interrater reliability) was estimated using intraclass

410   correlation coefficients (ICCs) for each scaling method. ICCs and their 95% confidence

411   intervals (CIs) were calculated using SPSS statistical package version 28 (SPSS Inc., Chicago,

412   IL) based on average-measures consistency, 2-way mixed-effects model with 20 listeners

413   across 68 samples. Reliability within raters (i.e., intrarater reliability) was judged using

414   Spearman's correlation coefficients.

415                                    **Results**

416           EAI arithmetic means plotted as a function of the DME geometric means revealed a

417   statistically significant result for a second-order polynomial curve of best fit for overall severity

418   ($F(2, 64) = 432.47$, $p < .001$, $R^2_{\text{Adjusted}} = .93$), articulatory imprecision ($F(2, 65) = 241.61$, $p <$

419   $.001$, $R^2_{\text{Adjusted}} = .88$), reduced loudness ($F(2, 65) = 85.98$, $p < .001$, $R^2_{\text{Adjusted}} = .72$), and short
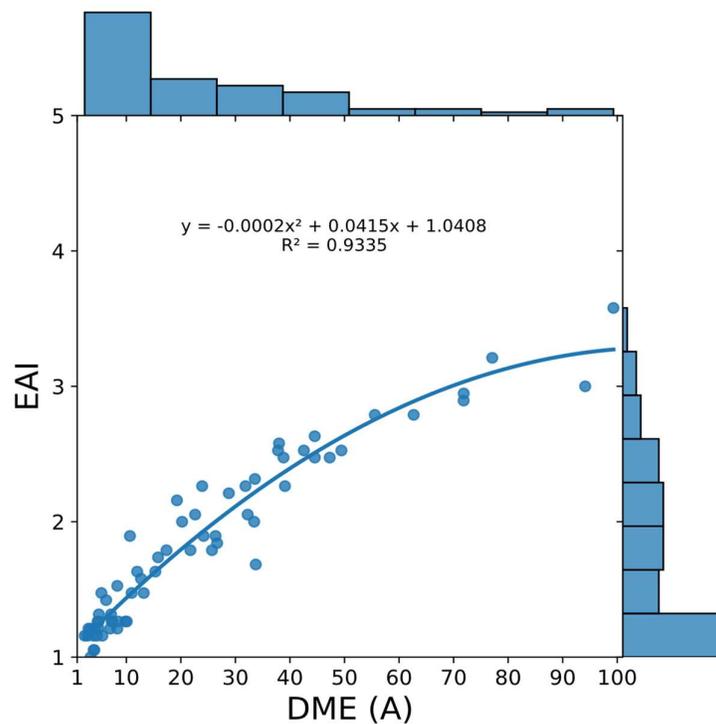
420    rushes of speech ($F(2, 65) = 197.50$, $p < .001$, $R^2_{Adjusted} = .86$). This indicates that the curvilinear

421    model accounted for a statistically significant amount of the variance above that observed with

422    a simple linear model. As shown in Figures 2 a-d, visual inspection of the model revealed a

423    downward bowing towards the end of the curve, a finding that agrees with other data for

424    prothetic continua reported in the literature (Eadie & Doyle, 2002 a & b; Zraick and Liss, 2000).

425        In contrast, a linear model accounted for a statistically significant amount of variance

426    for monotony ($F(1, 66) = 376.26$, $p < .001$, $R^2_{Adjusted} = .85$). The quadratic and cubic curvilinear

427    models revealed no significant improvement in the variance over the linear model. Visual

428    inspection of the data also revealed a good approximation to the raw data by the linear

429    regression, with no apparent downward bowing (see Figure 2 e).
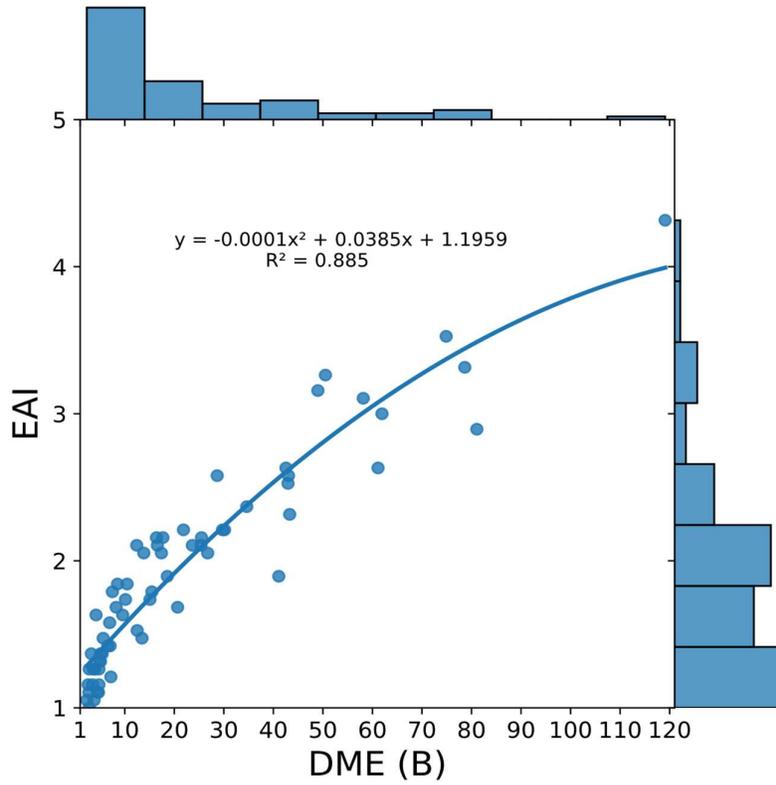
430

