

Is subjective perceptual similarity metacognitive?

Moharramipour Ali^{1*}, Zhou William¹, Rahnev Dobromir², Lau Hakwan^{1*}

¹ Center for Brain Science (CBS), RIKEN, Wako, Japan

² School of Psychology, Georgia Institute of Technology, Atlanta, United States

* Corresponding authors

Correspondence: alimoharrami1371@gmail.com, hakwan@gmail.com

ABSTRACT

Perceptual similarity is a cornerstone for human learning and generalization. However, in assessing the similarity between two stimuli differing in multiple dimensions, it is not well-defined which feature(s) one should focus on. The problem has accordingly been considered ill-posed. We hypothesize that similarity judgments may be, in a sense, metacognitive: The stimuli rated as subjectively similar are those that are in fact more challenging for oneself to discern in practice, in near-threshold settings (e.g., psychophysics experiments). This self-knowledge about one's own perceptual capacities provides a quasi-objective ground truth as to whether two stimuli 'should' be judged as similar. To test this idea, we measure perceptual discrimination capacity between face pairs, and ask subjects to rank the similarity between them. Based on pilot data, we hypothesize a positive association between perceptual discrimination capacity and subjective dissimilarity, with this association being importantly specific to each individual.

Keywords: similarity judgment, subjective perceptual similarity, perceptual discrimination capacity, metacognition, subjective perception

Introduction

30 Subjective perceptual similarity between stimulus pairs has long been studied in human
31 behavior. These studies explored various factors modulating similarity judgments, such as the
32 effects of knowledge and expertise, contextual cues, and the order of presenting the stimuli
33 (Shepard, 1964; Tversky, 1977; Smith, 1989; Smith & Heise, 1992; Medin et al., 1993). Different
34 theories and quantitative models of similarity have also been proposed (Nosofsky, 1984; Shepard,
35 1987; Smith, 1989). For example, Roger Shepard famously formulated the universal law of
36 generalization, according to which humans respond in the same way to stimuli of high similarity,
37 and the probability of this generalization decays exponentially as the distance increases within a
38 putative metric psychological space (Shepard, 1987). Later, Shepard's law was expanded to
39 encompass general non-metrically structured spaces (Tenenbaum & Griffiths, 2001) and different
40 accounts; notably, the rate-distortion theory was proposed to explain its nature (Sims, 2018).
41 Intriguingly, recent research has demonstrated that the exponential similarity decay, coupled with
42 a signal detection theory, can also effectively capture observations in visual working memory
43 (Schurgin et al., 2020). There is also a rich history of studies utilizing similarity judgments, in
44 combination with multidimensional scaling, to uncover the underlying perceptual dimensions of
45 stimuli (Borg & Groenen, 2005; Hebart et al., 2020).

46 Similarity judgments are subjective, in that it is up to the subject to report how they feel about
47 the stimuli. Accordingly, some researchers have argued that similarity judgments may reflect key
48 aspects of conscious perception (Clark, 2000; Rosenthal, 2010; Malach, 2021; Lau et al., 2022;
49 Tallon-Baudry, 2022; Zeleznikow-Johnston et al., 2023; Moharramipour & Lau, 2024). However, the
50 essentially subjective nature of these judgments also led to the well-known critique that similarity
51 is perhaps an ill-posed problem: there is, in a sense, no objective ground-truth as to how similar
52 two things really are (Goodman, 1972; Medin et al., 1993). For example, Joe Biden may look more
53 similar to Hillary Clinton than to Barack Obama, with respect to skin color. However, if we focus on
54 gender-related facial features, Joe Biden may look more similar to Barack Obama. From the outset,
55 it is unclear which visual features one should focus on. This presents a challenging obstacle to
56 understanding the processes underlying similarity judgments, as mechanistic explanations of
57 perception often rely on characterizing the observer as performing optimal inference, given
58 existing constraints (Rao, 1999; Shen & Ma, 2016).

59 Following previous theoretical work (Lau et al., 2022), we hypothesize that subjective similarity
60 judgments may be normative and rational, in the sense that they are made systematically based
61 on the metacognitive access of our own perceptual abilities. Stimuli pairs judged to be more similar
62 are, in fact, more challenging for oneself to discern in practice. If one judges two perceptual stimuli
63 to be highly dissimilar, and yet fails to distinguish them in psychophysical tasks, the said similarity
64 judgment can be regarded as 'incorrect' in a meaningful sense.

65 Revisiting the above example of how subjectively similar two faces are, the idea is that such
66 judgment would be made on a dimension in which all relevant features are optimally combined,
67 such that along this dimension, the two faces are maximally distinguishable. Specifically, for this
68 combination to be optimal, the choice of this dimension should be based on how perceptible each

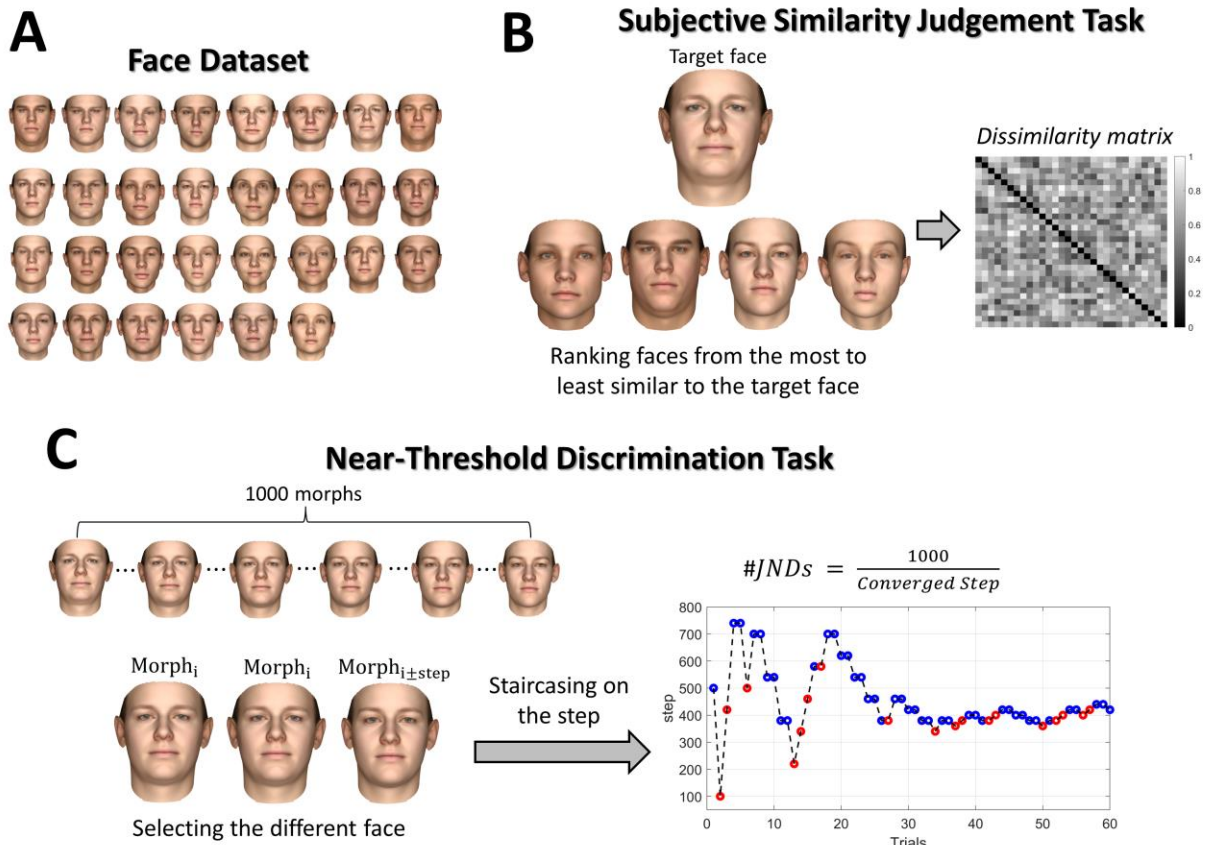
69 feature is to oneself. In other words, this process is not only about the physical stimulus itself, but
70 rather, it reflects (implicit) metacognitive knowledge of one's own perceptual abilities.

71 Metacognition is commonly defined as the monitoring (and control) of one's own cognitive
72 abilities_(Morales et al., 2018). In the present study, we hypothesize that similarity judgments
73 involve a type of implicit metacognition_(Lau et al., 2022). When we make a similarity judgment, it
74 reflects our own perceptual capacities.

