1

2

3

4        Perceived similarity ratings predict generalization success after traditional category learning and

5                                    a new paired-associate learning task

6                          Stefania R. Ashby, Caitlin R. Bowman, Dagmar Zeithamova

7                                              University of Oregon

8

9

10

11

12

13

14                                                  Author Note

15          Stefania R. Ashby, Department of Psychology, University of Oregon; Caitlin R. Bowman,

16   Department of Psychology, University of Oregon; Dagmar Zeithamova, Department of

17   Psychology, University of Oregon.

18          This work was supported by the National Institute on Aging Grant F32-AG-054204

19   awarded to Caitlin R. Bowman.

20          Correspondence concerning this article should be addressed to Dagmar Zeithamova,

21   Department of Psychology, 1227 University of Oregon, Eugene, OR 97403.

22   Email: dasa@uoregon.edu

23   Phone: 541-346-6731

24                                             Abstract

25          The current study investigated category learning across two experiments using face-blend

26   stimuli that formed face families controlled for within- and between-category similarity.

27   Experiment 1 was a traditional feedback-based category learning task, with three family names

28   serving as category labels. In Experiment 2, the shared family name was encountered in the context

29   of a face—full name paired-associate learning task, with a unique first name for each face. A

30   subsequent test that required participants to categorize new faces from each family showed

31   successful generalization in both experiments. Furthermore, perceived similarity ratings for pairs

32   of faces were collected before and after learning, prior to generalization test. In Experiment 1,

33   similarity ratings increased for faces within a family and decreased for faces that were physically

34   similar but belonged to different families. In Experiment 2, overall similarity ratings decreased

35   after learning, driven primarily by decreases for physically similar faces from different families.

36   The post-learning category bias in similarity ratings was predictive of subsequent generalization

37   success in both experiments. The results indicate that individuals formed generalizable category

38   knowledge prior to an explicit demand to generalize, and did so both when attention was directed

39   towards category-relevant similarities (Experiment 1) and when attention was directed towards

40   individuating faces within a family (Experiment 2). The results tie together research on category

41   learning and categorical perception and extend them beyond a traditional category learning task.

42          *Keywords:* category learning, perceived similarity, memory generalization

43      Perceived similarity ratings predict generalization success after traditional category learning and

44                          a new paired-associate learning task

45          Categorization helps us organize information from the world around us into meaningful

46      clusters relevant to behavior. A hallmark of category knowledge is the ability to categorize new

47      instances, allowing us to generalize prior knowledge and guide decisions in novel situations.

48      Category-learning tasks have thus been widely used to study memory generalization (Knowlton &

49      Squire, 1993; Nosofsky & Zaki, 1998; Poldrack et al., 2001; Reber, Stark, & Squire, 1998). Prior

50      work has contended that memory generalization relies on memory representations that form during

51      learning, linking information across related experiences (Knowlton & Squire, 1993; Schapiro,

52      Turk-Browne, Botvinick, & Norman, 2017; Schlichting, Zeithamova, & Preston, 2014; Shohamy

53      & Wagner, 2008; Zeithamova, Schlichting, & Preston, 2012). Other work maintains that specific

54      memory traces are formed during learning and that generalization judgements may be computed

55      from specific memories either on-the-fly at retrieval (Hintzman, 1984; Kruschke, 1992; Nosofsky,

56      1988) or by linking information across experiences in response to explicit generalization demands

57      (Carpenter & Schacter, 2017, 2018; Squire, 1992; Teyler & DiScenna, 1986; Winocur,

58      Moscovitch, & Sekeres, 2007). Finding ways to detect category knowledge in behavior outside of

59      generalization demands will help us to determine whether or not people spontaneously link related

60      experiences as they are encountered.