431    **Figure 2.** Arithmetic means of equal appearing interval scale (EAI) scores plotted against the
432    geometric means of the direct magnitude estimation (DME) scores for (A) overall severity, (B)
433    articulatory imprecision, (C) reduced loudness, (D) short rushes of speech, and (E) monotony.
434    The arithmetic means for reduced loudness were log transformed. The histograms based on the
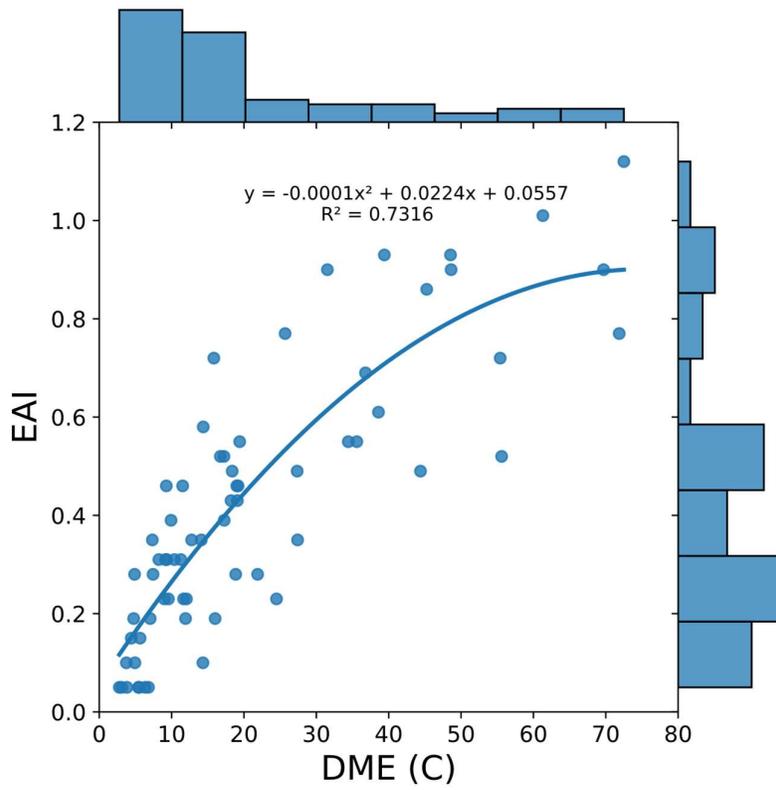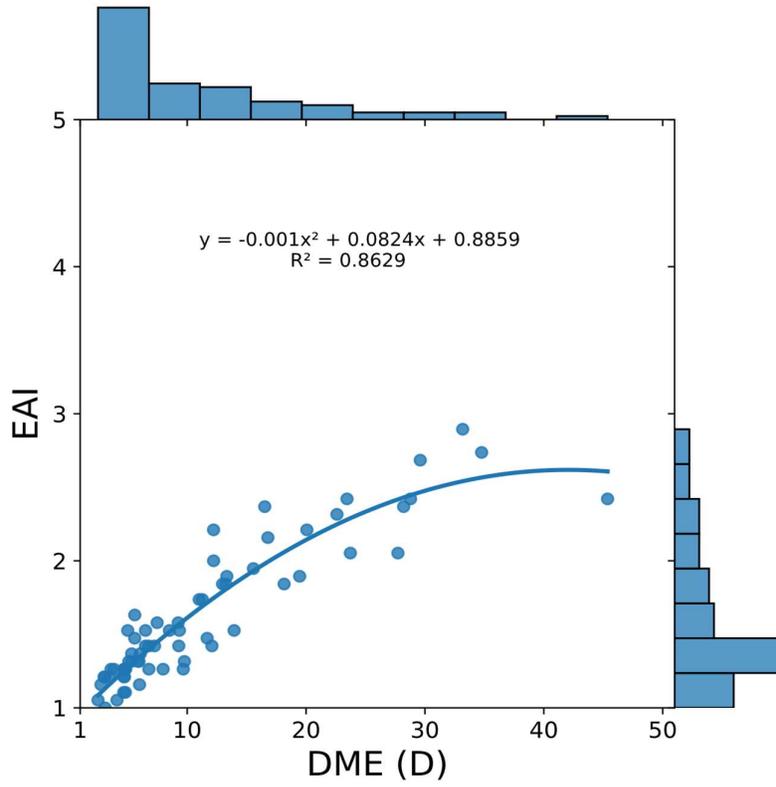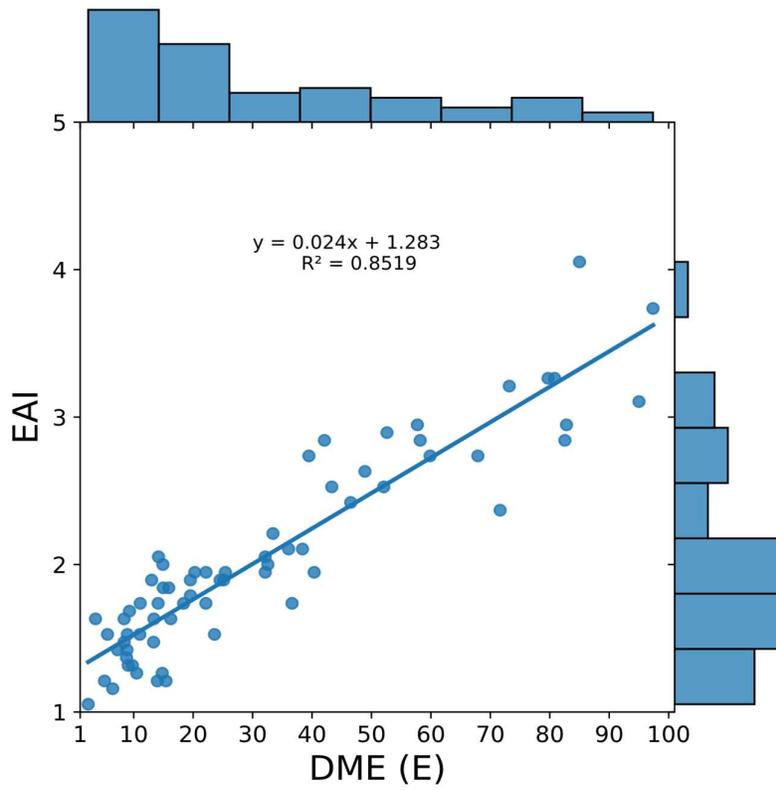435    means of the EAI and DME scores are displayed opposite each respective scale axis.