75 The above is a non-trivial prediction, because an alternative hypothesis is that subjective
76 similarity ratings may be made based on whatever visual features that happen to be more salient,
77 depending on one's fluctuating attentional states, or arbitrary preferences that aren't necessarily
78 related to one's own performance in near-threshold psychophysical tasks. This alternative
79 hypothesis is not implausible given that 'error' feedbacks are generally never given to subjects, to
80 'correct' them or train them, as they make these similarity ratings somewhat freely.

81 To test our hypothesis, we quantify the degree of subjective perceptual similarity between
82 stimulus pairs by having participants freely rank similarity, without being given specific criteria,
83 making it subjective, across a stimulus set (Figure 1A & 1B; subjective similarity judgment task). We
84 also assess participants' perceptual discrimination capacity between the pairs. The stimulus pairs
85 may be so obviously dissimilar that discriminating between them is just too easy (i.e. performance
86 under normal conditions would be at the ceiling). To address this problem, we propose to use a
87 psychophysical method to measure such discrimination capacity near perceptual threshold. We
88 measure the participants' discrimination performance within the morph set that spans between
89 the two stimuli (Figure 1C; near-threshold discrimination task). With this, we quantify the number
90 of just-noticeable-differences (#JNDs; see legends of Figure 1C for explanation) between a pair. The
91 #JNDs reflects the perceptual discrimination capacity, with its higher value indicating a higher
92 capacity. We use faces as stimuli in our study due to their high-dimensional (i.e. multi-features)
93 nature, and the fact that these are naturalistic stimuli commonly encountered in everyday life. In
94 subsequent sections, we use the notion "dissimilarity" instead of "similarity", so the hypothesized
95 association with discrimination capacity is positive.

96 In summary, we hypothesize that there is a correlation between perceptual discrimination
97 capacity (in near-threshold tasks) and subjective perceptual dissimilarity (as reflected by self-
98 ratings of supra-threshold stimulus pairs) within each individual (Hypothesis 1, first row in Table
99 1). Specifically, perceptual discrimination capacities are higher in face pairs that are subjectively
100 judged to be more dissimilar. Further – and critically– we hypothesize that this association is
101 specific to each individual (Hypothesis 2, second row in Table 1), meaning that one's subjective
102 perceptual dissimilarity is better explained by one's own perceptual discrimination capacity than
103 other participants' (average) discrimination capacity. This would support the notion that
104 subjective perceptual similarity may be metacognitive in nature, meaning that it concerns one's
105 own perceptual capacities, not just the general physical similarity between stimuli. A complete
106 overview of the hypotheses and their corresponding tests is provided in Table 1.



107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132

Figure 1. Experimental tasks for estimating subjective perceptual dissimilarity and perceptual discrimination capacity. (A) Illustration of 30 faces to be used in the present study. (B) The subjective similarity judgment task for estimating the level of subjective perceptual dissimilarity between face pairs. A target face on top and four other faces (candidates) on the bottom are shown to the participant in each trial. Participants are instructed to rank the candidate faces from the most to least similar with respect to the target face by clicking on them in order. Then, a 30x30 dissimilarity matrix is computed from the participant’s responses across trials, with the value in each cell of the matrix indicating the level of subjective dissimilarity between a face pair. (C) The near-threshold perceptual discrimination task for measuring the discrimination capacity between two faces. One thousand morphs are created as intermediate transitions between two faces (based on a computational face model; see Methods for details). In each trial, three faces are shown simultaneously to the participants. Two are identical, and one is different from the other two by a certain degree (number of morph steps within the 1000-morph series). Participants are instructed to click on the face that is different from the other two. Task difficulty is maintained by titrating the number of morph steps needed for the different face to be barely detectable, using a standard 1-up-2-down staircase method. The converged (i.e., stabilized) value of the staircase indicates the number of morph steps required to maintain near-threshold performance (71% correct); thus, this value reflects the just-noticeable-difference (JND). Because these morph steps come from a series of 1000 morphs between a face pair (e.g., any two faces in 1A), if e.g. JND = 250 morph steps, we can also describe the two faces concerned as being 4 JNDs apart from each other. This general notion of the number of JNDs (#JNDs) between face pairs, which is just the total number of morph steps (1000) divided by the measured JND, allows us to describe the psychophysical discriminability between two faces, free from the non-standard physical unit of ‘morph steps’ (which depends on the

133 arbitrary specifics of the morphing procedure such as total number of steps used).
134 Essentially, #JNDs indicates the perceptual distance between a face pair, in other words, how
135 many JNDs are in between the face pair; thus, its higher value corresponds to a higher
136 discrimination capacity.

137 **Methods**

138 Please note that the method section is written in the present tense as the experiment has not
139 been done yet. We will change the tenses to past tense in the second phase submission.

140 **Ethics information**

141 The study received approval from the Ethics Review Committee at RIKEN, complying with all
142 their ethical guidelines. Informed consent is obtained from participants before the experiment, and
143 in appreciation of their participation, they are compensated with 3000 yen (approximately 20 US
144 dollars) for each day of participation (roughly 90 minutes each day).

145 **Design**

146 Twelve participants are recruited for the study. They initially perform the subjective similarity
147 judgment task, twice over the course of two days. After all participants complete this task twice, a
148 set of 24 face pairs are selected for examination in the near-threshold discrimination task. The
149 criteria for selecting the face pairs are described in the subsequent sections. Then, all participants
150 are invited back to perform the near-threshold discrimination task over two days on these 24 face
151 pairs, randomly splitting the pairs between the days.

152 Note that there are 48 sessions in total, across 12 participants. Each participant performs four
153 sessions on different days with each session taking more than 60 minutes. This provides us with
154 enough data to perform our statistical analysis at the individual-level. The subjective similarity
155 judgment task consists of 300 efficiently crafted trials to estimate the level of subjective perceptual
156 dissimilarity between all face pairs. Participants perform this task twice, and the achieved
157 dissimilarity values are averaged to further enhance the robustness. The near-threshold
158 discrimination task comprises a total of 1440 trials, 60 trials per face pair, to effectively estimate
159 the perceptual discrimination capacity between a systematically selected set of 24 face pairs.
160 Furthermore, we recruit more participants if these initial 12 participants don't satisfy our data
161 collection stopping rule described in the Sampling plan section.

162 *Face data set*

163 The basal face model (BFM) (Paysan et al., 2009) is used to select our face dataset and generate
164 morphs between the faces for the near-threshold discrimination task. BFM is a widely used
165 morphable model for generating graphical faces with two embedded vectors describing the shape
166 and texture of the faces independently. We arbitrarily selected 30 faces from the BFM space, while
167 ensuring a diverse set that also includes faces positioned close to each other in the BFM space. The
168 top three shape dimensions were assigned systematically from a cylindrical coordinate with a 2.5
169 SD radius, and the subsequent top 47 shape dimensions and top 50 texture dimensions were

170 assigned randomly from a uniform distribution ranging between -1.5 and 1.5 SD. The remaining
171 less important shape and texture dimensions were set to zero. Figure 1A shows the selected 30
172 faces for the study.

173 *Subjective similarity judgment task*

174 In each trial, participants are presented with a visual arrangement consisting of one face
175 positioned at the top (target face) and four other faces positioned at the bottom (candidate faces)
176 (Figure 1B). Participants are instructed to rank the four bottom candidate faces based on their
177 perceived similarity to the top target face by mouse-clicking on the faces in the order of most to
178 least similar. Each clicked face immediately disappears from the screen, and the trial ends after all
179 candidate faces are clicked one by one. If participants fail to complete the trial within 30 seconds,
180 the trial is skipped, and any ranking assigned is disregarded. The aim of this task is to estimate the
181 level of subjective dissimilarity between each face pair and to construct a dissimilarity matrix
182 (Figure 1B) for each participant by analyzing their assigned rankings across trials.

183 The level of subjective dissimilarity (dissimilarity value) between two faces is estimated by
184 calculating the probability of one face being ranked lower than the rest of the faces when the other
185 face is the target, as outlined in the following. The rankings given in all trials are segmented into
186 sets of three, consisting of the target face and the combination of two of the four candidate faces
187 (i.e., six sets per trial). Within each set, the face that ranked lower is marked as the odd face.
188 Subsequently, the dissimilarity value between a face pair is determined by calculating the ratio of
189 instances where one of the faces is marked as the odd face across all sets that include the face pair
190 with either of them as the target face. It is noteworthy that when calculating this ratio, we account
191 for a non-tested set with an obvious outcome, where one of the faces repeats, by adding 0.5 to both
192 the numerator and the denominator. This fundamentally prevents getting a dissimilarity value of
193 zero, as only the diagonal value of the dissimilarity matrix should be zero.

