

Purdue University

Purdue e-Pubs

Department of Food Science Faculty
Publications

Department of Food Science

2019

Data approximation strategies between generalized line scales and the influence of labels and spacing

Jonathan C. Kershaw

Cordelia Running

Follow this and additional works at: <https://docs.lib.purdue.edu/foodscipubs>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

This is the author copy of an accepted manuscript, posted to the Purdue University Repository after an embargo period as permitted by the publisher.

The published copy can be found at:

[Data approximation strategies between generalized line scales and the influence of labels and spacing](#)

JC Kershaw, CA Running

Journal of Sensory Studies, e12507

<https://doi.org/10.1111/joss.12507>

1 Data approximation strategies between
2 generalized line scales and the influence of
3 labels and spacing
4

5 **Running title:** Data approximation strategies

6

7 Jonathan C. Kershaw^{1,2,3}

8 Cordelia A. Running^{1,2*}

9

10 ¹Department of Nutrition Science, 700 W State St, Purdue University, West Lafayette IN USA

11 ²Department of Food Science, 700 W State St, Purdue University, West Lafayette IN USA

12 ³Department of Public and Allied Health, 136 Health and Human Services, Bowling Green State
13 University, Bowling Green OH USA

14

15 *Corresponding author: crunning@purdue.edu

16 **Abstract**

17 Comparing sensory data gathered using different line scales is challenging. We tested whether
18 adding internal labels to a generalized visual analog scale (gVAS) would improve comparability
19 to a typical generalized labeled magnitude scale (gLMS). Untrained participants evaluated
20 cheeses using one of four randomly assigned scales. Normalization to a cross-modal standard
21 and/or two gLMS transformations were applied to the data. Response means and distributions
22 were lower for the gLMS than the gVAS, but no difference in resolving power was detected. The
23 presence of words, independent of internal lines, induced categorical behavior (marking close to
24 labels). Closer low-end label spacing for gLMS increased influenced participants to mark near
25 higher intensity labels when they were evaluating low-intensity samples. Although normalization
26 reduced differences between scales, neither transformation nor normalization was supported as
27 appropriate gLMS/gVAS approximation strategies. This study supports previous observations
28 that neither scale offers a systematic advantage and that participant usage differences limit direct
29 scale comparisons.

30

31 **Practical Applications**

32 Practitioners should exercise caution when comparing between gVAS and gLMS data, as neither
33 normalization nor transformations make these equivalent. The value of qualitative information
34 from internal labels (gLMS) and the expected intensity of samples should be considered when
35 choosing a scale (gVAS may be better for lower intensity samples, gLMS for high intensity).
36 While all scales in this study provide valid information regarding sample intensity, the impact of
37 scale on statistical assumptions, most notably the non-normal distribution of residuals from
38 gLMS data, should also be considered and corrected when necessary.

39

40 **Key words:** Scale selection, categorical behavior, transformation, normalization, gLMS, gVAS

41

42

43

44 **1. Introduction**

45 Continuous visual analog scales (VAS) are commonly used to quantify and compare sensory
46 experiences. The VAS is a continuous line anchored at either end by a minimum and maximum,
47 and is considered a practical and reliable tool for measuring and comparing human experiences
48 across different populations (Price, McGrath, Rafii, & Buckingham, 1983; Zealley & Aitken,
49 1969). Furthermore, the continuous nature of the scale allows for more sensitive statistical tests
50 compared with category scales. In some uses of the scale, internal semantic labels are added to
51 provide qualitative information about sensation intensity. However, VAS are not without
52 limitations (Bartoshuk et al., 2003). A desirable line scale would provide semantic information;
53 produce normally distributed, ratio-level data; and have high resolving power (i.e., the ability to
54 detect differences between distinct samples), among other desirable attributes.

55
56 In sensory research, the labeled magnitude scale (LMS) was developed to produce data with
57 ratio-level properties and qualitative meaning, which previously required the use of either
58 magnitude estimation or category scales, respectively (Green, Shaffer, & Gilmore, 1993). On the
59 LMS, internal labels are spaced quasi-logarithmically based on the quantitative derivation of
60 semantic descriptors. Building on findings of the LMS, other researchers suggested a top anchor
61 of “the strongest imaginable sensation of any kind,” to compare results across populations with
62 systematic differences in experiences and physiology, such as genetic differences in 6-n-
63 propylthiouracil sensitivity (Bartoshuk et al., 2003, 2004). This scale was designated as a
64 generalized labeled magnitude scale (gLMS). Further refinement of the scale revealed that the
65 descriptor “imaginable” does not improve across group comparisons (Bartoshuk, Fast, & Snyder,
66 2005). Others have applied a generalized, cross-modal top anchor to a VAS with no internal

67 labels and designated it as the general visual analog scale (gVAS) (Bartoshuk et al., 2005). Like
68 the gLMS, the gVAS can also produce ratio level data (Hayes, Allen, & Bennett, 2013). Of note,
69 generalized scales improve across group comparisons but do not necessarily offer an advantage
70 for within subject analyses (Kalva et al., 2014). Both of these generalized scales, as well as a
71 variety of derivatives from them, are commonly used in sensory evaluation.

72
73 The presence of internal labels, the key difference between the gLMS and the gVAS, provides
74 both a benefit and a limitation: although labels provide meaningful qualitative information, they
75 also influence participants to rate closely to the markings (i.e., categorical behavior), thus
76 producing clustered rather than continuous data (Hayes et al., 2013). As only one study has
77 directly compared the gLMS and the gVAS (Hayes et al., 2013), further research is helpful for
78 understanding how these scales compare in a variety of settings and sample sets, as well as for
79 selecting among these scales and other potential derivatives of these scales.

80
81 Although the gLMS and gVAS use identical end anchors, systematic differences in the way
82 participants use the scales limit inter-scale comparisons (Hayes et al., 2013). A method to
83 approximate sample means generated by the two scales would facilitate cross-study comparisons
84 and aid appropriate scale selection. To compare across scales, systematic transformations have
85 been applied in a number of contexts for both intensity and hedonic ratings (Green et al., 1993;
86 Lim, Wood, & Green, 2009; Schutz & Cardello, 2001). Normalizing individual ratings to a
87 cross-modal standard (e.g. the brightness of the sun, the intensity of a tone, or the heaviness of
88 jars of sand) has also been used to compare scales (Duffy, Peterson, & Bartoshuk, 2004; Hayes
89 et al., 2013; Webb, Bolhuis, Cicerale, Hayes, & Keast, 2015). In addition, a greater

90 understanding of how scale elements (i.e., anchors, label presence and spacing, etc.) influence
91 participant behavior could aid scale selection and optimization.

92
93 The purpose of this study is to explore the comparability of a gLMS and gVAS, following
94 normalization and/or transformations, and to assess how label presence and spacing influence
95 participant responses.

96

97 **2. Methods**

98 *2.1 Study participants and procedures*

99 Healthy participants were recruited at the Indiana State Fair (n=195, 130 women, 65 men, 0
100 other, ages 8-80, mean age 38.2) to evaluate three cheese samples (blue cheese crumbles, white
101 cheese cubes, and shaved parmesan cheese, donated by the American Dairy Association of
102 Indiana). All protocols were approved as exempt by Purdue University's Institutional Review
103 Board for Human Subjects Research under category 6, testing of foods and food ingredients.
104 Inclusion criteria included absence of food allergies and willingness to eat and evaluate cheeses;
105 acceptance of these criteria was done electronically and participation constituted consent.
106 Participants first answered six cross-modal "warm-up" questions about remembered or imagined
107 sensations to familiarize participants with the scale and to verify that they understood its usage
108 (Table 1, a reduced form of the warm-up from Hayes et al., 2013; responses from participants
109 that incorrectly used the scales were excluded data analyses, as described below). Following the
110 warm-up questions, participants evaluated the three cheese samples in a counter-balanced order.
111 A five second wait time was enforced between each sample. Participants were randomly
112 assigned one of four scales (described in detail below) and asked to rate the odor intensity of the

113 sample (participants were allowed to taste the samples as well, but not all participants chose to
114 do so; thus, we will only report on the odor data). In addition to a gLMS and a gVAS, two other
115 scales were created with the explicit purpose of investigating scale elements (Figure 1). The
116 gVAS labeled with equally-spaced gLMS labels is designated as the “generalized labeled VAS”
117 (gIVAS), and the same scale without the internal lines is designated as the “generalized words-
118 only VAS” (gwVAS). The purpose of adding these scales to the analysis was to see if spacing
119 out the internal lines would improve resolving power for lower intensity ratings (gIVAS), and if
120 removing the actual tick marks from the main line would reduce clustering of the ratings
121 (gwVAS); these scales have not been validated for general use.