436

$y = -0.0001x^2 + 0.0385x + 1.1959$
$R^2 = 0.885$

437



$y = -0.0001x^2 + 0.0224x + 0.0557$
$R^2 = 0.7316$

438

y = -0.001x² + 0.0824x + 0.8859
R² = 0.8629

439



y = 0.024x + 1.283
R² = 0.8519

440
441
442

443 **Rater Reliability**

444     Table 1 displays the interrater and intrarater reliability results. Interrater reliability

445 coefficients were highly acceptable (i.e., 0.89 or greater, indicating good or excellent

446 reliability) for all features and both scaling methods (Koo & Li, 2016). Intrarater reliability was

447 moderate (0.64-0.69) to strong (0.71-0.79) for both methods (Dancey & Reidy, 2004).

448

449 **Table 1.** Correlation coefficients for inter- and intrarater reliability of equal appearing interval
450 (EAI) and direct magnitude estimation (DME) ratings for all features.

| Dimensions | Inter-Rater Reliability Intraclass Correlation Coefficients (95% Confidence Intervals) | | Intra-Rater Reliability Spearman's Correlation Coefficient (95% Confidence Intervals) | |
|---|---|---|---|---|
| | EAI | DME | EAI | DME |
| Overall severity | 0.96 (.949-.975) | 0.93 (.895-.948) | 0.75 (.695-.803) | 0.76 (.697-.804) |
| Articulatory imprecision | 0.94 (.929-.975) | 0.93 (.898-.949) | 0.79 (.734-.829) | 0.69 (.621-.751) |
| Reduced loudness | 0.91 (.872-.937) | 0.89 (.846-.924) | 0.65 (.566-.712) | 0.65 (.568-.713) |
| Short rushes of speech | 0.93 (.895-.948) | 0.92 (.884-.942) | 0.71 (.638-.763) | 0.72 (.651-.772) |
| Monotony | 0.92 (.895-.948) | 0.90 (.860-.930) | 0.65 (.574-.718) | 0.64 (.560-.707) |