194 The tuple of five faces displayed in each trial is strategically selected using the InfoTuple
195 method (Canal et al., 2020). This method guarantees that each trial offers informative data, thereby
196 enhancing the estimation of the dissimilarity matrix. This essentially enables achieving a robust
197 estimation of the dissimilarity matrix over a smaller number of trials. The trial selection procedure
198 is similar to the one used by Canal et al. (Canal et al., 2020) and comprises the following steps:

- 199 1. The tuple set in the first 30 trials is selected at random while ensuring that each face is
200 selected once as the target face.
- 201 2. The dissimilarity matrix is calculated as described above, and a 5-dimensional metric
202 multidimensional scaling (MDS) (Borg & Groenen, 2005) is applied to the dissimilarity matrix
203 to find its embeddings.
- 204 3. A cycle of 30 trials, showcasing each face once as the target face, is selected by the InfoTuple
205 method using the embeddings. The InfoTuple method selects the tuple that maximizes a
206 mutual information estimate which involves two entropy terms: intuitively, one term favors
207 tuples whose rankings are uncertain given the current embeddings, while the other
208 discourages inherently ambiguous tuples that are expected to remain uncertain even if the
209 embeddings are revealed. So, it aims to select an informative tuple whose rankings are

210 unknown but yet can be answered reliably (consistently). Please refer to the original paper
211 for a detailed explanation and the mathematics of the InfoTuple (Canal et al., 2020).

212 4. The dissimilarity matrix is calculated given all the data collected, and the embeddings are
213 updated by applying a 5-dimensional metric MDS to the dissimilarity matrix and using the
214 previous embeddings as the seed in the MDS algorithm.

215 5. Steps 3 and 4 are repeated for 9 iterations. We stop after 9 iterations as the dissimilarity
216 matrix and embeddings reach a relatively stable state at this point. As completing 9
217 iterations is lengthy and can be exhausting, participants are given a short break after every
218 3 iterations.

219 The final obtained embeddings are used to recalculate the dissimilarity matrix by computing
220 the Euclidean distances between the faces. This process results in a more accurate version of the
221 dissimilarity matrix by refining inaccuracies in some cells of the original dissimilarity matrix due to
222 insufficient data attributed to them or due to noise in responses (i.e., response inconsistencies).
223 Therefore, this dissimilarity matrix, derived from the embeddings, is utilized in all the subsequent
224 stages instead of the original dissimilarity matrix. Please see the supplementary Figure 1 for a
225 schematic overview of the task design described above.

226 We also plan to test a 2-dimensional metric MDS for recalculating the dissimilarity matrix. Even
227 though there might be significant information loss in using such a low dimensional MDS, it can
228 further refine the matrix and also make the dissimilarity relations sharper (i.e., more distinct). A
229 sharp dissimilarity relation can potentially make its association with perceptual discrimination
230 capacity more salient. This was true in our pilot data. Both Hypothesis 1 and 2 (Figure 2) reached a
231 $p < 0.05$ at the individual-level in all four participants.

232 We note that in our algorithm, prior to applying a metric MDS, a nonmetric MDS is employed to
233 fill in the missing cells of the dissimilarity matrix (i.e., pairs with no ranking data). The missing cells
234 are filled in by the Euclidean distances computed from the embeddings derived by the nonmetric
235 MDS. This procedure is important in initial iterations in which there are a considerable number of
236 missing cells. Then, the metric MDS is applied to this filled-in dissimilarity matrix. We don't use the
237 nonmetric MDS directly because the nonmetric MDS, unlike the metric MDS, doesn't preserve the
238 magnitude of the dissimilarity between pairs.

239 Last, it is important to mention that each participant performs the above task twice, each time
240 on a different day. The average of the dissimilarity matrices obtained from each day forms the final
241 dissimilarity matrix. Further, to evaluate the reliability of the obtained dissimilarity matrix from a
242 session, we report the within-participant correlation between the dissimilarity matrices derived
243 from each day. As a reference, we also report the distribution of between-participant correlation
244 by randomly correlating a matrix from a session in one participant with that of another participant.
245 We expect that the within-participant correlation to be higher than the between-participant
246 correlation.

247 In addition to the above approach for deriving the embeddings and the dissimilarity matrix, we
248 plan to try a machine learning approach as an exploratory analysis (i.e., we don't use this method
249 in our main hypotheses testing and the stopping criterion). This approach starts with random
250 embeddings and iteratively updates them to minimize a loss function, which penalizes wrong

251 similarity rankings derived from the embeddings. The loss is constructed using a sigmoid activation
 252 function in a binary cross-entropy as follows:

$$253 \quad p = \frac{1}{1 + e^{-k(d_{dissim} - d_{sim})}}, \quad d_{sim} = \|x_{target} - x_{sim}\|_2, \quad d_{dissim} = \|x_{target} - x_{dissim}\|_2$$

$$254 \quad L = -\frac{1}{N} \sum_{i=1}^N \log(p_i) + \lambda_1 \sum_{m=1}^M \|x_m\|_1 + \lambda_2 \sum_{m=1}^M \|x_m\|_1^2$$

255 Where x_m represents the vector embedding of face m ; d_{sim} and d_{dissim} are the Euclidean
 256 distances between a target face and a face ranked as more similar and a face ranked as less similar
 257 to the target face by the participant, respectively; k corresponds to the ranking difference (e.g., 2
 258 for a face ranked first and a face ranked third), putting more emphasis on clearer similarity
 259 comparisons; N indicates the number of segmented trio comparisons (with six trio segments in a
 260 trial); λ_1 and λ_2 are the hyperparameters of L1 and L2 regularizations which help to control the
 261 sparsity and scale of the embeddings. We use the Keras library in Python, with Adam optimizer, to
 262 minimize the loss function.

263 We note that we don't use this machine learning approach during the task (i.e., in our online
 264 application) because it is slower, requiring cross-validations and careful selection of the
 265 hyperparameters. This approach is more sophisticated than our main approach, which involves
 266 estimating probabilities and running MDS, but it has the potential to yield a better estimation of
 267 the embeddings and the dissimilarity matrix. In our pilot study, using this approach, we got similar
 268 results to those shown in Figures 2 and 3.

269 *Near-threshold discrimination task*

270 The objective of this task is to estimate perceptual discrimination capacity in a face pair. A series
 271 of 1000 equally spaced morphs are generated along the line connecting a face pair in the BFM
 272 space. In each trial, three faces are shown to the participants: two identical faces, randomly
 273 selected from the morph set, and a third different face spaced by a certain number of morphs (step
 274 value) from the identical faces (e.g., morph 200, morph 200, and morph 300: here the step value is
 275 100). The faces are displayed simultaneously at the center of the screen next to each other, and
 276 their arrangement is randomized in each trial (Figure 1C). Participants are instructed to identify
 277 and click on the different face.

278 A staircase (Cornsweet, 1962) with a 1-up and 2-down protocol is applied to the step value (i.e.,
 279 the number of morphs between the different and identical faces), initiating from a step value of
 280 500. After each incorrect response trial, the step value is increased, and it is decreased after two
 281 consecutive correct trials. The magnitude of the change in the step value gradually decreases over
 282 trials, reaching a minimum change of 20 steps. The task is terminated after 60 trials, allowing
 283 precise convergence of the step value. The converged step value indicates the just-noticeable-
 284 difference (JND), the minimum degree of differences between the faces required to achieve near-
 285 threshold discrimination performance (71% correct response). The average of the steps achieved
 286 within the last 5 changes is defined as the converged step. Essentially, a small JND, for example,
 287 100, indicates that the two questioned faces are quite distinct, involving 10 JNDs (i.e., 1000 divided
 288 by 100) between them. We use the notion of the number of JNDs (#JNDs) to quantify perceptual

289 discrimination capacity. The #JNDs indicates the perceptual distance between a face pair, in other
290 words, how many JNDs are in between the face pair in a participant. Therefore, its higher value
291 corresponds to a higher discrimination capacity. The #JNDs is simply calculated as 1000 (i.e., total
292 number of morphs) divided by the JND.

293 In a session, there are 12 face pairs to undergo the staircase procedure. The staircases for each
294 of these face pairs are interleaved, progressing concurrently. There is a cycle of 12 trials, featuring
295 each staircase once in a random order. The trials are time-constrained, requiring participants to
296 respond within 8 seconds. If participants fail to respond within this time window, the trial is skipped
297 and reintroduced at the end of the cycle. To encourage participants to perform to the best of their
298 abilities, they are provided with feedback on their responses. A green circle is displayed on the
299 different face (i.e., correct answer) and a red cross on the identical faces (i.e., wrong answers) after
300 they provide their response. Since the session is lengthy, with a total of 720 trials, participants are
301 given a short break after every 180 trials.