122

123 **Table 1.** Instructions and questions presented to participants to familiarize them with generalized
124 scales.

First we'd like to familiarize you with our rating system. You will rate the strength of several imagined or remembered sensations. Please rate the sensations to the best of your ability. We use this information to verify that you are reading the directions! If you have not experienced a particular sensation, rate how intense you imagine it is. If you don't pay attention here, we cannot use ANY of your data from the rest of the test. This makes us very sad.
The brightness of this room
The brightness of the sun on a clear day
The loudness of a shout
The loudness of a whisper
The bitterness of black coffee
The sweetness of pure sugar

125

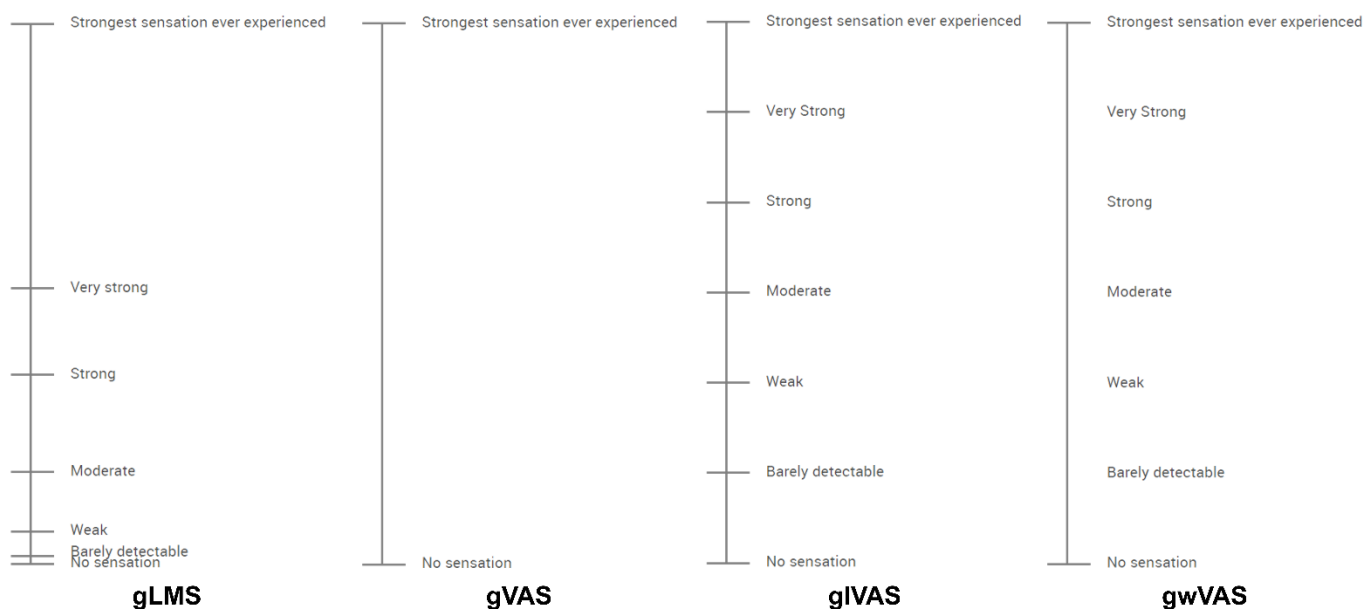
126

127

128

129

130



132 *Figure 1 Scales presented to study participants. gLMS: generalized labeled magnitude scale; gVAS: generalized visual anal*
 133 *scale; gIVAS: generalized labeled visual analog scale; gwVAS: generalized words-only visual analog scale.*

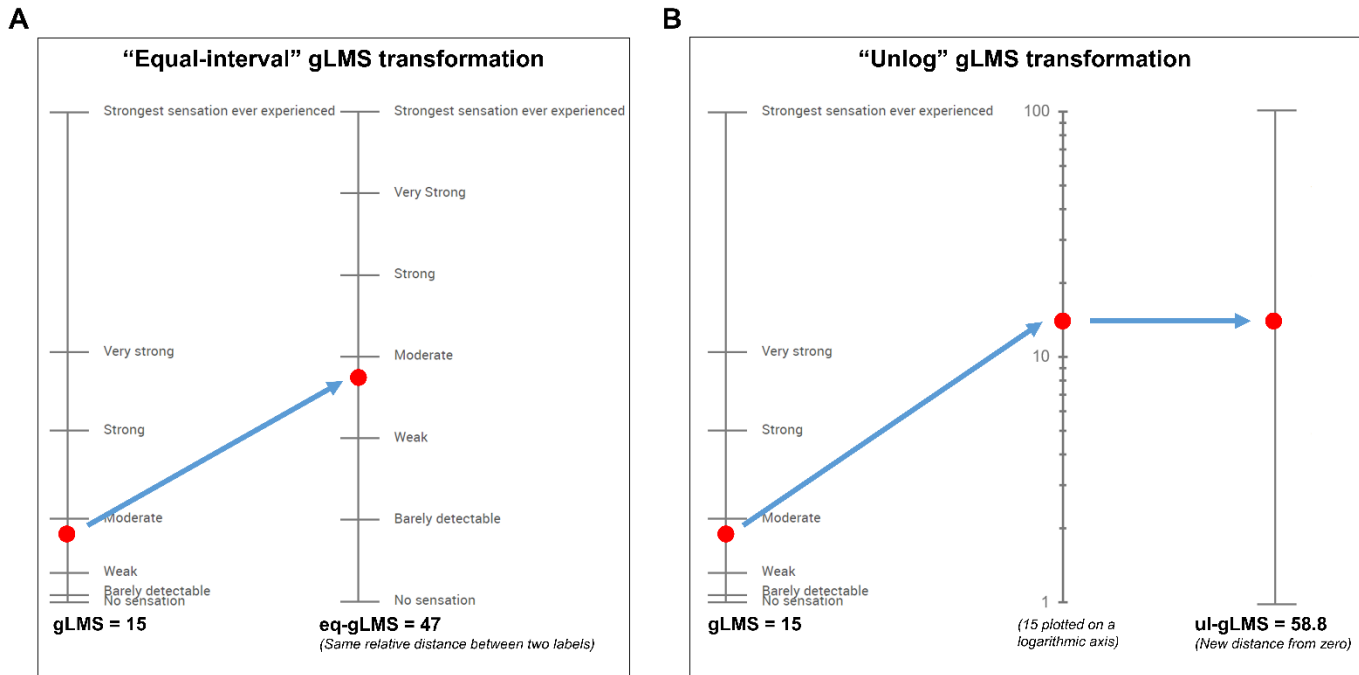
134

135 *2.2 Transformations and normalization*

136 To compare responses generated by the gLMS and the gVAS, two transformations were applied
 137 to gLMS data, as shown in Figure 2. In the equal-interval transformed gLMS (eq-gLMS), each
 138 rating was transformed to preserve the relative distance between its two closest internal
 139 markings, as if it were marked on a scale with equally spaced labels, based on the methods of
 140 Green et al. (Green et al., 1993). For example, 15 on the gLMS is 9/11ths between weak (6) and
 141 moderate (17), so 9/11ths between weak (16.33) and moderate (50) on an equally spaced scale
 142 would be 47. In the “unlog” transformation, gLMS response values were plotted on a 1-100
 143 logarithmic scale, then converted to the value that represented the distance from zero, as if it
 144 were rated using a standard 100mm scale ($[\log_{10}(\text{rating})]*50$) (Figure 2). Normalized values
 145 were obtained by dividing participant responses by their rating for “the brightness of the sun on a

146 clear day” then multiplying by 100. When applicable, normalization was performed before
147 transformation.