451

452

453                                            **Discussion**

454     The current research note sought to determine scale fit for several speech features

455 whose measurements are critical to the evaluation of hypokinetic dysarthria associated with

456 PD: overall speech severity, articulatory imprecision, reduced loudness, short rushes of speech,

457 and monotony. Specifically, we were interested in whether the commonly used EAI scale, or

458 the less popular DME scale would yield better validity and reliability. Schiavetti et al. (1981)

459 referred to determining whether a continuum is prothetic or metathetic as an aspect of construct

460 validity. As a reminder, three ways to determine continua class are by examining the (i)

461 relationship between interval and ratio scales, (ii) relationship between physical and perceptual

462 magnitude, and (iii) neural activation patterns. In this study, we compared the interval-ratio

463 scale relationship, and our findings indicate that DME scaling has greater construct validity

464 than EAI scaling for most of the features examined, except monotony. Our first hypothesis that

465    overall speech impairment severity would be a prothetic continua was supported by these

466    results.

467         Our hypotheses for reliability were not entirely supported because reliability was

468    comparable for both EAI and DME ratings for all features. We hypothesized that for prothetic

469    continua, like overall speech impairment, articulatory imprecision, reduced loudness, and short

470    rushes of speech, reliability would be lower for EAI than for DME. The reliability hypothesis

471    for metathetic continua was supported by the data for monotony ratings – the hypothesis was

472    that for metathetic continua, rater reliability would be similar for the two rating tasks. Herein,

473    we discuss these findings along with implications for clinical settings and research endeavors,

474    as well as directions for future research in this area.

475         Early work describing methods for evaluating scale construction highlighted the

476    importance of considering both reliability and validity of a scale to determine its usefulness,

477    even though these concepts are complex and challenging to fully evaluate (Dawis, 1987).

478    Therefore, as a brief, initial discussion point about the current study, inter-rater reliability was

479    excellent for both scaling methods across all five features (ICCs ranged from 0.89 to 0.96). The

480    reliability statistics in the current project are even slightly higher than those reported in recent

481    work from our group for VAS ratings of speakers with PD (Stipancic et al., 2023). This was a

482    somewhat unexpected finding, as rater reliability in scaling tasks has been questioned in the

483    speech perception literature (Miller, 2013; Schiavetti, 1992; Stipancic et al., 2016). One

484    difference among studies is whether single measures or average measures ICC is reported. The

485    former is based on a single measurement from a single observer, and the latter is based on the

486    average measurements of more than one observer. Average measures ICC tends to have higher

487    values because it accounts for the reduction in measurement error achieved by averaging

488    multiple measurements (Trevethan, 2017). Another unexpected finding was the slightly lower

489    intrarater reliability compared to interrater reliability across all the features, given that previous

490     literature has generally demonstrated the opposite (i.e., that listeners are more reliable *within*

491     themselves than *across* other listeners). Because the two types of scaling procedures used in

492     the current project have similar levels of reliability, scale choice should rely on other

493     considerations, such as those discussed in the following sections.

494     **Construct Validity is Better for DME than EAI for Most Dysarthric Features Explored**

495     Non-linear relationships between EAI and DME ratings were observed for four out of

496     five features (i.e., overall speech severity, articulatory imprecision, reduced loudness, and short

497     rushes of speech), indicating that they are prothetic dimensions. This finding suggests that these

498     four dimensions should not be scaled using EAI, which is what the Mayo Clinic Rating System

499     recommends. These findings point to DME scaling having better construct validity than EAI

500     scaling for the dimensions of overall speech severity, articulatory imprecision, reduced

501     loudness, and short rushes of speech. The current results are similar to seminal findings from

502     Schiavetti (1992) on the construct validity of scaling speech intelligibility. According to the

503     author, EAI scaling is not appropriate to measure intelligibility because it cannot be partitioned

504     linearly into equal intervals. Using an EAI scale for intelligibility results in an ordinal rather

505     than interval level of measurement. Similarly, it follows that overall speech severity,

506     articulatory imprecision, reduced loudness, and short rushes of speech, as rated in hypokinetic

507     dysarthria, cannot be partitioned into equal intervals by listeners, making an EAI scale

508     inappropriate to use for ratings.