302 The trajectory convergence of a staircase could indicate the reliability of the estimated #JNDs.
303 A staircase with a higher ratio of reversals in its later trials could be considered more reliable.
304 Therefore, we report the ratio of reversals in the last 20 trials of each pair's staircase and its
305 statistics across participants. In an absolute ideal case, given our 1-up and 2-down staircase
306 protocol, the ratio of reversals in the last 20 trials would be 0.6.

307 *Selection of the pairs for the near-threshold discrimination task*

308 Following the completion of the subjective similarity judgment task twice by all 12 participants,
309 24 face pairs are systematically selected to be examined in the near-threshold discrimination task.
310 Practical constraints (time limitations; it takes 4-5 minutes to complete the near-threshold
311 discrimination task for a face pair) limit us to examine only a small subset of the pairs. A sample
312 size of 24 pairs should be fairly adequate for detecting an effect, and it is further justified by our
313 pilot study, as we achieved reasonably robust results by examining only 13 pairs, almost half of our
314 planned 24 pairs. Participants are re-invited for two sessions to perform the near-threshold
315 discrimination task on these specific 24 pairs, completing 12 of the pairs in each session.

316 Measuring perceptual discrimination capacity, expressed as #JNDs, in a face pair, involves
317 running a near-threshold discrimination task dedicated to that specific pair. Thus, considering
318 practical constraints as described earlier, we have no choice but to examine only a limited subset
319 of pairs (24 out of 435 possible pairs). However, this subset is carefully chosen to provide the most
320 informative data for testing our hypotheses while fairly covering different ranges in the group-
321 averaged dissimilarity matrix. The pairs with a controversial subjective dissimilarity degree across
322 participants are particularly promising candidates. If the hypothesis holds true, these pairs should
323 also exhibit controversial discrimination capacity across participants. Considering the inherent
324 noise in our methods in estimating the dissimilarity values and the #JNDs, any effect should be
325 better detectable on pairs with larger standard deviations, those that are more distinct, across
326 participants. So, we select 18 pairs with controversial dissimilarity values across participants, and
327 for the sake of comparison, we select 6 pairs with less (non) controversial values.

328 First, the dissimilarity matrix is z-normalized within each participant to ensure that its scale is
329 consistent across participants. Subsequently, the mean and SD of the dissimilarity matrix are
330 computed across participants, and the quantiles of the mean values are derived. Within the first
331 and the last quantiles, 3 pairs with the highest and 1 with the lowest SD are selected. Additionally,
332 6 pairs with the highest and 2 pairs with the lowest SD are chosen within the second and the third
333 quantiles. This systematic selection ensures choosing 18 controversial and 6 non-controversial
334 pairs that cover a diverse range in the group-averaged dissimilarity matrix.

335 **Sampling plan**

336 Participants who meet the following criteria are excluded from the analysis: Those who don't
337 complete all four experimental sessions and those who show a lack of attentiveness to the task in
338 any of the sessions. The lack of attentiveness in the near-threshold discrimination task is identified
339 by non-converging staircases, indicated by a non-fluctuating increment in the step value over
340 trials. Specifically, a session in which there are more than 4 (out of 12) staircases with less than
341 three downs in their last 20 trials is considered bad with lacking sufficient attentiveness. In the
342 subjective similarity judgment task, the lack of attentiveness is judged by comparing the
343 consistency of responses between the first and second half of the session. Specifically, if the
344 correlation between the dissimilarity matrices estimated from each half falls below 0.2, the session
345 is considered bad with inadequate attentiveness. This correlation was 0.56 ± 0.086 (mean \pm SD) in
346 our pilot data. The data collection continues until we have 12 participants who successfully
347 complete the experiment without meeting any of the exclusion criteria. If a participant meets the
348 exclusion criteria, a new participant is recruited to replace the excluded participant. Note that the
349 second phase of the study, involving the selection of the pairs and the near-threshold
350 discrimination task, does not start until the quality of the data from all 12 participants in the
351 subjective similarity judgment task is confirmed as not meeting the above mentioned exclusion
352 criteria. Following this, any subsequent exclusion and recruitment of new participants do not
353 modify the initially selected pairs for the near-threshold discrimination task.

354 After analyzing the data from these 12 participants, if the statistics fail to meet the following
355 stopping criterion, we recruit more participants until the criterion is satisfied. The individual-level
356 statistic is converted to z-values, and the 95% confidence interval of the group-mean z-value is
357 derived (See the Analysis Plan section). We stop the experiment, if, in both Hypothesis 1 and 2, the
358 width of this 95% confidence interval is less than 1. Moreover, we consider a hypothesis to be
359 confirmed, if the group-mean z-value is significantly above zero, specifically, if the 95% confidence
360 interval is above zero. Note that we set our stopping criterion independent of the significance
361 testing and solely based on the precision of the effect (i.e., the confidence interval). We do not stop
362 our experiment until achieving a high precision, so that we are confident that the effect is not being
363 confirmed or rejected because of some extreme observations (Cumming, 2008; Lakens, 2014) .
364 Given our sample size scale, we expect a considerable effect to have a group-mean z-value of at
365 least above 0.5. So, a minimally significant scenario involves a group-mean z-value of 0.5 with a
366 95% confidence interval width of less than 1. Considering this, we set our stopping criterion as the

367 width of the 95% confidence interval being smaller than 1. Thus, this seems like enough precision
368 to safely reject or accept a hypothesis.

369 We note that after our initial 12 participants, we recruit three more participants, each time the
370 stopping criterion is not met. We repeat this until reaching a maximum of 24 participants. Given
371 that our pilot data with only four participants show a 95% confidence interval with a width of
372 around 1.5 (see Figure 2B & 2D), it is unlikely not meeting the stopping criterion before reaching
373 our maximum sample size of 24 (see the supplementary Figure 2). It is also worth noting that the
374 recruitment of new participants does not alter the pairs used in the near-threshold discrimination
375 task. The newly recruited participants perform the task on the same pairs selected based on our
376 initial 12 participants.

377 Analysis plan

378 Hypothesis 1

379 Spearman correlation coefficient and its p-value are computed between the dissimilarity
380 values and #JNDs of the examined 24 pairs in each individual. The Spearman correlations are
381 converted to z-values using the Fisher z-transformation (Fieller et al., 1957) to conduct group-level
382 statistical tests. The distribution of the group-mean z-value is computed by bootstrapping, iterated
383 100,000 times, and then its 95% confidence interval is derived by obtaining the 2.5th and 97.5th
384 percentile of the distribution. The hypothesis is confirmed if this confidence interval is above zero.
385 The following statistics are reported as complementary information: the p-value and the Bayes
386 factor of a t-test applied to z-values (BayesFactor Matlab package is used:
387 <https://zenodo.org/badge/latestdoi/162604707>), the p-value of a Fisher's combined probability
388 test, combining individual-level p-values (Brown, 1975), and a Bayesian posterior distribution of
389 population prevalence (Ince et al., 2021) and its 95% highest posterior density interval, considering
390 the p-value of 0.05 as the individual-level significance threshold. The Bayesian posterior
391 distribution quantitatively summarizes how prevalent a particular effect would be in the
392 population, based on the number of participants tested in a study and their proportion showing
393 the effect significantly. Furthermore, we note that if the 95% confidence interval of the group-mean
394 z-value includes zero (i.e., a non-significant finding), we incorporate an equivalence test (Lakens et
395 al., 2018) to further assess the significance of the null result (i.e., how confidently the effect can be
396 considered non-significant). We assume a group-mean z-value range of -0.5 to 0.5 as the smallest
397 effect size of interest; thus, if the 95% confidence interval falls within this range, we regard the
398 effect to be confidently null.

399 Hypothesis 2

400 Each participant's #JNDs is z-normalized to ensure that the #JNDs range is consistent across
401 participants. This normalization is crucial, given that some participants may exhibit generally
402 higher #JNDs than others. Subsequently, a nonparametric permutation test is applied to each
403 individual to assess the specificity of the relationship between their #JNDs and dissimilarity values,
404 as follows:

- 405 1. A permutation set of #JNDs is constructed by randomly permuting the #JNDs across
406 participants (i.e., for each pair, selecting the value in one of the participants at random),
407 excluding the participant in question. Essentially, the permutation set simulates a new
408 participant by mixing the existing participants.
- 409 2. The Spearman correlation coefficient is calculated between the permuted #JNDs and the
410 dissimilarity values of the participant in question. It is noteworthy that with 12 participants
411 and 24 pairs, there are an enormous number of possible permutations (i.e., 11^{24} unique
412 permutations), which ensure constructing a reliable null distribution.
- 413 3. Steps 1 and 2 are repeated 100,000 times to derive the distribution of the Spearman
414 coefficients. This distribution represents the null hypothesis distribution in which there is
415 no individual specificity.
- 416 4. The actual Spearman coefficient between the #JNDs and dissimilarity values of the
417 participant in question is tested against the null distribution, and the p-value, indicating
418 the significance level, is derived. The z-value is also calculated by subtracting the actual
419 Spearman correlation from the null distribution's mean and then dividing it by the null
420 distribution's SD.