61          Category knowledge alters perception such that items learned to belong to the same

62      category are perceived as more similar while items learned to belong to different categories are

63      perceived as less similar after learning  (Beale & Keil, 1995; Folstein, Palmeri, & Gauthier, 2013;

64      Goldstone, 1994a; Goldstone, Lippa, & Shiffrin, 2001; Livingston, Andrews, & Harnad, 1998;

65      Rosch & Mervis, 1975). Thus, measures of perceived similarity may be useful for assessing

66  category knowledge without creating an explicit generalization demand. In the current report, we

67  conducted two experiments testing whether measures of perceived similarity can reveal the

68  formation of category knowledge prior to an explicit generalization test. Participants were shown

69  faces that belonged to three categories (families), designated by a family name. Face stimuli were

70  created as blends of never-seen "parent" faces, resulting in increased physical similarity between

71  faces that shared a parent (Figure 1). Some physically similar faces were members of the same

72  family while others were members of different families, allowing us to dissociate the effect of

73  category membership from physical similarity.

74      In Experiment 1, faces were encountered in the context of a traditional feedback-based

75  category learning task, emphasizing similarities among faces belonging to the same family. We

76  tested participants' ability to extract commonalities across faces belonging to the same family and

77  generalize family names to new face-blend stimuli. We also measured category knowledge

78  indirectly, using perceived similarity ratings immediately before and after learning to determine to

79  what degree people link related experiences prior to explicit demands to generalize. The category

80  bias in perceived similarity ratings after learning was related to subsequent generalization success

81  to determine the utility of using perceived similarity ratings as a measure of category learning. The

82  same measures were collected in Experiment 2, where faces were encoded through observational,

83  face—full name paired-associate learning. While family names were identical to Experiment 1,

84  with each family name shared across several faces, first names were unique for each face, requiring

85  the participant to differentiate faces within each family. This allowed us to test to what degree the

86  results from Experiment 1 replicate outside of a traditional category learning task context.

87                                                   **Methods**

88     **Participants**

89          Healthy participants—N = 39 in Experiment 1 and N = 43 in Experiment 2—were recruited

90     from the University of Oregon community via the university SONA research system and received

91     course credit for their participation. Except for the learning phase, all procedures were identical

92     across experiments and will be presented together. All participants provided written informed

93     consent, and experimental procedures were approved by Research Compliance Services at the

94     University of Oregon. From Experiment 1, four participants were excluded due to chance

95     performance (accuracy $\leq$ .33) in categorizing the training faces. From Experiment 2, participants

96     were excluded for failing to make responses on more than 25% of categorization trials (n = 3) and

97     incomplete data (n = 1). After exclusions, analyses were carried out with the remaining 35

98     participants for Experiment 1 ($M_{age}$ = 20.43, $SD_{age}$ = 2.58, 18-32 years, 21 females) and 39

99     participants for Experiment 2 ($M_{age}$ = 19.26, $SD_{age}$ = 1.13, 18-23 years, 21 females). These sample

100    sizes provide 80% power for detecting medium size (d $\geq$ 0.5) effects using planned one-sample

101    and paired t-tests and strong (r $\geq$ .5) correlations, as determined in G-Power (Faul, Erdfelder,

102    Buchner, & Lang, 2009; Faul, Erdfelder, Lang, & Buchner, 2007).

103    **Stimuli**

104         Stimuli were grayscale images of blended faces constructed by morphing two unaltered

105    face images together using FantaMorph Version 5 by Abrosoft. Prior work has shown that category

106    effects differ based on whether morphed faces are constructed from parents within one race versus

107    across two races (Levin & Angelone, 2002). Thus, we restricted all parent faces to be Caucasian

108    to ensure that the resulting face-blend stimuli were comparably similar to all other faces with a

109    shared parent. Additionally, all parent faces were of a single gender (male) to ensure that face-

110    blends maintained a realistic appearance.

111         The stimulus structure is presented in Figure 1. For each participant, three category-

112    relevant parent faces and three category-irrelevant parent faces were randomly selected from a

113    total set of twenty faces. Each of the three category-relevant parent faces were individually

114    morphed with each of the three category-irrelevant parent faces with equal weight given to each

115    parent face (50/50 blend). The resulting nine blended faces were then used as training stimuli.