509     **Scaling Monotony Appears to Differ from the Other Dysarthric Features**

510     In contrast to the other four dimensions, there was a linear relationship between EAI

511     and DME ratings for monotony, indicating that it is a metathetic dimension. Interestingly, in

512     recent work (Stipancic et al., 2023), our group identified ratings of monotony as having the

513     poorest criterion validity and reliability compared to three other features of dysarthria (i.e.,

514     overall speech impairment, articulatory imprecision, and slow rate). Findings of the current

515     study suggest that monotony can be rated using an EAI scale and, therefore, that listeners

516     perceive monotony more linearly than the other dimensions. As a reminder, listeners were not

517     asked to make a distinction between monopitch and monoloudness because it is difficult for

518     non-expert listeners to distinguish these categories (Kim, 1994). In future studies, we propose

519     that construct validity be determined for each sub-feature separately, regardless of the fit of the

520     "root" feature. In this case, for example, ratings of monoloudness, monopitch, and

521     monoduration should all have scale fit independently determined, outside of their contributions

522     to monotony. The relative relationships of each of these features to monotony should also be

523     considered in future work. Overall, taking into account the current study and our previous work

524     (Stipancic et al., 2023), monotony appears to be a challenging dimension for non-expert

525     listeners to rate, whether it be because it is a composite feature (i.e., made up of other features),

526     or because its rating is limited due to psychophysical and neural restrictions in the ability to

527     perceive monotony, or because listeners have a poor working definition of monotony.

528     Improving the rating of monotony will be an interesting area of future research, given its

529     importance as a feature that has long been considered central to the diagnosis of dysarthria

530     subtypes like hypokinetic dysarthria (Darley et al., 1969 a & b).

531     **Clinical and Research Implications**

532          Results of this work suggest that DME is the best fit for scaling several hypokinetic

533     dysarthria features, and not the conventionally used EAI scale. Error related to incorrect scale

534     use has implications for both assessment and treatment tracking of individuals with dysarthria.

535     This is especially true for the feature of overall severity which is ubiquitous in the field both in

536     clinical practice and in research studies (Stipancic et al., 2021). Further, for individuals with

537     PD participating in programs like the Lee Silverman Voice Treatment, multiple loudness

538     ratings are typically carried out during assessment (e.g., baseline and stimulability check) and

539     treatment (e.g., daily for four weeks). Similarly, features tied to speech rate, such as short

540    rushes of speech, are likely to receive additional attention from a clinician during treatment

541    sessions. If an EAI scale is used in these instances, listeners will be unable to place stimuli

542    along the scale in an unbiased manner due to the propensity to divide the lower end of the

543    continuum into finer segments than the upper end (Stevens, 1974; Whitehill et al., 2002). EAI

544    scales, therefore, limit listeners from fully communicating their auditory perception when it

545    comes to prothetic continua. This is thought to result in the ordinal partition of the interval

546    scale, and if used in research studies, different descriptive and inferential statistics are

547    recommended for ordinal versus interval variables to decrease the chance of erroneous

548    conclusions from distorted effect sizes, inflated false alarm rates, etc. According to Stevens

549    (1958), more statistical operations, particularly parametric statistics, are permissible only for

550    higher measurement levels (i.e., ratio and interval) and not for ordinal and nominal levels.

551    Based on this understanding, ordinal scaling disallows the use of metric models for interpreting

552    results; however, this idea is debated (Torrin & Kruschke, 2018).

553         Several authors have noted that although psychometric properties of scales, such as

554    reliability and validity, are crucial to determine, other considerations are equally important

555    (Dawis, 1987). These other factors include administration concerns and usability of the scale

556    in the setting(s) of interest. Even though the current results suggest that DME is more

557    appropriate than EAI for rating most features considered here, DME is known to be

558    cumbersome to use in both clinical and research settings – it takes longer to use, requires

559    preparation to identify moduli, and poses challenges in conveying the meaningfulness of DME

560    scores to clients. VAS is a very popular alternative to both approaches (as used in Stipancic et

561    al., 2023), largely due to ease of implementation and use, but its scale properties are unknown,

562    and we need to establish the validity associated with VAS. Particularly, as related to the current

563    work, it will be important to determine whether VAS functions more like a DME or an EAI

564  scale, and whether construct validity is comparable for VAS and DME since DME can be used

565  with both prothetic and metathetic dimensions.

566       Final thoughts from Dawis in 1987 are still applicable to the current work: "In scale

567  construction, as in much of human endeavor, there can be no single "best" method. One method

568  may be best for one research purpose but not for another. Purpose, context, and limitations on

569  the researcher must be considered. Trade-offs in advantages and disadvantages seem to be the

570  rule…" (Dawis, 1987, p. 488). Therefore, it is important to consider the context in which a

571  measure might be used. If, for example, as the current results suggest, DME has better validity

572  than EAI for overall speech severity, but a practicing clinician is reluctant to use DME because

573  of the inconvenience or the lack of a standard modulus, then using an easier-to-implement scale

574  that might not have the best validity, may be appropriate. In contrast, a researcher who is

575  interested in making very precise, valid measurements to evaluate treatment efficacy, and has

576  the time and resources for a more complex procedure, should consider using DME. Overall,

577  the current findings call for rethinking the widespread use of EAI scales for rating perceptual

578  features as part of assessment and treatment of motor speech disorders.