421 Then, similar to Hypothesis 1, the 95% confidence interval of the group-mean z-value is derived
422 through bootstrapping, and if it is above zero, the hypothesis is confirmed. The complementary
423 statistics, outlined in Hypothesis 1, are also reported for Hypothesis 2.

Table 1 Experimental Design Table

Question	Hypothesis	Outcome Measures	Sampling plan	Analysis Plan	Interpretation given to different outcomes
How does perceptual discrimination capacity relate to subjective perceptual dissimilarity?	There are higher perceptual discrimination capacities in pairs that are perceptually more dissimilar. In other words, we hypothesize that perceptual discrimination capacity is positively correlated with subjective perceptual dissimilarity.	Subjective perceptual dissimilarity and perceptual discrimination capacity between stimuli pairs are measured in each participant through two different psychophysical tasks called the subjective similarity judgment task (Figure 1B) and the near-threshold discrimination task (Figure 1C). The subjective similarity judgment task assesses the level of subjective dissimilarity (dissimilarity value) between stimulus pairs. The near-threshold discrimination task measures the number of just-noticeable-differences (#JNDs) between stimulus pairs, quantifying the perceptual discrimination capacity. A higher #JNDs indicates a higher capacity in distinguishing a stimulus pair.	Twelve participants are recruited, each completing four sessions over four days, spending two days on the subjective similarity judgment task and two days on the near-threshold discrimination task. Participants failing to complete all four sessions and those displaying a lack of attentiveness to the task in any of the sessions are excluded from the analysis. The criteria for a lack of attentiveness are described in the method section. After completing the data collection on 12 participants who don't meet any of the exclusion criteria, if the results don't satisfy our following experiment's stopping criterion, we recruit more participants until the criterion is satisfied. However, we end the experiment once we reach a maximum of 24 participants, regardless. The stopping criterion is met if the width of the 95% confidence interval of the group-mean z-value (see the method section) is less than 1, and the hypothesis is confirmed if this 95% confidence interval is above zero. * The hypothesis was confirmed in the pilot study (Figure 2A & 2B)	In each participant, the Spearman correlation between the #JNDs and dissimilarity values is computed. Then, the Spearman correlations are transformed into z-values using the Fisher z-transformation. Subsequently, the 95% confidence interval of the group-mean z-value is derived through bootstrapping with 100,000 iterations.	If the described group test doesn't reach the level of significance, yet at the individual level, the correlation reaches significance ($p < 0.05$) within a certain few participants, we can interpret that the positive association between perceptual discrimination capacity and subjective perceptual dissimilarity holds true for certain individuals and doesn't generalize to the entire population. Failure of the group test may indicate that subjective perceptual dissimilarity is made rather arbitrarily, based on subjective preferences, and does not reflect underlying psychophysical capacities.
Is the association between perceptual discrimination capacity and subjective perceptual dissimilarity specific to each individual?	A participant's subjective perceptual dissimilarity is better explained by their own perceptual discrimination capacity than by a group-averaged perceptual discrimination capacity. To put it differently, we hypothesize that the positive association between perceptual discrimination capacity and subjective perceptual dissimilarity is specific to each individual. Essentially, subjective dissimilarity reflects a metacognitive assessment of one's own perceptual discrimination capacity, rather than general knowledge about the physical differences of the stimuli.	The level of individual-specificity of the relationship between perceptual discrimination capacity and subjective perceptual dissimilarity is measured by a nonparametric permutation test, which assesses whether one's own dissimilarity value is more strongly correlated with one's own #JNDs than with others' #JNDs. Essentially, the above test reflects how specific this relationship is in each individual in terms of z-values.	The participants' exclusion and the experiment's stopping strategies remain the same as above. Similar to the above hypothesis, the hypothesis is confirmed if the 95% confidence interval of the group-mean z-value is above zero, and the experiment stopping rule is met if the width of this confidence interval is smaller than 1. * The hypothesis was confirmed in the pilot study (Figure 2C & 2D)	First, the #JNDs is z-normalized within each individual to ensure that its scale is consistent across participants, and then a nonparametric permutation test is applied to each participant separately as follows: Briefly, a permutation #JNDs set is constructed by randomly permuting the #JNDs of all participants, excluding the participant in question. The Spearman coefficient is then computed between the permuted #JNDs set and the dissimilarity values of the participant in question. This process is iterated 100,000 times to establish the distribution of the Spearman coefficient, representing the null hypothesis distribution. Finally, the actual Spearman coefficient between the dissimilarity values and the #JNDs of the participant in question is compared against the null distribution, and its z-value is computed. Subsequently, the 95% confidence interval of the group-mean z-value is derived by bootstrapping, similar to the above hypothesis.	If the described group test fails to reach the significance level, but there are certain individuals with significant statistics ($p < 0.05$), we can draw a similar interpretation as the above that the hypothesis holds true only for certain people. Failure of the group test, if the first hypothesis holds true, may suggest that subjective perceptual dissimilarity is made based on general stimulus properties that predict the psychophysical performance of human subjects and is not metacognitive in the sense of reflecting direct access to one's own perceptual capacities.

Pilot data

427 We conducted a pilot version of the study with 4 participants. The experimental design was
428 similar to the currently proposed experiment (Figure 1), with a few differences. Participants
429 performed the subjective similarity judgment task only once. After all of them completed the task,
430 13 pairs of faces were selected, and the participants were all invited back to perform, in a different
431 session, the near-threshold discrimination task on these pairs. In this pilot study, we randomly
432 selected the pairs while ensuring that most of them have high dissimilarity values SD across
433 participants, indicating that the degree of subjective dissimilarity is quite ‘controversial’, i.e.,
434 individual-specific. The subjective similarity judgment task and the near-threshold discrimination
435 task were conducted similarly to the proposed experiment, except that there was no time
436 constraint on both tasks, and no trial-by-trial behavioral feedback was provided during the near-
437 threshold discrimination task. Additionally, in one participant, the near-threshold discrimination
438 task comprised 50 instead of 60 trials. We applied the same analysis approach described in the
439 method section on the pilot data.

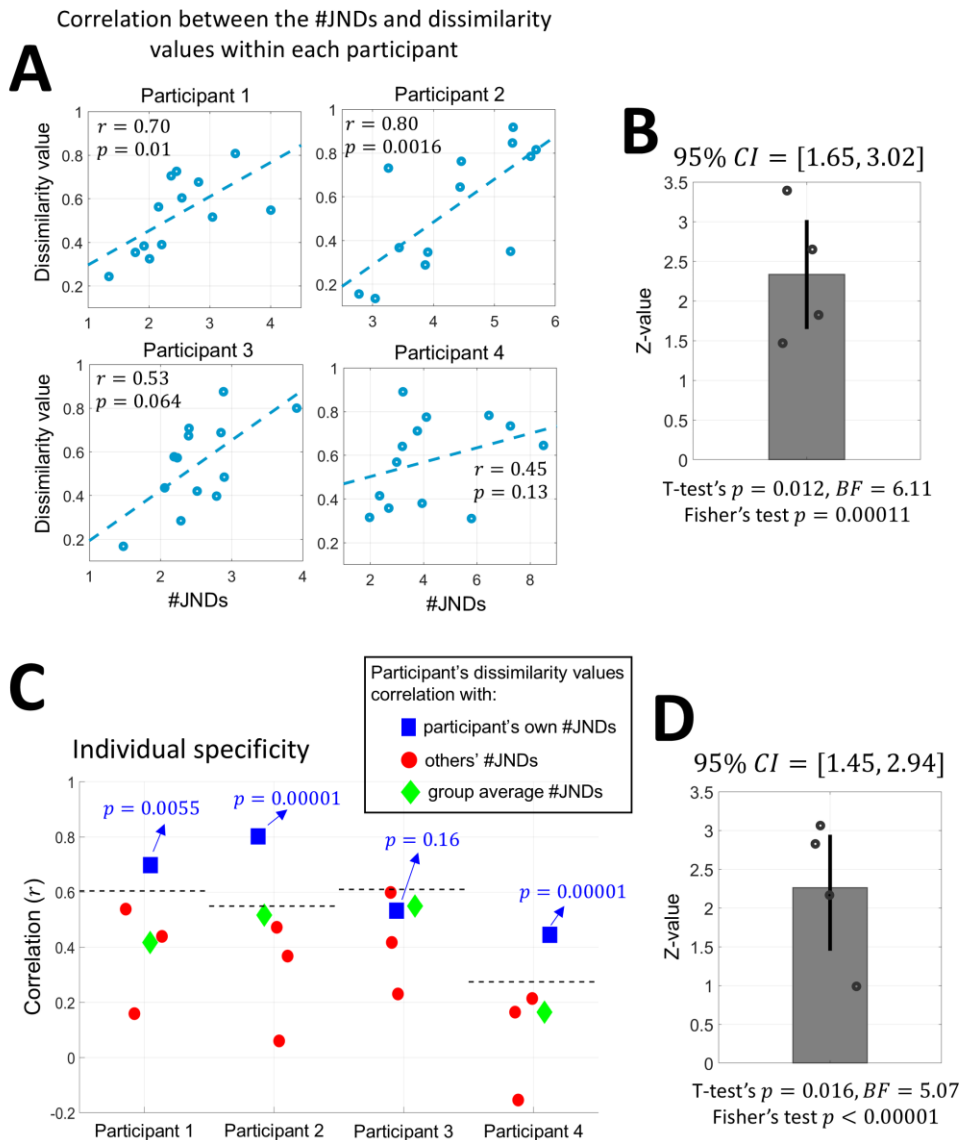
440 Hypothesis 1 was confirmed in the pilot study (Figure 2A & 2B), indicating that subjective
441 perceptual dissimilarity and perceptual discrimination capacity are highly correlated. The
442 correlation was significant ($p < 0.05$) at the individual-level in 2 out of 4 participants. At the group
443 level, the mean z-value across participants was 2.33, with a 95% confidence interval between 1.65
444 and 3.02. A t-test on the z-values yielded a p-value of 0.012, and a Bayes factor (BF) of 6.11.
445 Moreover, a Fisher’s test, combining the individual-level p-values, resulted in a p-value of 0.00011.