116    Faces that shared a category-relevant parent shared a family name (belonged to the same category).

117    Faces that shared a category-irrelevant parent belonged to different families. As faces sharing any



*Figure 1*. Example face-blend stimuli. Parent faces on the leftmost side are designated "category relevant parents" as these parents determined family membership—Miller, Wilson, or Davis—during learning and generalization. Parent faces across the top are designated "category-irrelevant parents" as these parents introduced physical similarity across families but did not determine categories. Three category-irrelevant parents were used for learning. The rightmost three category-irrelevant parents are a subset of new faces used for generalization. Parent faces were never viewed by participants, only the resulting blended faces. The face blending procedure produced pairs of faces that shared a category-relevant parent and belonged to the same family (shared parent - same family name; example indicated with dark grey box), pairs of faces that shared a category-irrelevant parent and belonged to different families (shared parent- different family name; example indicated with medium grey box). Non-adjacent pairs did not share a parent and were not related (example indicated with light grey boxes).

118    parent (category-relevant or category-irrelevant) shared physical traits, physical similarity alone

119    was not diagnostic of category membership. Generalization stimuli were new faces created by

120    blending category-relevant parent faces with fourteen remaining parent faces not used for creation

121    of the training faces.

122    **Procedure**

123    Both experiments consisted of the following phases: passive viewing, pre-learning

124    similarity ratings, learning (different in each Experiment), passive viewing, post-learning

125    similarity ratings, and category generalization. Additionally, Experiment 2 included cued-recall of

126    face-name associations before the category generalization phase. Self-paced breaks separated the

127    phases.

128    **Passive viewing.** To familiarize participants with the stimuli and give them an idea of the

129    degree of similarity between all faces before collecting perceived similarity ratings, participants

130    first viewed each of the nine training stimuli individually, once in a random order without any

131    labels and without making any responses. Face-blends were shown for 3s with a 1s inter-stimulus-

132    interval (ISI). Passive viewing of the face-blends immediately before the pre- and post-learning

133    similarity rating phases was also included as a pilot of a future neuroimaging experiment. No

134    responses were collected during viewing.

135    **Pre-learning similarity ratings.** To validate that participants were sensitive to the

136    similarity structure among faces introduced by the blending process and to obtain baseline

137    similarity ratings, participants rated the subjective similarity of pairs of faces to be used during the

138    learning phase. All possible 36 pairwise comparisons of the 9 training faces were presented and

139    participants rated the similarity of the two faces on a scale from one to six (1 = two faces appeared

140    very dissimilar, 6 = two faces appeared very similar). Face pairs and the similarity rating scale

141    were displayed for 5s with a 1s ISI. Face pairs were then binned into three conditions for analyses

142    depending on whether they 1) shared a parent and a family name, 2) shared a parent face but did

143    not share a family name, or 3) did not share a parent face (see example pairs in Figure 1).

144    **Learning phase.**

145    *Experiment 1: Feedback-based category learning.* On each trial, a training face was

146    presented on the screen along with family names (Miller, Wilson, Davis) as response options.

147    Participants were instructed to indicate family membership via a button press and received

148    corrective feedback after each trial. Each face was viewed simultaneously with the family name

149    response options on the screen for 4s, received corrective feedback for 1s, and trials were separated

150    by a 1s ISI. Each face was presented 16 times total, evenly split across 2 blocks.

151    *Experiment 2: Observational learning of face—full name associations.* To test the

152    robustness of category learning outside of a traditional categorization task, Experiment 2 provided

153    an opportunity to link faces from the same families in the context of a face—full name associative

154    learning task. On each trial, participants studied a face-name pair and then made a prospective

155    memory judgement on a scale from one to four (1 = definitely will not remember, 4 = definitely

156    will remember). Prospective memory judgments were included to facilitate participant engagement

157    with the observational learning task and were not considered further. Family names were identical

158    to Experiment 1 and shared across faces whereas first names were unique to each face. While the

159    inclusion of face-specific first names required participants to differentiate individual faces, the

160    inclusion of the shared family names provided an opportunity to form links between related faces.