148



150 *Figure 2. Examples of transformations applied to data generated by the gLMS. A) In the equal-interval gLMS (eq-gLMS)*
151 *transformation, 15 becomes 47, as both are 9/11ths between weak and moderate on their respective scales. B) In the unlog gLMS*
152 *(ul-gLMS) transformation, 15 becomes 58.8 by plotting it on a 1-100 logarithmic scale and then converting it to the value that*
153 *represents the distance from zero as if it were a 100mm scale ($[\log_{10}(15)] * 50$).*

154

155 2.3 Statistical analysis

156 Data from participants that incorrectly used the scale (i.e., rated “the brightness of this room”
157 higher than “the brightness of the sun on a clear day” or “the loudness of a whisper” higher than
158 “the loudness of a shout”) were excluded from statistical analysis (final n=183, 120 women, 63
159 men, 0 other, ages 8-80, mean age 38.5). The number of analyzed participants (and fails) for
160 gLMS, gVAS, glVAS, and gwVAS, respectively was 48(2), 48(2), 46(4), and 43(4). Differences
161 in sample means, resolving power, residual distributions, overall sample distribution, and

162 categorical behavior for both cheese odor intensity and warm-up questions were compared as
163 follows: 1) gLMS vs. gVAS; 2) gVAS vs. gLMS transformation/normalization to assess
164 strategies for inter-scale comparisons; and 3) gVAS, gIVAS, gwVAS, and eq-gLMS to assess the
165 impact of label presence and spacing. Statistical significance was set at $p < 0.05$, with no
166 adjustments for multiple comparisons in order to set stricter standards for accepting the null
167 hypothesis of no differences between scales.

168

169 Differences in sample means were assessed using linear mixed models (“proc mixed”) in SAS
170 9.4. Participant was set as the repeated factor using the autoregressive covariance structure (data
171 sorted by cheese type, then scale, participant ID, and cheese testing order) and the Kenward-
172 Roger approximation for denominator degrees of freedom. Pair-wise comparisons within a
173 cheese type (or question, in the case of analyzing the warm-up data) were assessed using the
174 LSmeans function. To observe how scale type influenced resolving power, sample means were
175 compared using a similar linear mixed model, with the exception that data was analyzed by scale
176 rather than by cheese (following a sorting first by scale, then cheese). Residuals plots for each
177 scale type were assessed qualitatively based on SAS residual outputs. Differences in response
178 distributions between scales were determined using the Kolmogorov-Smirnov test in R-Studio, R
179 version 3.4.3.

180

181 Categorical behavior was defined as ratings falling within 1 point of the scale label (Hayes et al.,
182 2013; Lawless, Popper, & Kroll, 2010). Values for internal labels were rounded to the nearest
183 whole number, as half points were not recorded. A Fisher’s Exact test comparing the proportion
184 of categorical behavior to the proportion that would occur by chance alone (0.19) was used to

185 determine the significance of categorical behavior. To visualize categorical behavior, participant
186 responses were assigned to the label category it most closely approximated (Schutz & Cardello,
187 2001).

188

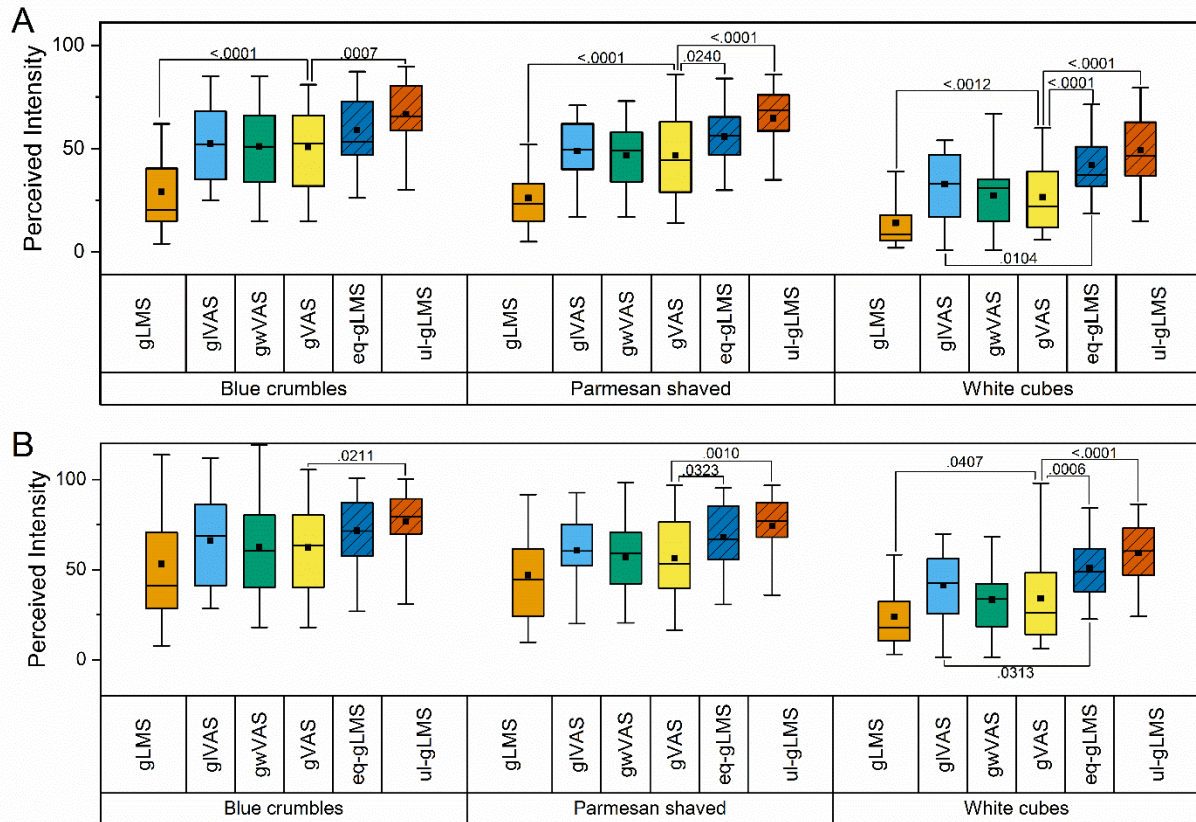
189 **3. Results**

190 *3.1 gLMS and gVAS comparison and the effects of data approximation strategies on sample*
191 *means, response distribution, resolving power, and residuals.*

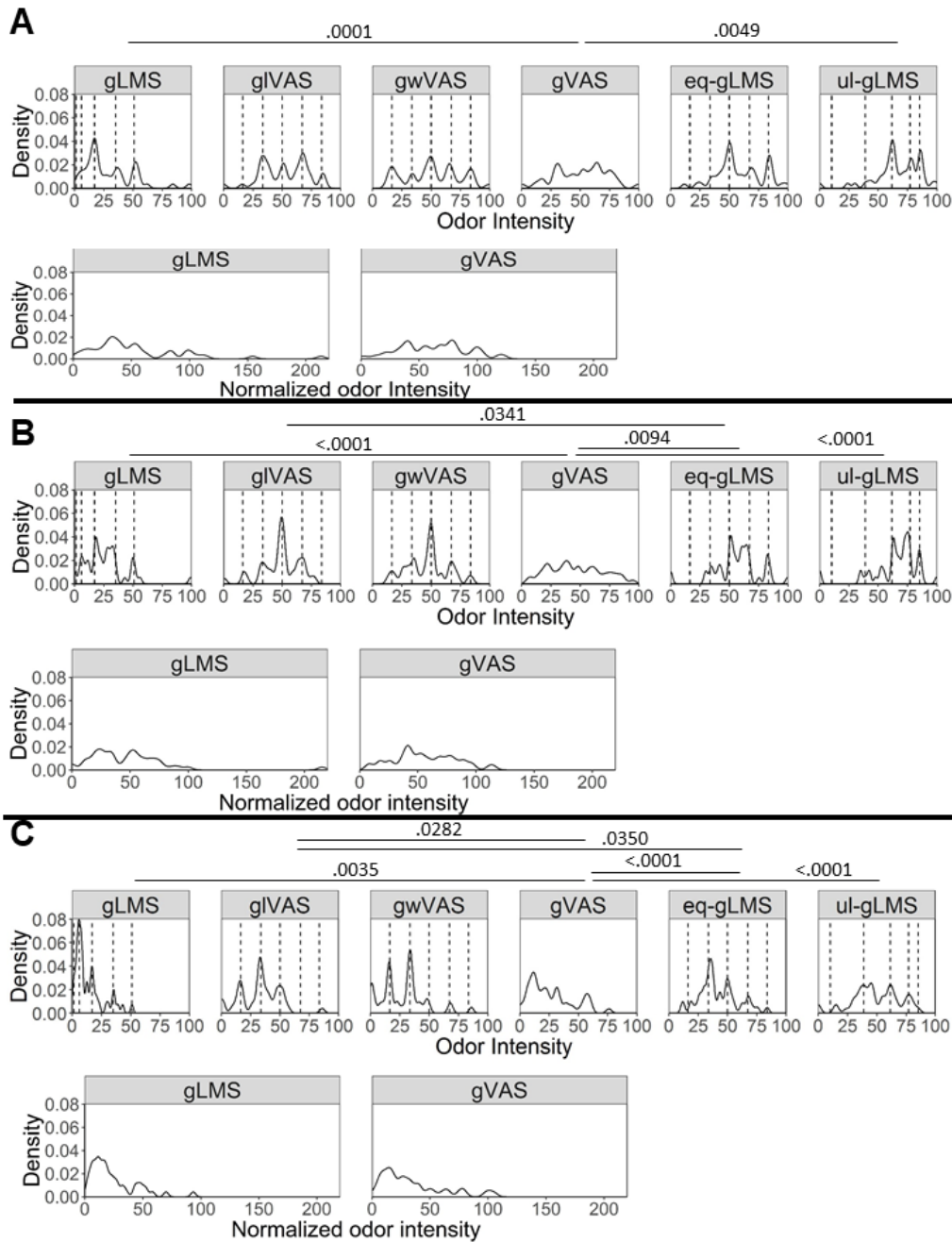
192 Consistent with other reports (Hayes et al., 2013), we observed significantly lower mean odor
193 intensity ratings generated by the gLMS compared with the gVAS (Figure 3, additional details
194 found in Supplemental table 1). Overall response distribution was also significantly different
195 between the gVAS and gLMS for all samples (Figure 4). Similar trends were observed in the
196 warm-up data (Supplemental figure 1, Supplemental table 2). Sample rank order and
197 discrimination sensitivity did not noticeably change between the two scales (Supplemental table
198 1), consistent with previous reports (Hayes et al., 2013). Residuals from the gLMS data showed
199 greater positive skew than those from the gVAS (Figure 5).