579  **Limitations and Future Directions**

580       Only construct validity was examined in this project and future work could examine

581  other types of validity along with construct validity to determine key scale properties. Robust

582  analyses of validity *and* reliability (both interrater and intrarater) in the same study would be

583  beneficial for determining psychometric properties of similar types of scales since validity and

584  reliability are separate concepts and do not necessarily vary in tandem (i.e., a scale can be

585  reliable, but not valid and vice versa). The limitations of reliability analyses in the current work

586  may have contributed to the unexpected findings and could be addressed in future studies.

587  Specifically, results revealed that intrarater reliability was lower than interrater reliability,

588  when the opposite is typically true (i.e., reliability within a rater, is typically better than

589   reliability across a group of raters). However, the small number of data points used in the

590   reliability analyses, the difference in statistical analyses used for calculating intra vs. interrater

591   reliability, and an inability to determine statistical or clinically meaningful differences between

592   reliability estimates, may have contributed to this unexpected finding. More work examining

593   scale fit, other types of validity, and reliability of the ratings used commonly in research and

594   in clinic is critically needed. For example, only five of the 38 possible dysarthria features from

595   the Mayo Clinic Rating System were included in the current study, and the psychophysical

596   properties of the majority tied to other dysarthria types remain largely unexplored. This study

597   used inexperienced listeners to rate the dysarthric speech features of interest and future work

598   could examine the impact of listener characteristics on validity and reliability of scaling tasks.

599   Given that continuum type is determined by neural activation patterns and how physical

600   stimulus properties relate to perceptual magnitude, it remains to be determined if listener

601   experience will alter how perceptual continua are classified. Lastly, the level of measurement

602   represented by VAS needs to be established. Some voice researchers suggest that it behaves

603   similar to EAI, without the fixed, pre-defined points along the scale (Wuyts, De Bodt, &Van

604   de Heyning, 1999). Other groups suggest that VAS can function as a ratio scale like DME

605   (Price et al., 1983). Since VAS has gained popularity for research and clinical use, this is a

606   critical next step.

607   **Conclusions**

608       The validity of different scaling methods depends upon the continua class of the

609   dimensions being rated. Four of the five cardinal features of hypokinetic dysarthria (i.e., overall

610   severity, articulatory imprecision, reduced loudness, and short rushes of speech) were

611   determined to be prothetic dimensions best scaled using DME and only monotony behaved as

612   a metathetic dimension suited either to EAI or DME scaling. Given the unsuitability of EAI

613   scales for prothetic dimensions, and the cumbersome nature of the optimal alternative, DME,

614     it is recommended that researchers and clinicians consider the purpose and context of use while

615     also weighing the advantages and disadvantages of each of these factors. To this end, inter- and

616     intrarater reliability are comparable between the two scaling methods when rating features of

617     hypokinetic dysarthria included in this study.

### Data Availability Statement

629     Data supporting the results reported in this manuscript are available for interested researchers

630     on request from the authors.

### References

632     Boersma, P., & Weenink, D. (2022): Praat: doing phonetics by computer [Computer program].

633         Version 6.2.06, retrieved 23 January 2022 from https://www.praat.org.

634     Bunton, K., Kent, J., Duffy, J., Rosenbek, J., & Kent, R. (2008). Listener agreement for

635         auditory perceptual ratings of dysarthria. *Journal of Speech, Language and Hearing*

636         *Research, 50*(6), 1481-1495.

637    Chiu, Y.-F., Neel, A., & Loux, T. (2021). Exploring the acoustic perceptual relationship of

638        speech in Parkinson's disease. *Journal of Speech, Language and Hearing Research*, *64*,

639        1560-1570.

640    Darley, F. L., Aronson, A. E., & Brown, J. R. (1969). Clusters of deviant speech dimensions

641        in the dysarthrias. *Journal of Speech and Hearing Research*, *12*, 462-496.

642    Dawis, R. V. (1987). Scale construction. *Journal of Counseling Psychology, 34*(4), 481-489.

643    DeMaagd, G., & Philip, A. (2015). Parkinson's disease and its management. *Pharmacy &*

644        *Therapeutics*, *40*(8), 504-510.

645    Duffy, J. R. (2013). *Defining, understanding, and categorizing motor speech disorders.*

646        Elsevier Mosby.

647    Eadie, T. L., & Doyle, P. C. (2002a). Direct magnitude estimation and interval scaling of

648        naturalness and severity in tracheoesophageal (TE) speakers. *Journal of Speech,*

649        *Language and Hearing Research*, *45*(6), 1088-1096.