161    The fact that family names were repeated across faces or that there was a category structure among

162    faces was not explicitly emphasized to participants. Each face-name pair was presented on screen

163    for 2s after which the prospective memory judgment scale appeared beneath the face-name pair

164    for an additional 2s. Trials were separated by a 4s ISI. Participants viewed each face-name pair

165    twelve times, evenly split across 3 blocks.

166            **Post-learning similarity ratings.** Perceived similarity ratings were repeated after the

167    learning phase with the same timing as pre-learning ratings. Of main interest was a potential

168    category bias in perceived similarity, i.e., whether faces that shared a parent would be rated as

169    more similar when they had the same family name than when they had different family names.

170            **Cued recall of face-name associations.** Experiment 2 included a self-paced cued-recall

171    task of face-name associations. Participants viewed each training face individually on a computer

172    screen and handwrote the full name of each face on a sheet of paper. Participants advanced the

173    trials at their own pace but were not able to skip faces or go back and look at faces already named.

174    Participants were encouraged to make their best guess as to the first and family names of each face

175    even if they were not confident in their memory.

176            **Generalization phase.** As the last phase of both Experiments, category knowledge was

177    tested directly using categorization of old and new faces. In addition to the nine training faces,

178    participants categorized 42 never-seen faces, consisting of 14 new blends of each of the three

179    category-relevant parent faces. Participants were asked to select via button press the family name

180    for each face, which were presented individually for 4s, from the three options (Miller, Wilson,

181    Davis) presented on the screen. Trials were separated by an 8s ISI. No feedback was provided, and

182    participants were encouraged to make their best guess when unsure of family membership.

183                                              **Results**

184    **Learning Phase**

185            **Experiment 1: Feedback-based category learning**. Overall percent correct across

186    training was 76% (SD = 14%), which was well above chance (.33 for three categories; one-sample

187    t(34) = 17.66, p < .001, d = 3.01).  Categorization accuracy improved across training, from 66%

188    in the first half to 85% in the second half (t(34) = 9.72, p < .001, d = 1.63), demonstrating learning

189    over time.

190         **Experiment 2: Observational learning of full name—face associations.** Observational

191    learning provided no measure of accuracy from the learning phase. Therefore, in Experiment 2 a

192    cued-recall task was included to assess how well participants learned the face-full name pairs.

193    Participants recalled on average 52% of first names and 65% of family names.

194    **Similarity Ratings**

195         We compared mean face similarity ratings in each pair-type (shared parent-same family

196    name, shared parent-different family name, not related) using repeated-measures ANOVA.

197    Analyses were performed separately in each phase (pre-learning, post-learning). We also assessed

198    learning-related rating changes by comparing ratings across phases. For all ANOVAs, a

199    Greenhouse-Geisser correction for degrees of freedom (denoted as *GG*) was used wherever

200    Mauchly's test indicated a violation of the assumption of sphericity.

201         **Experiment 1.** Pre-learning ratings (Fig. 2A) demonstrated that participants were sensitive

202    to the physical similarity structure introduced with the face-blending procedure. A one-way,

203    repeated measures ANOVA showed a significant effect of pair type (F(2, 68) = 58.74, p < .001,

204    $\eta_p^2$ = .63), driven by lower perceived similarity for faces that did not share a parent compared to

205    those that shared a parent (with or without shared family name, both t > 9.17, p < .001, d > 1.50).

206    Faces that shared a parent were perceived as equally similar to one another irrespective of whether

207    they also shared the same—not yet presented—family name (t(34) = -0.17, p = .87, d = 0.03).

208         Post-learning ratings (Fig. 2B) revealed a category bias on perceived similarity: pairs of

209    faces sharing a parent and family name were perceived as significantly more similar than faces

210    that shared a parent but not a family name ($M_{diff}$ = 0.72, $SD_{diff}$ = 1.41, t(34) = 3.02, p = .005, d =

211    0.51). Faces that shared a parent remained rated as more similar than unrelated faces (both t > 6.85,

212    p < .001, d > 1.15).