200

201



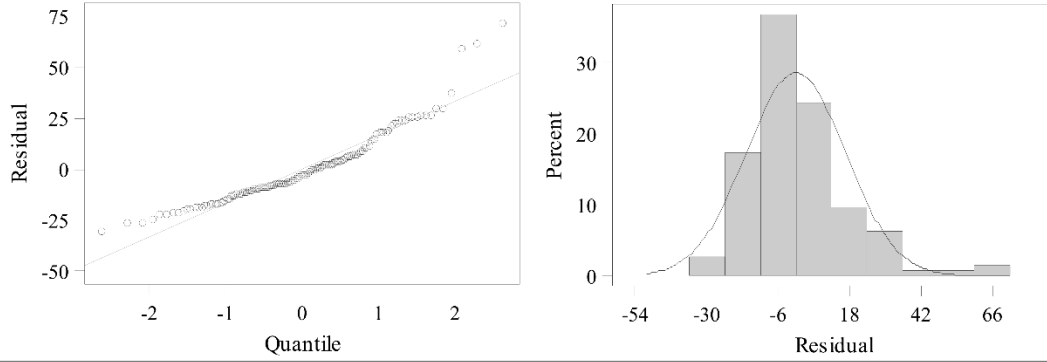
202
 203 *Figure 3. Box plots of raw a) and normalized b) cheese odor. Two gLMS transformations (equally spaced-gLMS and unLog-*
 204 *gLMS, indicated by hashed fill pattern) are also included. Boxes enclose the middle 50% of responses; the median is indicated by*
 205 *the center line. Scale means are represented by a square. Whiskers show 5th and 95th percentiles. Relevant pair-wise comparisons*
 206 *where $p < 0.05$ are displayed; no adjustments were made for multiple comparisons. Cheese odor differences within a scale are*
 207 *displayed in Supplemental table 1.*



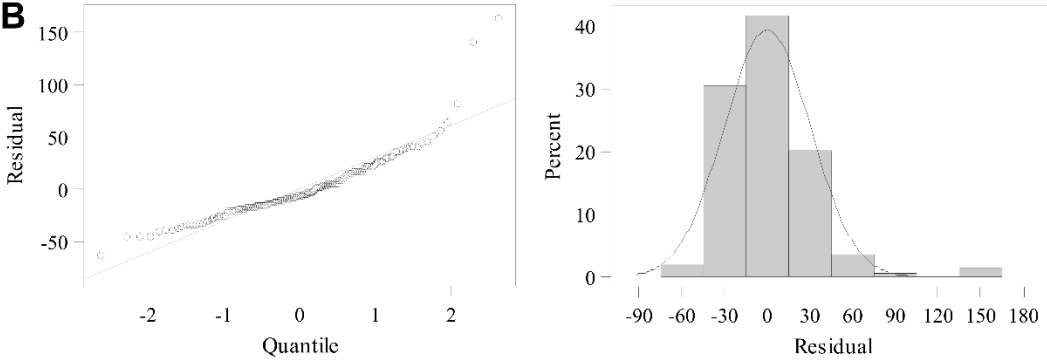
208

209 *Figure 4. Kernel density estimates for cheese odor from raw data, transformed gLMS data and selected normalized responses for*
 210 *a) blue, b) parmesan, and c) white cheese. Dashed lines indicate the location of internal labels, when present. Comparisons*
 211 *within non-normalized data where $p < 0.05$ using the Kolmogorov-Smirnov test are displayed; no adjustments were made for*
 212 *multiple comparisons. Kolmogorov-Smirnov tests between normalized gLMS vs. gVAS were greater than 0.05 for all cheese types.*

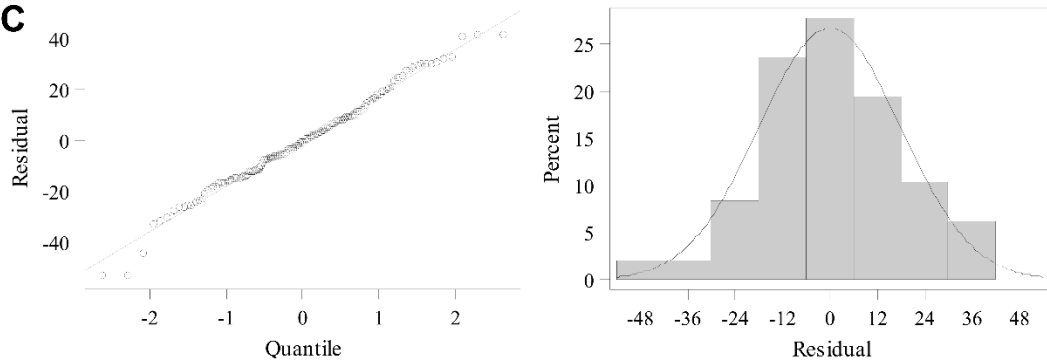
A



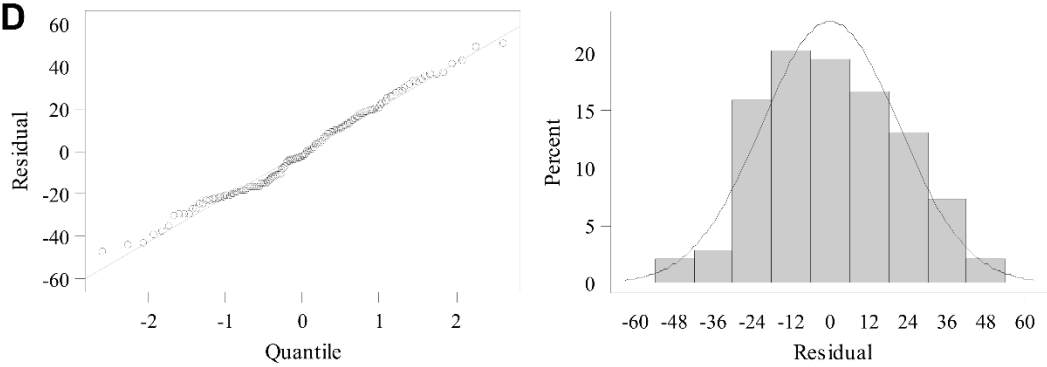
B



C



D



213
214
215

Figure 5 Residual plots of a) raw, b) normalized, and c) equal-interval-transformed gLMS and d) gVAS data, as generated using PROC MIXED in SAS 9.4. Data from all cheeses within a scale were combined to generate a single set of plots.