446 More importantly, Hypothesis 2 was also confirmed in the pilot study (Figure 2C & 2D),
447 suggesting that the association between subjective perceptual dissimilarity and perceptual
448 discrimination capacity is specific to each individual. To put it differently, others’ perceptual
449 discrimination capacity cannot account for one’s subjective perceptual dissimilarity as well as their
450 own perceptual discrimination capacity. The statistic was highly significant ($p < 0.05$) at the
451 individual-level in 3 out of 4 participants. At the group-level, the mean z-value across participants
452 was 2.26, with a 95% confidence interval between 1.45 and 2.94. A t-test on the z-values resulted in
453 a p-value of 0.016 and a BF of 5.07, and a Fisher’s test yielded a p-value smaller than 0.00001.

454 In the main study, we anticipate observing even stronger statistics not only at the group-level
455 but also at the individual-level. We expect that testing more stimulus pairs and having more
456 participants lead to observing stronger results at the individual-level for Hypothesis 1 and
457 Hypothesis 2, respectively.

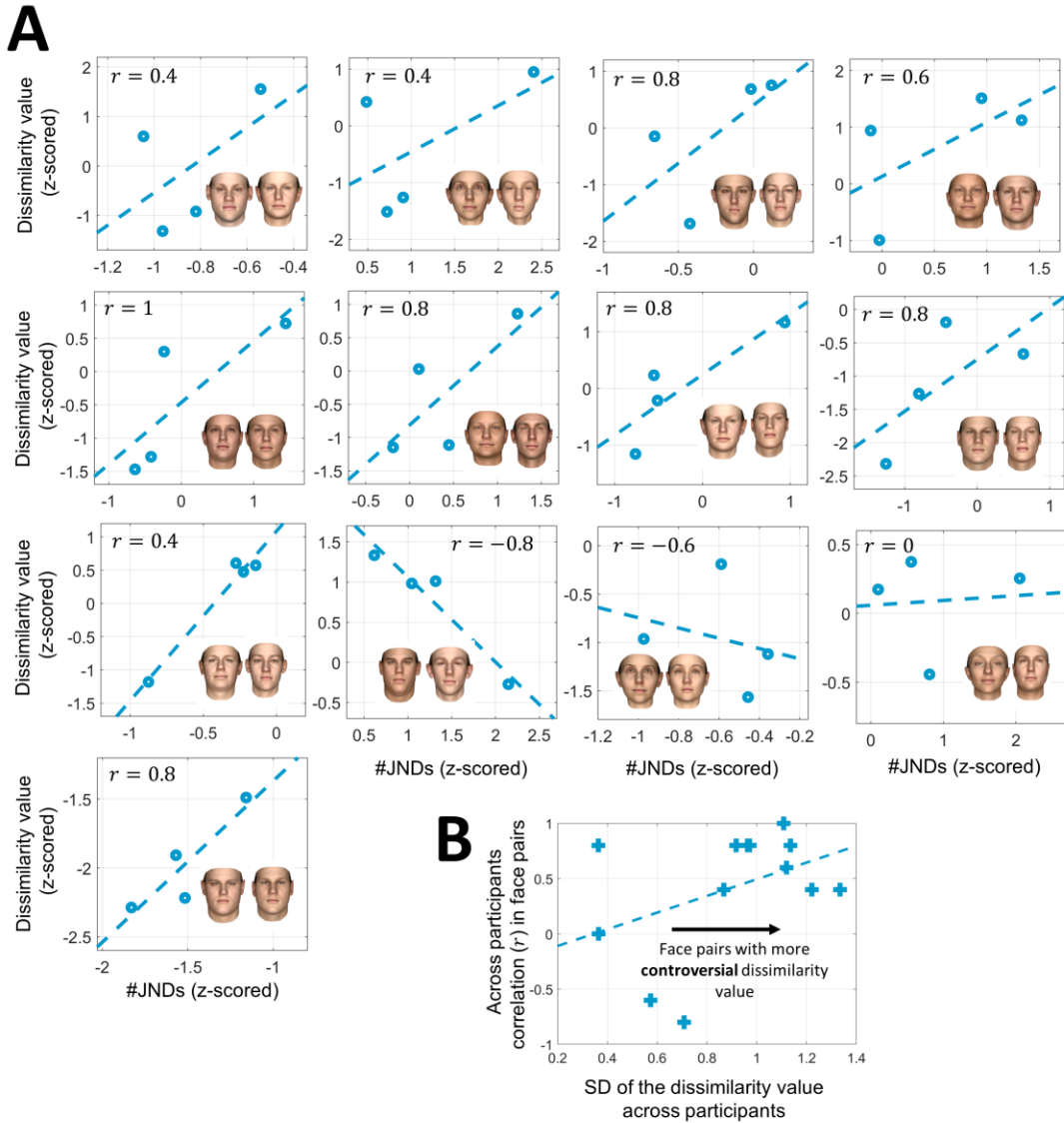
458 We further explored the correlation between subjective perceptual dissimilarity and perceptual
459 discrimination capacity in each face pair across participants in our pilot study (Figure 3). Given the
460 small sample size (i.e., four participants), no meaningful statistical conclusions can be inferred.
461 However, it is notable that the correlations were strongly positive, particularly in the controversial
462 pairs: those with controversial degrees of subjective dissimilarity across participants. Essentially,
463 Figure 3 also indicates that one’s perceptual discrimination capacity can explain one’s subjective
464 perceptual dissimilarity. Similar inter-individual differences observable in the subjective
465 perceptual dissimilarity could also be found in the perceptual discrimination capacity. However,

466 this is not apparent in the less controversial pairs. Nonetheless, this may not necessarily suggest
 467 that the association doesn't exist in the less controversial pairs. The measures obtained from our
 468 psychophysical tasks inevitably contain some noise which may make them to be not precise
 469 enough to capture the subtle differences across participants in the less controversial pairs. In the
 470 main experiment, we expect to obtain a clearer picture by having more participants and
 471 experimental sessions.