213          To further test the effect of learning, we conducted a 2 x 3 (timepoint [pre-learning, post-

214    learning] x pair-type [shared parent-same family name, shared parent-different family name, not

215    related]) repeated-measures ANOVA. There was no main effect of timepoint (F(1, 34) = 0.04, p =

216    .85, $\eta_p^2$ = .001). There was a significant main effect of pair-type (F(1.63, 55.38) = 61.21, p < .001,

217    $\eta_p^2$ = .64, *GG*), and a significant interaction between timepoint and pair-type (F(1.64, 55.88) =

218    11.85, p < .001, $\eta_p^2$ = .25, *GG*). Follow-up pre-post comparisons within each pair-type (Fig. 2C)

219    revealed that this interaction was driven by both a significant *increase* in similarity ratings for

220    faces sharing a parent and a family name (t(34) = 3.02, p = .005, d = 0.51) and a significant

221    *decrease* in similarity ratings for faces only sharing a parent but not a family name (t(34) = -2.33,

222    p = .026, d = -0.39). There was no significant change in similarity ratings for faces that did not

223    share a parent (t(34) = -0.18, p = .86, d = -0.03).

224          **Experiment 2.** As in Experiment 1, participants were sensitive to the face similarity

225    structure. Pre-learning similarity ratings (Fig. 2E) differed significantly among pair types (F(1.46,

226    55.47) = 72.22, p < .001, $\eta_p^2$ = .655, *GG*), driven by lower perceived similarity of faces that did not

227    share a parent compared to faces that shared a parent (with and without shared family names, both

228    t > 10.65, p < .001, d > 1.70). For faces that shared a parent, ratings did not significantly differ

229    when face pairs had the same or different—not yet presented—family names (t(38) = 1.82, p =

230    .077, d = 0.29). A category bias was found in post-learning ratings (Fig. 2F) with pairs of faces

231    sharing a parent and family name perceived as significantly more similar than faces that shared a

232    parent but not a family name ($M_{diff}$ = 0.58, $SD_{diff}$ = 1.52; t(38) = 2.39, p = .022, d = 0.38).

233        Testing the effect of learning, the 2 x 3 (timepoint x pair-type) repeated-measures ANOVA

234    revealed a significant main effect of timepoint (F(1, 38) = 5.20, p = .028, $\eta_p^2$ = .120), with overall

235    similarity ratings being lower post-learning than pre-learning ($M_{pre}$ = 3.49, $SD_{pre}$ = 0.51; $M_{post}$ =

236    3.33, $SD_{post}$ = 0.59; t(38) = -2.28, p = .028, d = 0.37). There was also a significant main effect of

237    pair-type (F(1.28, 48.60) = 60.42, p < .001, $\eta_p^2$ = .614, *GG*), and a significant interaction between

238    timepoint and pair-type (F(1.67, 63.37) = 4.21, p = .03, $\eta_p^2$ = .10, *GG*). Follow-up pre-post

239    comparisons within each pair-type (Fig. 2G) revealed that the interaction was driven by a

240    significant *decrease* in similarity ratings for faces sharing a parent but not a family name (t(38) =

241    -3.71, p = .001, d = -0.59), but there were no significant changes in similarity ratings for other pair-



*Figure 2.* Top panel are results from the traditional category learning experiment. Bottom panel (shaded grey) are results from the face-name paired associate learning experiment. *A & E.* Average similarity ratings for faces that share a parent and family name, faces that only share a parent, and faces that don't share any parents before learning. *B & F.* Average similarity ratings for the same pairwise comparisons after learning. Asterisk represents a significant (p < .05) difference in post-learning similarity ratings for faces that belong to the same family vs. faces that share physical similarity but belong to different families (i.e. a category bias in perception). *C & G.* Changes in similarity ratings from pre- to post-learning. Asterisk denotes significant (p < .05) increases and decreases in perceived similarity for faces. *D & H.* Positive relationship between indirect (category bias in perception) and direct (categorization accuracy for new faces) measures of memory generalization.