216

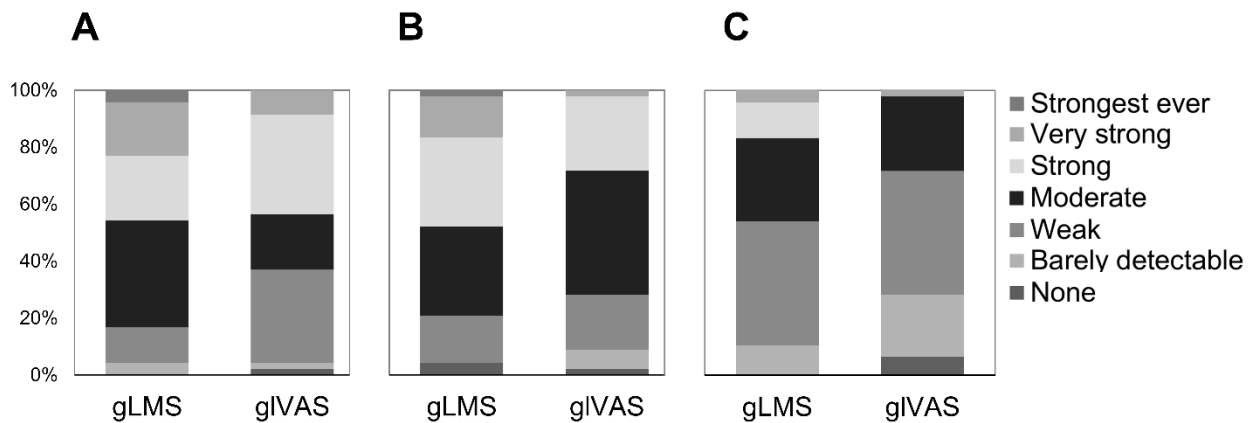
217 To investigate gVAS and gLMS data approximation strategies that could facilitate inter-scale
218 comparisons, we next compared gVAS means with transformed and/or normalized gLMS values
219 (Figures 3 and 4); fewer differences could indicate an appropriate comparison. Following the
220 equal-interval transformation, sample means and response distribution between the gLMS and
221 gVAS were significantly different in two of three cheese samples. Following the “unlog”
222 transformation, response distribution and sample means were still different for all cheese types.
223 Although fewer differences between the gLMS and gVAS data were observed following
224 normalization, in line with findings by others (Hayes et al., 2013), higher variance was also
225 observed; thus, statistical power to find differences was also reduced. Transforming normalized
226 data increased the number of differences between the gLMS and gVAS compared to non-
227 transformed normalized data. Normalization did not alter conclusions in warm-up data; mean
228 responses and response distributions between gVAS and gLMS were still different. The equal-
229 interval transformation, more than normalization, reduced the skewness of gLMS residuals
230 (Figure 5).

231

232 *3.2 Effect of scale elements on sample means, response distribution, resolving power, and*
233 *categorical behavior.*

234 To evaluate the effect of label spacing on participant behavior, we compared the gIVAS with the
235 eq-gLMS transformation. Consistent with findings by Green et. al (Green et al., 1993), we
236 observed a tendency for higher sample means in eq-gLMS responses compared with gIVAS
237 responses (Figure 3). However, this difference was only statistically significant for the cheese
238 with the lowest odor intensity (white cheese). Likewise, differences in mean warm-up responses
239 between the eq-gLMS and gIVAS were only observed for “the loudness of a whisper.”

240 Normalization did not alter conclusions regarding sample mean comparison of these two scales.
 241 Differences between the response distribution of the eq-gLMS and gIVAS reached statistical
 242 significance for two of three cheese samples (Figure 4). Relative to the eq-gLMS, participants
 243 using the gIVAS showed a greater use of the lower end of the scale, as illustrated by shifted
 244 kernel density estimates (Figure 4) and redisplayed as a categorical data visualization for added
 245 clarity (Figure 6).
 246



247
 248 *Figure 6. Visualization of “categories” following conversion of raw data to a categorical scale for a) blue, b) parmesan, and c)*
 249 *white cheese odor.*

250
 251 We next investigated the effect of equally-spaced internal labels with (gIVAS) or without
 252 (gwVAS) lines on participant responses compared with no labels (gVAS). Both labeling
 253 approaches had a minimal effect on mean intensity rating and response distribution; only the
 254 response distribution of white cheese was statistically significant between the gVAS and the
 255 gIVAS (Figures 3 and 4). Similarly, minimal differences were observed when labels were added
 256 to the warm-up questions (Supplemental figure 1).
 257

258 Consistent with others' observations (Hayes et al., 2013), we observed categorical behavior in
 259 gLMS responses for cheese odor intensity and warm-up questions (qualitatively demonstrated by
 260 kernel density estimates in Figure 4). Indeed, 50% of gLMS responses fell within 1 point of the
 261 labeled lines (Table 2), which accounts for only 19% of the scale. The presence of words, with or
 262 without tick marks, was sufficient to induce categorical behavior, as categorical behavior was
 263 observed for both the gIVAS and gwVAS.

264
 265 **Table 2.** Categorical behavior (defined as marking within one point of an internal label) for each
 266 scale for both a) cheese odor and b) warm-up questions. A Fisher's Exact test comparing the
 267 proportion of categorical behavior to the proportion that would occur by chance alone (0.19) was
 268 used to determine the significance of categorical behavior.

A	gLMS	gIVAS	gwVAS
Categorical marks	72	53	61
Total responses	144	138	129
% Categorical	50.0	38.4	47.3
Lower Wald's CI	41.8	30.3	38.7
High Wald's CI	58.2	46.5	55.9
Test statistic¹	<.0001	<.0001	<.0001
¹ Fisher's exact test vs. 0.19, two-sided			
B	gLMS	gIVAS	gwVAS
Categorical marks	134	116	105
Total responses	276	276	258
% Categorical	48.6	42.0	40.7
Lower Wald's CI	0.4265	0.3621	0.347
High Wald's CI	0.5445	0.4785	0.4669
Test statistic¹	<.0001	<.0001	<.0001
¹ Fisher's exact test vs. 0.19, two-sided			

269

270 **4. Discussion**

271 In this study, we systematically compared strategies to approximate gLMS and gVAS data to
272 explore the practicality of inter-scale comparisons. Our findings provide limited support for the
273 use of cross-modal normalization or transformation in this context. Furthermore, we explored the
274 effect of scale elements (i.e., internal label presence and spacing) on sample means, response
275 distribution (including categorical behavior), resolving power, and residual distribution. We
276 observed that for lower intensity samples, participants marked closer to more intense descriptors
277 on the gLMS (which compresses the descriptors on the lower end of the scale) compared to a
278 scale where those same descriptors were equally spaced. Additionally, the presence of labels,
279 with or without lines, induced categorical behavior.

280

281 *4.1 gLMS vs. gVAS*

282 The first objective of this study was to compare two widely used scales in sensory analysis, the
283 gLMS and gVAS, and whether transformation and/or normalization facilitated their comparison.
284 Our observation of lower sample means in the gLMS is consistent with previous observations
285 that compression of gLMS sample means is due in part to internal labeling (Hayes et al., 2013),
286 as anchors are identical between these two scales. Among our limited sample set, we detected no
287 differences in resolving power between the gLMS and gVAS, consistent with other reports
288 (Hayes et al., 2013). Although some have proposed that a gLMS may have reduced
289 discrimination sensitivity due to response compression (Cardello, Lawless, & Schutz, 2008;
290 Lawless, Popper, et al., 2010), scale compression may also increase resolving power due to
291 compressed variance (Cardello et al., 2008). Our findings are consistent with others that have
292 failed to detect differences in rank order and sensitivity among scales with differences in

293 response compression (Cardello et al., 2008; Hayes et al., 2013; Lawless, Popper, et al., 2010;
294 Ludy & Mattes, 2011). Although labeled magnitude scales may offer advantages when extreme
295 ratings are expected (Bartoshuk et al., 2003; Cardello et al., 2008; Lim et al., 2009; Schutz &
296 Cardello, 2001), for intensities encountered in everyday experiences, there may be no advantage
297 of one scale over others (Lawless, Sinopoli, & Chapman, 2010; Ludy & Mattes, 2011). Clearly,
298 testing context must be considered when selecting the most appropriate scale.

299

300 Parametric statistical tests, which are commonly used to analyze results from both the gLMS and
301 gVAS, require an assumption of equal variances and normally distributed residuals. However,
302 these assumptions may be violated more often when using the gLMS, as we observed that
303 residuals from the gVAS were less skewed than those from the gLMS. Additional investigation
304 is needed to quantify the extent that violations of statistical assumptions influence conclusions
305 from analyzing such data with parametric tests; however, from a technical perspective,
306 researchers are advised to check the distributions of residuals if a gLMS is selected and
307 transform the data as appropriate to correct the issue.