472 **Figure 2. Relationship between perceptual discrimination capacity and subjective**
 473 **perceptual dissimilarity in the pilot study (N = 4 participants) at the individual**
 474 **participant level.** (A) Correlation between the perceptual discrimination capacity expressed
 475 by #JNDs and the subjective perceptual dissimilarity within each participant. Each subplot
 476 illustrates the correlation in an individual participant, with each data point corresponding to
 477 a face pair. r indicates the Spearman correlation coefficient, and p denotes its associated p -
 478 value. (B) Individuals' Spearman correlations of A were transformed to z-values for group-
 479

480 level hypothesis testing. The bar plot shows the mean z-value across participants, the vertical
481 line represents the 95% confidence interval of the group-mean z-value, calculated through
482 bootstrapping, and each dot on the plot corresponds to a participant. As displayed under the
483 plot, analysis of the group-level effect by applying a t-test on the z-values yielded a p-value
484 of 0.012, and a Bayes factor (BF) of 6.11. Alternatively, employing a Fisher's combined
485 probability test, combining the individual-level p-values, resulted in a p-value of 0.00011. (C)
486 Individual specificity of the relationship between perceptual discrimination capacity and
487 subjective perceptual dissimilarity. Blue square, red circle, and green diamond indicate the
488 Spearman correlation coefficient between each participant's dissimilarity values and the
489 participant's own #JNDs, other participants' #JNDs, and the group averaged #JNDs,
490 respectively. The dotted horizontal black line denotes the Spearman correlation value
491 corresponding to a permutation test's p-value of 0.05, rejecting the null hypothesis and
492 indicating that the correlation is specific to each individual. The correlation would not be
493 specific to each individual (i.e., null hypothesis) if one participant's dissimilarity values are as
494 equally correlated to the other participants' #JNDs as the own participant's #JNDs. The p-
495 value rejecting the null hypothesis in each participant is shown in blue at the top of the blue
496 squares. The result of the permutation test suggests that the relationship between
497 perceptual discrimination capacity and subjective perceptual dissimilarity was highly
498 specific to each individual in three out of four participants. (D) Individuals' specificity
499 statistics of C were converted to z-values for group-level hypothesis testing. The remaining
500 descriptions of the plot are similar to those in B.



501
 502
 503 **Figure 3. Relationship between perceptual discrimination capacity and subjective**
 504 **perceptual dissimilarity across participants (pilot study N = 4 participants).** (A) Across
 505 participants' correlation between the #JNDs and the dissimilarity values in different face
 506 pairs. Each panel corresponds to a face pair, sorted based on their controversy in the level of
 507 subjective dissimilarity across participants. The top left panel shows the most controversial
 508 pair (i.e., one with the highest dissimilarity value SD across participants), and the bottom
 509 right panel showcases the least controversial pair. Both the #JNDs and the dissimilarity
 510 values were z-normalized within each participant. r indicates the Spearman correlation
 511 coefficient, and each dot on the plots corresponds to a participant. (B) Relationship between
 512 the face pairs Spearman correlation coefficient shown in A and their controversy in the level
 513 of subjective dissimilarity across participants (i.e., SD of the dissimilarity values across
 514 participants). The relationship between perceptual discrimination capacity and subjective
 perceptual dissimilarity across participants was more salient in highly controversial pairs.

Discussion

515

516 In the present study, we use a near-threshold psychophysical task to quantify perceptual
517 discrimination capacity, which indicates one's capability to distinguish two stimuli (Figure 1C). We
518 aim to examine whether this perceptual discrimination capacity measured at near-threshold is
519 associated with subjective perceptual similarity rankings (Figure 1B) given at suprathreshold
520 (Hypothesis 1). More critically, we seek to explore whether this association is specific to each
521 individual, meaning that one's perceptual discrimination capacity can best explain one's own
522 subjective perceptual similarity compared to that of others' (Hypothesis 2).

523 We conducted a pilot version of the study and confirmed both Hypotheses on our pilot data
524 (Figure 2). However, to further and more precisely investigate our hypotheses, we intend to
525 conduct a larger-scale study with more participants and experimental sessions. Given the high
526 significance level observed in our pilot data, we expect a high likelihood of confirming the
527 hypotheses in the main experiment.

528 If our hypotheses hold true, it may suggest that subjective similarity judgment is, in a specific
529 sense, metacognitive: The self-knowledge of one's perceptual capacity guides one's subjective
530 similarity judgment, and this may occur automatically and implicitly. In essence, perceptual
531 discrimination capacities serve as a ground truth basis for making similarity judgments. A more
532 accurate perceptual similarity judgment could be defined as the one with a more precise
533 metacognitive read-out of one's own perceptual discrimination capacities. Similarly, the instability
534 in perceptual similarity judgments could be considered as the result of inaccurate metacognitive
535 assessment of one's own perceptual capacities.

536 Consequently, higher cortical brain areas, particularly the prefrontal cortex, may play a critical
537 role in perceptual similarity judgments, given that its activity has been demonstrated to be
538 associated with perceptual metacognition (McCurdy et al., 2013; Fleming et al., 2014; Morales et al.,
539 2018). Of course, the current study does not directly test this hypothesis about neural mechanisms.
540 Others have suggested that perceptual similarity information resides within the sensory cortices
541 (Malach, 2021). In light of this, we are currently investigating whether perceptual similarity
542 representations can be found beyond the visual areas, such as the lateral prefrontal cortex, using
543 fMRI.

544 Finally, if our hypotheses are correct, perhaps it could shed light on one conundrum regarding
545 large language models and consciousness. Recently, it has been reported that these models built
546 with current technology in artificial intelligence can give human-like similarity ratings (Kawakita et
547 al., 2023; Marjeh et al., 2023). If the qualitative characters of conscious perception are determined
548 by the relevant similarity relations, as some researchers assume (Clark, 2000; Rosenthal, 2010;
549 Malach, 2021; Lau et al., 2022; Tallon-Baudry, 2022; Zeleznikow-Johnston et al., 2023;
550 Moharramipour & Lau, 2024), does it mean that these artificial agents are conscious
551 (Moharramipour & Lau, 2024)? Or, at least, does it mean that they contain the essential information
552 that is encapsulated within human perceptual experiences? The answer is probably no, if the
553 metacognitive perspective described above is correct. That is, for the similarity judgment to be
554 relevant for subjective experiences, according to our hypothesis, they need to reflect one's own

555 perceptual capacities. What these models do is simply to mimic what humans say in general, and
556 as such, their similarity judgments at best reflect common world knowledge about the physical
557 characteristics of the stimuli, but they are not about one's own perceptual capacities (of which
558 these models have none). There is, thus, a critical difference between humans and those models,
559 in terms of what the similarity judgments mean for them.

560

561

Author contributions

562 A.M. contributed to conceptualization, project planning and design, methodology application,
563 data collection, data analysis, visualization, writing, review, and editing.

564 W.Z. contributed to methodology application, data analysis, review, and editing.

565 D.R. contributed to conceptualization, project design, review, and editing.

566 H.L. contributed to conceptualization, supervision, project planning, design, analysis, review,
567 and editing.

568

Funding

569 This study is supported by internal funding at the RIKEN Center for Brain Science (to H.L.).

570 The funders had no role in study design, data collection and analysis, decision to publish, or
571 preparation of the manuscript.

572

Conflict of interest disclosure

573 The authors declare no competing interests.

574

Data, scripts, code, and supplementary information availability

575 The pilot data, the code used in the pilot study and the code that will be used in the main
576 experiment are publicly accessible from the GitHub repository below:

577 <https://github.com/AliMoharramipour/Subjective-Dissimilarity-and-Discrimination-Capacity->

578

References

579 Borg I, Groenen PJ (2005) *Modern multidimensional scaling: Theory and applications*. Springer

580 Science & Business Media.

581 Brown MB (1975) 400: A method for combining non-independent, one-sided tests of

582 significance. *Biometrics*, 987–992.

583 Canal G, Fenu S, Rozell C (2020) Active ordinal querying for tuplewise similarity learning. In:
584 pp. 3332–3340.

585 Clark A (2000) *A theory of sentience*. Clarendon press.

586 Cornsweet TN (1962) The staircase-method in psychophysics. *The American journal of*
587 *psychology*, **75**, 485–491.

588 Cumming G (2008) Replication and p intervals: p values predict the future only vaguely, but
589 confidence intervals do much better. *Perspectives on psychological science*, **3**, 286–
590 300.

591 Fieller EC, Hartley HO, Pearson ES (1957) Tests for rank correlation coefficients. I. *Biometrika*,
592 **44**, 470–481.

593 Fleming SM, Ryu J, Golfinos JG, Blackmon KE (2014) Domain-specific impairment in
594 metacognitive accuracy following anterior prefrontal lesions. *Brain*, **137**, 2811–2822.

595 Goodman N (1972) Seven strictures on similarity.

596 Hebart MN, Zheng CY, Pereira F, Baker CI (2020) Revealing the multidimensional mental
597 representations of natural objects underlying human similarity judgements. *Nature*
598 *human behaviour*, **4**, 1173–1185.

599 Ince RA, Paton AT, Kay JW, Schyns PG (2021) Bayesian inference of population prevalence.
600 *Elife*, **10**, e62461.

601 Kawakita G, Zeleznikow-Johnston A, Tsuchiya N, Oizumi M (2023) Comparing color similarity
602 structures between humans and LLMs via unsupervised alignment. *arXiv preprint*
603 *arXiv:2308.04381*.