242    types (both t < -1.04, p > .30, d < -0.18). Thus, changes in perceived similarity were affected by

243    category membership in both experiments.

244          Although not significant (p = .077), we noted a numerical tendency towards a category bias

245    in pre-learning similarity ratings. Parent faces were randomly selected for each participant to serve

246    as category-relevant or category-irrelevant parents, but some of the category-relevant parent faces

247    may have been more salient, leading to a numerically greater pre-learning similarity rating. Thus,

248    we tested whether the post-learning category bias on perceived similarity was reliably greater than

249    pre-learning bias. A 2 x 2 (timepoint [pre-learning, post-learning] x pair-type [shared parent-same

250    family name, shared parent-different family name]) repeated-measures ANOVA showed only a

251    marginal interaction between timepoint and condition (F(1, 38) = 2.87, p = .098, $\eta_p^2$ = .07). We

252    thus controlled for pre-learning similarity rating differences in subsequent analyses that assessed

253    the relationship of post-learning ratings and generalization performance.

**Category Generalization**

255          **Experiment 1.** Participants correctly categorized 85% of training faces (SD = 17%) and

256    74% of new faces (SD = 13%), which was well above chance (.33 for three categories; both one-

257    sample t(34) > 18.12, p < .001, d > 3.06). A paired-samples t-test showed higher categorization

258    accuracy for the training faces than for the new faces (t(34) = 5.48, p < .001 , d = 0.93). We next

259    tested whether the category bias on perceived similarity ratings (an indirect measure of category

260    knowledge) was related to subsequent generalization success. A Pearson's correlation showed a

261    significant positive relationship between the category bias on perceived similarity ratings and

262    generalization accuracy (r(33) = .64, p < .001; Fig. 2D). The category bias on perceived similarity

263    in the post-learning phase was a significant predictor of subsequent generalization performance

264    even when pre-learning similarity ratings were considered (multiple regression: pre-learning

265 differences in perceived similarity $\beta = .30$, $t(34) = 1.80$, $p = .08$; post-learning category bias $\beta =$

266 $.46$, $t(34) = 2.75$, $p = .01$).

267 **Experiment 2.** Participants correctly categorized 70% of training faces (SD = 23%) and

268 64% of new faces (SD = 22%), which was well above chance (.33 for three categories; both one-

269 sample $t(38) > 8.65$, $p < .001$, $d > 1.38$). A paired-samples t-test showed higher categorization

270 accuracy for the training faces than for new faces ($t(38) = 2.12$, $p = .04$, $d = 0.34$). The post-

271 learning category bias on perceived similarity ratings was significantly correlated with

272 generalization accuracy (Pearson's $r(37) = .48$, $p = .002$; Fig. 2H). Further, the category bias was

273 a significant predictor of subsequent generalization performance even when pre-learning similarity

274 ratings were controlled for (multiple regression: pre-learning category bias $\beta = -.22$, $t(38) = -0.86$,

275 $p = .40$; post-learning category bias $\beta = .66$, $t(38) = 2.57$, $p = .01$).

276 **Discussion**

277 The current study investigated category learning using measures of perceived similarity

278 and category generalization across two experiments. Face-blend stimuli were used to control

279 physical similarity within and across categories (families). Experiment 1 was a traditional

280 feedback-based category learning task, with three family names serving as category labels. In

281 Experiment 2, the shared family name category label was encountered in the context of a face-full

282 name paired-associate learning task, where first names were unique for each face. We were

283 interested in how well people generalize family names to new faces in the two tasks and to what

284 degree category bias in perceived similarity ratings indicates the formation of category knowledge

285 prior to explicit generalization demands.

286 Participants were able to successfully apply category labels to new faces in both

287 experiments, demonstrating that category information can be extracted in support of generalization

288    even when task goals do not emphasize learning categories at encoding. Past work has shown that

289    individuals can extract category structures when not instructed using patterns of physical similarity

290    as category cues (Aizenstein et al., 2000; Love, 2002; Reber, Gitelman, Parrish, & Mesulam,

291    2003). We extend these prior findings by showing that category structure can also be extracted

292    when category membership is dissociable from physical similarity and further when individuals

293    are actively learning information that differentiates individual items even within the same

294    category.