308

309 *4.2 Data approximation strategies*

310 Transformations and cross-modal normalizations of raw data have been performed as strategies
311 to facilitate inter-scale comparisons (Hayes et al., 2013; Lim et al., 2009; Schutz & Cardello,
312 2001). In evaluating two transformation strategies for gVAS-gLMS comparisons, we observed
313 that the equal-interval transformation eliminated statistical differences between sample means for
314 only a limited number of samples; thus, our results do not support the use of either attempted
315 transformation strategy for inter-scale comparisons. While normalizing responses resulted in

316 fewer differences in sample means and response distributions between the gLMS and the gVAS,
317 greater variance was also observed; thus we conclude that our failure to find differences was
318 likely due to reduced statistical power rather than improved data approximation. Although some
319 have proposed that normalizing raw data to a participant's own rating of a cross-modal standard
320 may control for idiosyncratic or systematic differences in scale usage (Bartoshuk et al., 2004),
321 more studies are needed to determine the most appropriate situations for the use of normalized
322 data, and caution should be used when using normalized data to draw conclusions regarding
323 sample differences.

324

325 *4.3 Influence of label spacing and presence*

326 Early psychophysicists noted that intensity labels such as “weak” and “strong” do not possess
327 interval-level qualities (Borg, 1982; Lasagna, 1960). Furthermore, others have argued that
328 spacing labels consistent with the relative strength of their perceptual intensity is important for
329 proper scale usage (Green et al., 1993). Consequently, many researchers have attempted to
330 quantitatively determine appropriate spacing of semantic descriptors to generate ratio-level level,
331 including those that developed labeled magnitude scales (Green et al., 1993; Moskowitz, 1977;
332 Schutz & Cardello, 2001). In the present work, we revisited the effect of label spacing by
333 comparing responses from a scale with equally spaced labels (gIVAS) and an equal-interval
334 transformation of semi-logarithmically spaced labels (eq-gLMS). As eq-gLMS sample means
335 were significantly higher than gIVAS means only for weaker modalities (odor intensity of white
336 cheese, loudness of a whisper), we suspect that end-use avoidance (Lawless & Heymann, 2010)
337 influenced participant behavior. Similarly, Green and others observed higher responses from
338 their LMS compared to a transformed equal-interval scale; although, in contrast to our findings,

339 statistical significance was only observed for higher intensity samples (Green et al., 1993). As
340 the gLMS was designed in part to accommodate “ceiling” effects (L. M. Bartoshuk et al., 2004),
341 we suggest that the gLMS may not be appropriate for the evaluation of low-intensity samples by
342 untrained participants.

343

344 In addition to label spacing, we also considered the effect of label presence on categorical
345 behavior, or the tendency of participants to mark primarily where labels occur (Cardello et al.,
346 2008; Hayes et al., 2013; Lawless, Popper, et al., 2010). We observed categorical behavior in all
347 scales containing semantic labels, independent of internal line presence. If participants make
348 category rather than ratio-level judgements, label spacing becomes less relevant because the
349 scale no longer truly possesses continuous or ratio properties (Lawless, 2013). When participants
350 make category judgements, the issue of end-use avoidance becomes more salient, as end-use
351 avoidance is higher in category scales than continuous scales (Schutz & Cardello, 2001). We
352 note that labels may not be necessary to generate ratio-level data, as ratio-level data has been
353 generated using the gVAS (Hayes et al., 2013). Taken together, these observations further
354 support the exercise of caution when using the gLMS for assessment of low-intensity samples.

355

356 Despite categorical behavior, we observed similar sample means, response distributions, and
357 resolving power between the label-less gVAS and the labeled glVAS and gwVAS. Together, this
358 suggests a minimal impact of categorical behavior on study conclusions, at least in our sample
359 set. Although the presence of labels may influence behaviors that violate assumptions of truly
360 normally distributed responses, whether this categorical behavior actually alters conclusions
361 from data analysis is still unknown (Hayes et al., 2013).

362

363 Although labels can induce categorical behavior, internal markings may also make scales easier
364 to use and provide meaningful qualitative information (Hayes et al., 2013). However, qualitative
365 semantic labels also possess inherent limitations, as the same descriptors can have different
366 meanings to different people (Bartoshuk et al., 2005; Bartoshuk et al., 2003). Thus, when
367 selecting the appropriate scale for sensory research, the researcher should consider the value of
368 qualitative information, expected sample intensity, and potential effects of categorical behavior
369 on data analysis.

370

371 Conclusions from the current study must be interpreted within context of several limitations.
372 First, each participant only used one scale to evaluate cheese odor intensity rather than using all
373 scales. Additionally, only the odor of three quite distinct cheeses were assessed, so results may
374 not apply to alternative sample types or a more internally similar sample set. Although the
375 somewhat noisy testing environment of a fair booth is a study limitation, it also reflects a more
376 realistic setting for how people actually consume food.

377

378 **5. Conclusion**

379 After exploring methods to compare data generated by the gLMS and gVAS, we found limited
380 evidence to support transformation or cross-modal normalization as appropriate data
381 approximation strategies. Although normalization may provide a very rough approximation to
382 compare response distributions, the greater variance suggests caution should be used when
383 drawing conclusions regarding resolving power. Despite identical scale anchors, differences in
384 internal labels influence participant response behavior and thus hinder inter-scale comparisons.
385 These data support previous recommendations that the gLMS and gVAS are not one-to-one

386 comparable (Hayes et al., 2013), even when transformed or normalized. Our investigation of
387 label spacing and presence revealed that words alone are sufficient to induce categorical
388 behavior, and suggested that compressed labels (as in the gLMS) may cause participants to mark
389 near higher intensity descriptors when low-intensity samples are evaluated. Consistent with
390 several other reports, we did not find any systematic advantage of one scale over another
391 (Cardello et al., 2008; Hayes et al., 2013; Lawless, Popper, et al., 2010; Ludy & Mattes, 2011).
392 Overall, we conclude that all of these scales are valid and useful tools, but we urge researchers to
393 carefully select the scale and/or transformation method best suited for their samples, participant
394 group, and data analysis approach.

395

396 **Acknowledgments**

397 The authors thank the American Dairy Association of Indiana for providing the cheese samples
398 and Miguel Odrón for his assistance in carrying out the study. The authors declare no conflicts of
399 interest.

400 **References**

401 Bartoshuk, L.M., Fast, K., & Snyder, D. J. (2005). Differences in our sensory worlds:
402 Invalid comparisons with labeled scales. *Current Directions in Psychological Science*,
403 *14*(3), 122–125. <https://doi.org/10.1111/j.0963-7214.2005.00346.x>

404

405 Bartoshuk, L. M., Duffy, V. B., Fast, K., Green, B. G., Prutkin, J., & Snyder, D. J.
406 (2003). Labeled scales (e.g., category, Likert, VAS) and invalid across-group
407 comparisons: What we have learned from genetic variation in taste. *Food Quality and*
408 *Preference*, *14*(2), 125–138. [https://doi.org/10.1016/S0950-3293\(02\)00077-0](https://doi.org/10.1016/S0950-3293(02)00077-0)

409

410 Bartoshuk, L. M., Duffy, V. B., Green, B. G., Hoffman, H. J., Ko, C. W., Lucchina, L.
411 A., ... Weiffenbach, J. M. (2004). Valid across-group comparisons with labeled scales:
412 the gLMS versus magnitude matching. *Physiology and Behavior*, *82*(1), 109–114.
413 <https://doi.org/10.1016/j.physbeh.2004.02.033>

414

415 Borg, G. (1982). A category scale with ratio properties for intermodal and interindividual
416 comparisons (pp. 25–34). Presented at the XXII International Congress of Psychology,
417 North-Holland Pub. Co.