650    Eadie, T. L., & Doyle, P. C. (2002b). Direct magnitude estimation and interval scaling of

651        pleasantness and severity in dysphonic and normal speakers. *The Journal of the*

652        *Acoustical Society of America*, *112*. https://doi.org/https://doi.org/10.1121/1.1518983

653    Fereshtehnejad, S.-M., Yao, C., Pelletier, A., Montplaisir, J. Y., Gagnon, J.-F., & Postuma, R.

654        B. (2019). Evolution of prodromal Parkinson's disease and dementia with Lewy bodies:

655        a prospective study. *Brain*, *142*(7), 2051-2067.

656    Ho, A. K., Bradshaw, J. L., Iansek, R., & Alfredson, R. (1999). Speech volume regulation in

657        Parkinson's disease: Effects of implicit cues and explicit instructions.

658        *Neuropsychologia*, *37*(13), 1453-1460.

659     Kempster, G. B., Kistler, D. J., & Hillenbrand, J. (1991). Multidimensional scaling analysis of

660         dysphonia in two speaker groups. *Journal of Speeech and Hearing Research*, *34*, 534-

661         543.

662     Kent, R. D. (1996). Hearing and believing: Some limits to the auditory-perceptual assessment

663         of speech and voice disorders. *American Journal of Speech-Language Pathology*, *5*(3),

664         7-23. https://doi.org/https://doi.org/10.1044/1058-0360.0503.07

665     Kim, H. (1994). *Monotony of speech production in Parkinson's disease: Acoustic*

666         *characteristics and their perceptual relations* (Doctoral dissertation).

667     Kreiman, J., & Gerratt, B. R. (1996). The perceptual structure of pathologic voice quality. *The*

668         *Journal of the Acoustical Society of America*, *100*, 1787-1795.

669     Kreiman, J., Gerratt, B. R., & Berke, G. S. (1994). The multidimensional nature of pathologic

670         vocal quality. *The Journal of the Acoustical Society of America*, *96*(3), 1291-1302.

671     Kreiman, J., Gerratt, B. R., Kempster, G. B., Erman, A., & Berke, G. (1993a). Perceptual

672         evaluation of voice quality: Review, tutorial, and a framework for future research.

673         *Journal of Speech and Hearing Research*, *36*, 21-40.

674     Kreiman, J., Gerratt, B. R., Kempster, G. B., Erman, A., & Berke, G. S. (1993b). Perceptual

675         evaluation of voice quality: Review, tutorial, and a framework for future research.

676         *Journal of Speech and Hearing Research*, *36*, 21-40.

677     Logemann, J. A., Fisher, H. B., Boshes, B., & Blonsky, E. R. (1978). Frequency and co-

678         occurence of vocal tract dysfunctions in the speech of a large sample of Parkinson

679         patients. *Journal of Speech and Hearing Disorders*, *43*(1).

680     Ma, A., Lau, K. K., & Thyagarajan, D. (2020). Voice changes in Parkinson's disease: What

681         are they telling us? *Journal of Clinical Neuroscience*, *72*, 1-7.

682         https://doi.org/https://doi.org/10.1016/j.jocn.2019.12.029

683  Metz, D. E., Schiavetti, N., & Sacco, P. R. (1990). Acoustic and psychosocial dimensions of

684      the perceived speech naturalness of nonstutterers and posttreatment stutterers. *Journal*

685      *of Speech and Hearing Disorders*, *55*, 516-525.

686  Miller, N. (2013). Review: Measuring up to speech intelligibility. *International Journal of*

687      *Language & Communication Disorders, 48*, 601-612.

688  Patel, S., Shrivastav, R., & Eddins, D. A. (2010). Perceptual distances of breathy voice quality:

689      A comparison of psychophysical methods. *Journal of Voice, 24*(2), 168-177.

690  Price, D. D., McGrath, P. A., Rafii, A., & Buckingham, B. (1983). The validation of visual

691      analogue scales as ratio scale measures for chronic and experimental pain. *Pain*, *17*(1),

692      45–56. https://doi.org/10.1016/0304-3959(83)90126-4

693  Ryan, M. L. (1971). *Investigation of the psychophysical relationship of kinesthetic extent of*

694      *arm movement* (Doctoral dissertation, University of British Columbia).

695  Schiavetti, N. (1992). Scaling procedures for the measurement of speech intelligibility. In R.

696      Kent (Ed.), *Intelligibility in speech disorders* (pp. 11-34). Philadelphia, PA: John

697      Benjamins.

698  Schiavetti, N., Metz, D. E., & Sitler, R. W. (1981). Construct validity of direct magnitude

699      estimation and interval scaling of speech intelligbility: Evidence from a study of the

700      hearing impaired. *Journal of Speech and Hearing Research, 24*, 441-445.