604 Lakens D (2014) Performing high-powered studies efficiently with sequential analyses.
605 *European Journal of Social Psychology*, **44**, 701–710.

606 Lakens D, Scheel AM, Isager PM (2018) Equivalence testing for psychological research: A
607 tutorial. *Advances in methods and practices in psychological science*, **1**, 259–269.

608 Lau H, Michel M, LeDoux JE, Fleming SM (2022) The mnemonic basis of subjective experience.
609 *Nature Reviews Psychology*, **1**, 479–488.

610 Malach R (2021) Local neuronal relational structures underlying the contents of human
611 conscious experience. *Neuroscience of consciousness*, **2021**, niab028.

612 Marjeh R, Sucholutsky I, van Rijn P, Jacoby N, Griffiths TL (2023) Large language models predict
613 human sensory judgments across six modalities. *arXiv preprint arXiv:2302.01308*.

614 McCurdy LY, Maniscalco B, Metcalfe J, Liu KY, De Lange FP, Lau H (2013) Anatomical coupling
615 between distinct metacognitive systems for memory and visual perception. *Journal of*
616 *Neuroscience*, **33**, 1897–1906.

617 Medin DL, Goldstone RL, Gentner D (1993) Respects for similarity. *Psychological review*, **100**,
618 254.

619 Moharramipour A, Lau H (2024) Open Review of Kawakita et al’s “Is my ‘red’ your ‘red’?” (2023)
620 PsyArXiv.

621 Morales J, Lau H, Fleming SM (2018) Domain-general and domain-specific patterns of activity
622 supporting metacognition in human prefrontal cortex. *Journal of Neuroscience*, **38**,
623 3534–3546.

624 Nosofsky RM (1984) Choice, similarity, and the context theory of classification. *Journal of*
625 *Experimental Psychology: Learning, memory, and cognition*, **10**, 104.

626 Paysan P, Knothe R, Amberg B, Romdhani S, Vetter T (2009) A 3D face model for pose and
627 illumination invariant face recognition. In:., pp. 296–301. IEEE.

628 Rao RP (1999) An optimal estimation approach to visual perception and learning. *Vision*
629 *research*, **39**, 1963–1989.

630 Rosenthal D (2010) How to think about mental qualities. *Philosophical Issues*, **20**, 368–393.

631 Schurgin MW, Wixted JT, Brady TF (2020) Psychophysical scaling reveals a unified theory of
632 visual memory strength. *Nature human behaviour*, **4**, 1156–1172.

633 Shen S, Ma WJ (2016) A detailed comparison of optimality and simplicity in perceptual decision
634 making. *Psychological review*, **123**, 452.

635 Shepard RN (1964) Attention and the metric structure of the stimulus space. *Journal of*
636 *mathematical psychology*, **1**, 54–87.

637 Shepard RN (1987) Toward a universal law of generalization for psychological science. *Science*,
638 **237**, 1317–1323.

639 Sims CR (2018) Efficient coding explains the universal law of generalization in human
640 perception. *Science*, **360**, 652–656.

641 Smith LB (1989) A model of perceptual classification in children and adults. *Psychological*
642 *review*, **96**, 125.

643 Smith LB, Heise D (1992) Perceptual similarity and conceptual structure. In: *Advances in*
644 *psychology*, pp. 233–272. Elsevier.

645 Tallon-Baudry C (2022) The topological space of subjective experience. *Trends in Cognitive*
646 *Sciences*, **26**, 1068–1069.

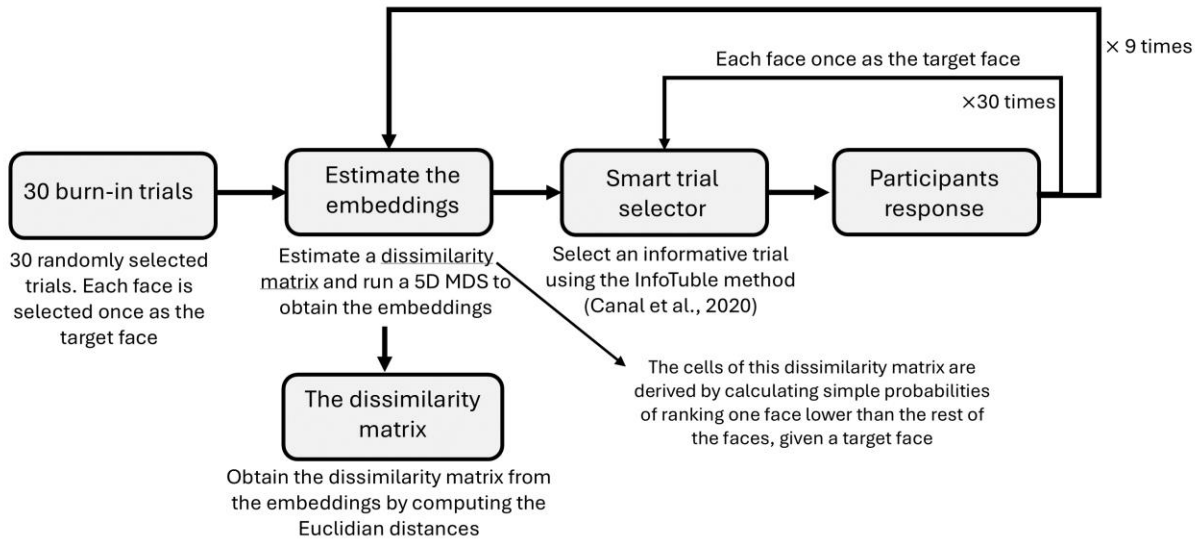
647 Tenenbaum JB, Griffiths TL (2001) Generalization, similarity, and Bayesian inference. *Behavioral*
648 *and brain sciences*, **24**, 629–640.

649 Tversky A (1977) Features of similarity. *Psychological review*, **84**, 327.

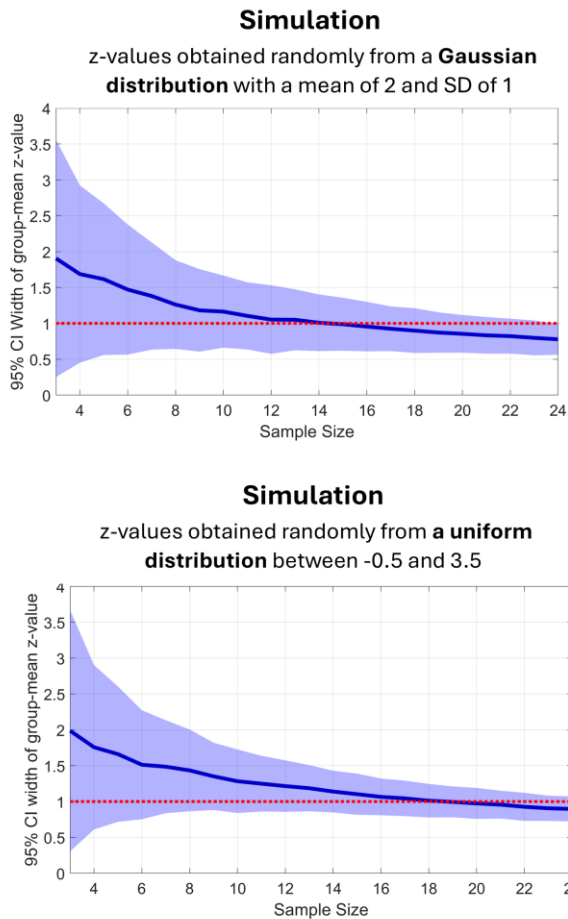
650 Zeleznikow-Johnston A, Aizawa Y, Yamada M, Tsuchiya N (2023) Are color experiences the
651 same across the visual field? *Journal of Cognitive Neuroscience*, **35**, 509–542.

652

Supplementary Figures



Supplementary Figure 1. Schematic of the subjective similarity judgment task design



Supplementary Figure 2. Simulations to determine a feasible sample size for our stopping criterion.

We ran two simulations: one generated z-values from a Gaussian distribution with a mean of 2 and SD of 1 and another from a uniform distribution ranging between -0.5 and 3.5. These example distributions seem reasonable given our expectations based on our pilot data and seem conservative enough. For example, the 95% CI width in our pilot data with four participants was around 1.5, however, in the presented simulations, the 95% CI width is, on average, around 1.7 for the same sample size of four. The shaded area indicates the 2.5th and 97.5th percentile of the 95% CI width obtained over 1000 simulations. Assuming the used distributions are realistic, there is a high likelihood of hitting the stopping criterion by reaching a maximum sample size of 24.