418

419 Cardello, A., Lawless, H. T., & Schutz, H. G. (2008). Effects of extreme anchors and
420 interior label spacing on labeled affective magnitude scales. *Food Quality and*
421 *Preference*, *19*(5), 473–480. <https://doi.org/10.1016/j.foodqual.2008.02.003>

422

423 Duffy, V. B., Peterson, J. M., & Bartoshuk, L. M. (2004). Associations between taste
424 genetics, oral sensation and alcohol intake. *Physiology and Behavior*, 82(2–3), 435–445.
425 <https://doi.org/10.1016/j.physbeh.2004.04.060>

426
427 Green, B. G., Shaffer, G. S., & Gilmore, M. M. (1993). Derivation and evaluation of a
428 semantic scale of oral sensation magnitude with apparent ratio properties. *Chemical*
429 *Senses*, 18(6), 683–702. <https://doi.org/10.1093/chemse/18.6.683>

430
431 Hayes, J. E., Allen, A. L., & Bennett, S. M. (2013). Direct comparison of the generalized
432 visual analog scale (gVAS) and general labeled magnitude scale (gLMS). *Food Quality*
433 *and Preference*, 28(1), 36–44. <https://doi.org/10.1016/j.foodqual.2012.07.012>

434
435 Kalva, J. J., Sims, C. A., Puentes, L. A., Snyder, D. J., & Bartoshuk, L. M. (2014).
436 Comparison of the hedonic general labeled magnitude scale with the hedonic 9-point
437 scale. *Journal of Food Science*, 79(2), S238–S245. [https://doi.org/10.1111/1750-](https://doi.org/10.1111/1750-3841.12342)
438 [3841.12342](https://doi.org/10.1111/1750-3841.12342)

439
440 Lasagna, L. (1960). The clinical measurement of pain. *Annals of the New York Academy*
441 *of Sciences*, 86(1), 28–37. <https://doi.org/10.1111/j.1749-6632.1960.tb42788.x>

442
443 Lawless, H. T. (2013). *Quantitative Sensory Analysis*. West Sussex, UK: Wiley-
444 Blackwell.

445

446 Lawless, H. T., & Heymann, H. (2010). *Sensory evaluation of food: principles and*
447 *practices*. Springer Science & Business Media.

448

449 Lawless, H. T., Popper, R., & Kroll, B. J. (2010). A comparison of the labeled magnitude
450 (LAM) scale, an 11-point category scale and the traditional 9-point hedonic scale. *Food*
451 *Quality and Preference*, 21(1), 4–12. <https://doi.org/10.1016/j.foodqual.2009.06.009>
452

453 Lawless, H. T., Sinopoli, D., & Chapman, K. W. (2010). A comparison of the labeled
454 affective magnitude scale and the 9-point hedonic scale and examination of categorical
455 behavior. *Journal of Sensory Studies*, 25(s1), 54–66. [https://doi.org/10.1111/j.1745-](https://doi.org/10.1111/j.1745-459X.2010.00279.x)
456 [459X.2010.00279.x](https://doi.org/10.1111/j.1745-459X.2010.00279.x)

457

458 Lim, J., Wood, A., & Green, B. G. (2009). Derivation and evaluation of a labeled hedonic
459 scale. *Chemical Senses*, 34, 34(9, 9), 739, 739–751.
460 <https://doi.org/10.1093/chemse/bjp054>, [10.1093/chemse/bjp054](https://doi.org/10.1093/chemse/bjp054)

461

462 Ludy, M.-J., & Mattes, R. D. (2011). Noxious stimuli sensitivity in regular spicy food
463 users and non-users: Comparison of visual analog and general labeled magnitude scaling.
464 *Chemosensory Perception*, 4(4), 123–133. <https://doi.org/10.1007/s12078-011-9100-x>
465

466 Moskowitz, H. R. (1977). Magnitude estimation: Notes on what, how, when, and why to
467 use it. *Journal of Food Quality*, 1(3), 195–227. [https://doi.org/10.1111/j.1745-](https://doi.org/10.1111/j.1745-4557.1977.tb00942.x)
468 [4557.1977.tb00942.x](https://doi.org/10.1111/j.1745-4557.1977.tb00942.x)

469

470 Price, D. D., McGrath, P. A., Rafii, A., & Buckingham, B. (1983). The validation of
471 visual analogue scales as ratio scale measures for chronic and experimental pain. *Pain*,
472 *17*(1), 45–56.

473

474 Schutz, H. G., & Cardello, A. V. (2001). A labeled affective magnitude (LAM) scale for
475 assessing food liking/disliking. *Journal of Sensory Studies*, *16*(2), 117–159.
476 <https://doi.org/10.1111/j.1745-459X.2001.tb00293.x>

477

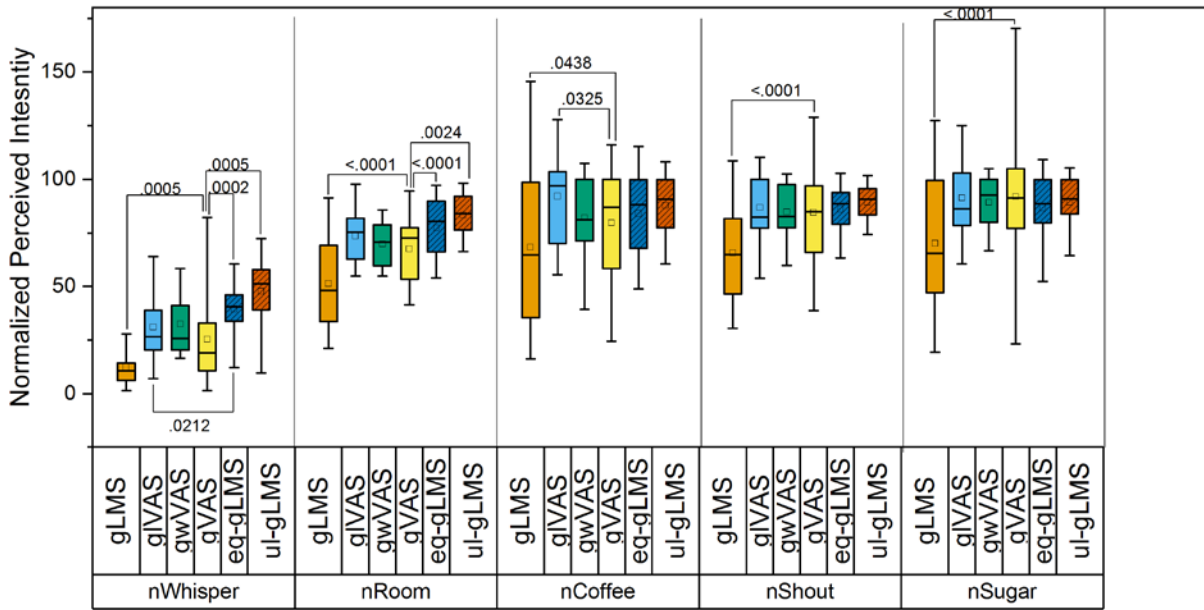
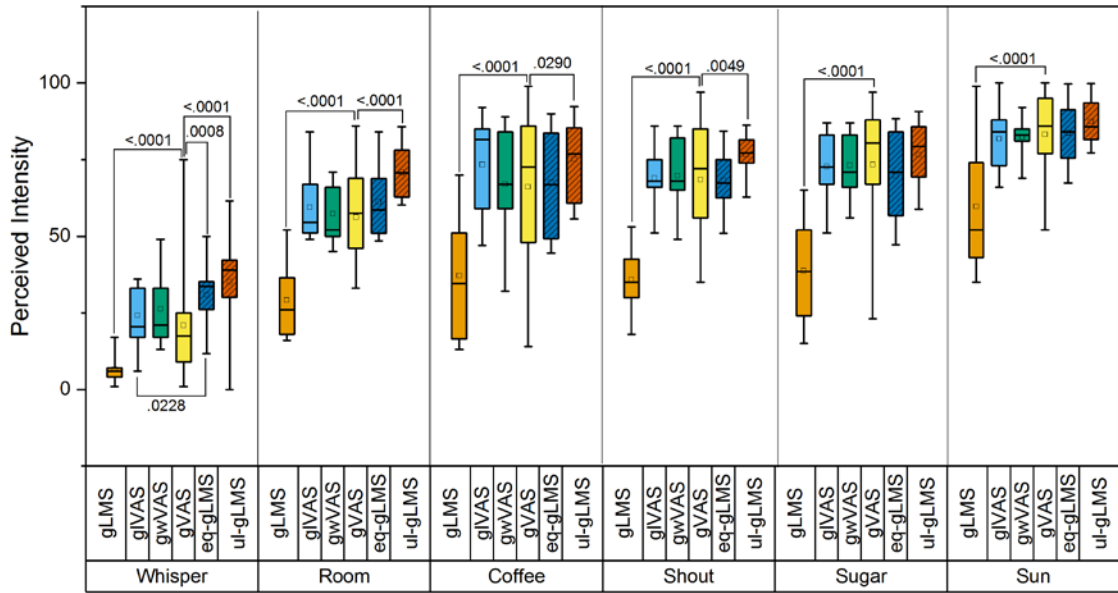
478 Webb, J., Bolhuis, D. P., Cicerale, S., Hayes, J. E., & Keast, R. (2015). The relationships
479 between common measurements of taste function. *Chemosensory Perception*, *8*(1), 11–
480 18. <https://doi.org/10.1007/s12078-015-9183-x>

481

482 Zealley, A. K., & Aitken, R. C. (1969). Measurement of mood. *Proceedings of the Royal*
483 *Society of Medicine*, *62*(10), 993–996.