701  Schiavetti, N., Sacco, P. R., Metz, D. E., & Sitler, R. W. (1983). Direct magnitude estimation

702      and interval scaling of stuttering serverity. *Journal of Speech and Hearing Research,*

703      *26*, 568-573.

704  Sewall, A., Weglarski, A., Metz, D. E., Schiavetti, N., & Whitehead, R. L. (1999). A

705      methodolohgical  control  study  of  scaled  vocal  breathiness  measurements.

706      *Contemporary Issues in Communication Sciences and Disorders*, *26*, 168-172.

707   Siegel, S. (1956). A method for obtaining an ordered metric scale. *Psychometrika, 21*(2), 207-

708       216.

709   Snijders, B. & Bosker, J. (2012). *Multilevel Analysis: An Introduction to Basic and Advanced*

710       *Multilevel Modeling*, 2nd Ed. Sage, London.

711   Southwood, M. H., & Weismer, G. (1993). Listener judgements of the bizarreness,

712       acceptibility, naturalness, and normalcy of the dysarthria associated with amyotrophic

713       lateral sclerosis. *Journal of Medical Speech-Language Pathology*, *1*, 151-161.

714   Stepp, C. E., Heaton, J. T., Rolland, R. G., & Hillman, R. E. (2009). Neck and face surface

715       electromyography for Prosthetic voice control after total laryngectomy. *IEEE*

716       *Transactions on Neural Systems and Rehabilitation Engineering*, *17*(2), 146– 155.

717       https://doi.org/10.1109/TNSRE.2009.2017805

718   Stevens, S. S. (1958). Measurement and man. *Science, 127*(3295), 383-389.

719   Stevens, S. S. (1969). Sensory scales of taste intensity. *Perception & Psychophysics, 6*, 302-

720       308.

721   Stevens, S. S. (1975). *Psychophysics: Introduction to its perceptual, neural, and social*

722       *prospects.* Wiley.

723   Stevens, S. S., & Galanter, E. H. (1957). Ratio scales and category scales for a dozen perceptual

724       continua. *Journal of Experimental Psychology*, *54*, 377-411.

725   Stipancic, K. L., Golzy, M., Zhao, Y., Pinkerton, L., Rohl, A., & Kuruvilla-Dugdale, M. (2023).

726       Improving Perceptual Speech Ratings: The Effects of Auditory Training on Judgments

727       of Dysarthric Speech. *Journal of Speech, Language, and Hearing Research*, *66*(11),

728       4236–4258. https://doi.org/10.1044/2023_JSLHR-23-00322.

729 Stipancic, K. L., Palmer, K. M., Rowe, H. P., Yunusova, Y., & Green, J. R. (2021). "You say

730 severe, I say mild": Toward an empirical classification of dysarthria severity. *Journal*

731 *of Speech, Language, and Hearing Research, 64*, 4718-4735.

732 Stipancic, K. L., Tjaden, K., & Wilding, G. (2016). Comparison of intelligibility measures for

733 adults with Parkinson's disease, adults with multiple sclerosis, and healthy controls.

734 *Journal of Speech, Language, and Hearing Research, 59*(2), 230-238.

735 Trevethan, R. (2017). Intraclass correlation coefficients: clearing the air, extending some

736 cautions, and making some requests. *Health Services and Outcomes Research*

737 *Methodology*, *17*(2), 127-143.

738 Whitehill, T. L., Lee, A. S. Y., & Chun, J. C. (2002). Direct magnitude estimation and interval

739 scaling of hypernasality. *Journal of Speech, Language and Hearing Research*, *45*(1),

740 80-88.

741 Wolfe, V. I., & Ratusnik, D. L. (1988). Acoustic and perceptual measurements of roughness

742 influencing judgements of pitch. *Journal of Speech and Hearing Disorders*, *53*, 15-22.

743 Yiu, E. M.-L., & Ng, C.-Y. (2004). Equal appearing interval and visual analogue scaling of

744 perceptual roughness and breathiness. *Clinical Linguistics and Phonetics*, *18*(3), 211-

745 229. https://doi.org/https://doi.org/10.1080/0269920042000193599

746 Yorkston, K. M., Beukelman, D., & Hakel, M. (2007). *Speech Intelligibility Test (SIT) for*

747 *Windows [Computer software]*. Madonna Rehabilitation Hospital.

748 Zraik, R. I., & Liss, J. M. (2000). A comparison of equal-appearing interval scaling and direct

749 magnitude estimation of nasal voice quality. *Journal of Speech, Language and Hearing*

750 *Research*, *43*, 979-988. https://doi.org/1092-4388/00/4304-0979

751 Zyski, B. J., & Weisiger, B. E. (1987). Identification of dysarthria types based on perceptual

752 analysis. *Journal of Communication Disorders*, *20*, 367-378.

753