484

485



490 Supplemental figure 1. Box plots of raw a) and normalized b) warm-up data. Two gLMS transformations (equally spaced-gLMS
 491 and unLog-gLMS, indicated by hashed fill pattern) are also included. Boxes enclose the middle 50% of responses; the median is
 492 indicated by the center line. Scale means are represented by a square. Whiskers show 5th and 95th percentiles. Relevant pair-wise
 493 comparisons were $p < 0.05$ are displayed; no adjustments were made for multiple comparisons. Question differences within a
 494 scale are displayed in table 1.

497 **Supplemental table 1.** Differences between response means within a scale (displayed in figure
 498 and supplemental figure 1). Normalized values are indicated by “n”. Samples with different
 499 superscripts indicate a p-value < 0.05 using the least squares difference between two
 500 samples/questions within a scale; no adjustments were made for multiple comparisons.

501

	gLMS	gIVAS	gwVAS	gVAS	eq-gLMS	unlog-gLMS
Blue	29.0 ^a	52.8 ^a	50.8 ^a	50.8 ^a	58.7 ^a	66.5 ^a
Parmesan	26.1 ^a	48.9 ^a	45.8 ^a	46.6 ^a	55.7 ^a	64.5 ^a
White	14.1 ^b	32.4 ^b	27.8 ^b	26.8 ^b	42.3 ^b	49.4 ^b
nBlue	48.3 ^a	66.9 ^a	61.9 ^a	62.1 ^a	99.6 ^a	113.8 ^a
nParmesan	41.7 ^a	60.6 ^a	55.8 ^a	56.4 ^a	90.2 ^a	105.9 ^a
nWhite	21.9 ^b	40.6 ^b	33.9 ^b	34.7 ^b	68.1 ^b	79.4 ^b
Whisper	6.7 ^a	24.2 ^a	26.3 ^a	20.9 ^a	30.8 ^a	35.2 ^a
Room	29.2 ^b	59.5 ^b	57.3 ^b	56.3 ^b	61.2 ^b	70.8 ^b
Coffee	37.2 ^{bc}	73.4 ^c	66.9 ^c	66.1 ^{bc}	67.5 ^{bc}	75.2 ^b
Shout	35.9 ^{bc}	69.0 ^c	69.8 ^c	68.5 ^c	68.2 ^{bc}	76.7 ^b
Sugar	38.8 ^c	72.7 ^c	73.1 ^c	73.4 ^{cd}	69.8 ^c	76.6 ^b
Sun	59.8 ^d	81.7 ^d	82.4 ^d	83.2 ^d	83.8 ^d	87.5 ^c
nWhisper	12.4 ^a	31.1 ^a	32.3 ^a	25.4 ^a	39.5 ^a	47.7 ^a
nRoom	51.4 ^b	73.6 ^b	69.8 ^b	67.4 ^b	77.4 ^b	83.4 ^b
nCoffee	68.2 ^c	92.1 ^{cd}	82.0 ^c	79.7 ^{bc}	84.0 ^{bc}	87.8 ^b
nShout	65.7 ^c	86.8 ^c	84.8 ^c	84.5 ^c	86.1 ^c	89.3 ^b
nSugar	70.1 ^c	91.3 ^{cd}	89.2 ^c	91.8 ^{cd}	86.6 ^c	89.2 ^b
nSun	100.0 ^d	100.0 ^d	100.0 ^d	100.0 ^d	100.0 ^d	100.0 ^c

502

503

504 **Supplemental table 2.** Differences in response distributions of raw, transformed, and
 505 normalized (n, gLMS and gVAS only) warm-up questions, as measured using the Kolmogorov-
 506 Smirnov test. No adjustments were made for multiple comparisons.

Question	Comparison		KS statistic
Whisper	gLMS	gVAS	<.0001
Whisper	glVAS	gVAS	0.0037
Whisper	gwVAS	gVAS	0.0049
Whisper	glVAS	gwVAS	0.9123
Whisper	eq-gLMS	gVAS	0.0000
Whisper	ul-gLMS	gVAS	0.0000
Whisper	glVAS	eq-gLMS	0.0036
Whisper	n gLMS	n gVAS	0.0002
Room	gLMS	gVAS	<.0001
Room	glVAS	gVAS	0.0874
Room	gwVAS	gVAS	0.2515
Room	glVAS	gwVAS	0.8571
Room	eq-gLMS	gVAS	0.1308
Room	ul-gLMS	gVAS	<.0001
Room	glVAS	eq-gLMS	0.5892
Room	n gLMS	n gVAS	0.0006
Coffee	gLMS	gVAS	<.0001
Coffee	glVAS	gVAS	0.2271
Coffee	gwVAS	gVAS	0.5702
Coffee	glVAS	gwVAS	0.2715
Coffee	eq-gLMS	gVAS	0.4695
Coffee	ul-gLMS	gVAS	0.0469
Coffee	glVAS	eq-gLMS	0.1447
Coffee	n gLMS	n gVAS	0.0155
Shout	gLMS	gVAS	<.0001
Shout	glVAS	gVAS	0.2271
Shout	gwVAS	gVAS	0.6179
Shout	glVAS	gwVAS	0.3472
Shout	eq-gLMS	gVAS	0.1335
Shout	ul-gLMS	gVAS	0.0033
Shout	glVAS	eq-gLMS	0.6408
Shout	n gLMS	n gVAS	0.0021
Sugar	gLMS	gVAS	<.0001
Sugar	glVAS	gVAS	0.2271
Sugar	gwVAS	gVAS	0.1625
Sugar	glVAS	gwVAS	0.9894
Sugar	eq-gLMS	gVAS	0.1976
Sugar	ul-gLMS	gVAS	0.6260
Sugar	glVAS	eq-gLMS	0.2090
Sugar	n gLMS	n gVAS	0.0003
Sun	gLMS	gVAS	<.0001
Sun	glVAS	gVAS	0.1440
Sun	gwVAS	gVAS	0.0134
Sun	glVAS	gwVAS	0.3063
Sun	eq-gLMS	gVAS	0.5530
Sun	ul-gLMS	gVAS	0.1308
Sun	glVAS	eq-gLMS	0.7648
Sun	n gLMS	n gVAS	--

507


```

509 Code
510
511 proc sort data=Cheese;
512 by cheese scale ID order;
513 run;
514 *Comparing SCALE by cheese odor;
515 ods graphics on;
516 title 'Scale differences by cheese odor';
517 ods output tests3=effects diffs=MeansDiffs LSmeans=LSmeans;
518 proc mixed data=Cheese;
519     by cheese;
520     class ID scale order;
521     model Odor= scale order scale*order / residual s outp=pred;
522     repeated / subject = id type= ar(1);
523     lsmeans scale / diff adjust=tukey;
524 run;
525 proc print data=effects;
526 run;
527 proc print data=MeansDiffs;
528 run;
529 proc print data=LSmeans;
530 run;
531
532 proc sort data=Cheese;
533 by scale cheese ID order;
534 run;
535 *Comparing cheese odor by SCALE;
536 ods graphics on;
537 title 'Odor by scale';
538 ods output tests3=effects diffs=MeansDiffs LSmeans=LSmeans;
539 proc mixed data=Cheese;
540     by scale;
541     class ID Cheese order;
542     model Odor= cheese order cheese*order / residual s
543 outp=pred;
544     repeated / subject = id type= ar(1);
545     lsmeans cheese / diff adjust=tukey;
546 run;
547 proc print data=effects;
548 run;
549 proc print data=MeansDiffs;
550 run;
551 proc print data=LSmeans;
552 run;
553

```

```
554 proc freq data= ms_elms order=data;
555     tables scale*YN / fisher;
556     exact fisher;
557     weight freq;
558 run;
559
560 proc freq data= Categ order=data;
561     by scale;
562     tables YN / binomial(equiv p=.19 margin=.05);
563     weight freq;
564 run;
565
```