295         Learning-related changes in perceived similarity ratings were observed in both

296    experiments. In Experiment 1, consistent with prior studies (Beale & Keil, 1995; Goldstone,

297    1994a, 1994b; Goldstone et al., 2001; Rosch & Mervis, 1975), similarity ratings for faces within

298    a family increased while similarity ratings for faces that were physically similar but belonged to

299    different families decreased. These shifts in perceived similarity may reflect allocation of selective

300    attention to features that are category-relevant while diverting attention away from category-

301    irrelevant features (Goldstone & Steyvers, 2001; Kruschke, 1996; Nosofsky, 1991). In contrast, in

302    Experiment 2 the face-name paired-associate learning was associated with an overall decrease in

303    similarity ratings from pre- to post-encoding, driven primarily by decreased similarity for faces

304    that were physically similar but belonged to different families. This decrease in similarity ratings

305    could reflect learning-related differentiation of representations to minimize confusability and

306    interference (Chanales, Oza, Favila, & Kuhl, 2017; Favila, Chanales, & Kuhl, 2016; Hulbert &

307    Norman, 2015; Kim, Norman, & Turk-Browne, 2017; Lohnas et al., 2018).  Changes in perceived

308    similarity ratings were modulated by category membership of the faces in both experiments,

309    indicating that people tended to link faces with a shared last name even outside the context of a

310    traditional category learning task.

311          The inclusion of similarity ratings also allowed us to address the question of whether or

312    not people spontaneously link related information in service of generalization prior to explicit

313    generalization demands. The category bias in similarity ratings observed after learning predicted

314    subsequent generalization of category information to new examples in both experiments,

315    indicating that both measures index the same category knowledge formation. Critically, the

316    category bias was measured *after* learning but *before* the explicit generalization test, indicating

317    that people likely linked related faces at encoding (see also Shohamy & Wagner, 2008;

318    Zeithamova, Dominick, & Preston, 2012) rather than in response to generalization demands. Our

319    results also extend prior studies on changes in perceived similarity as a result of explicit instruction

320    where attention is directed towards category-relevant similarities (Goldstone, 1994b, 1994a;

321    Livingston et al., 1998) to a novel task where attention was directed towards individuating

322    differences. Observation of the category bias after the face-name paired-associate learning

323    indicates that the mere presence of a shared piece of information biased perceived similarity in

324    many participants.

325          In summary, our findings indicate that generalizable category representations form at

326    encoding, prior to explicit generalization demands. Individuals spontaneously linked related

327    experiences to form conceptual knowledge even when learning goals required participants to learn

328    individuating differences between stimuli. The relationship between category bias in similarity

329    ratings and subsequent generalization further indicates that measures of perceived similarity are

330    useful for measuring category learning without explicit demands to generalize. Building upon long

331    lines of research on category learning (*for reviews see* Ashby & Maddox, 2011; Seger, 2008) and

332    categorical perception (Etcoff & Magee, 1992; Liberman, Harris, Hoffman, & Griffith, 1957;

333    Livingston et al., 1998; *for reviews see* Goldstone & Hendrickson, 2010; Harnad, 2006), the

334    current work links generalization and perception together and extends prior findings beyond

335    traditional category learning paradigms.

**Open Practices**

337        None of the experiments discussed in the current report were preregistered. Data and

338    materials for all experiments are freely available in the *Blended-Face Similarity Ratings and*

339    *Categorization Tasks* repository on the Open Science Framework

340    (https://osf.io/e8htb/?view_only=ca5a189813b14dfebd9804151bc1a1ed).

341                                    References

342    Aizenstein, H., MacDonald, A., Stenger, V., Nebes, R., Larson, J., Ursu, S., & Carter, C. (2000).

343        Complementary category learning systems identified using fMRI. *Journal of Cognitive*

344        *Neuroscience*, *12*(6), 977–987.

345    Ashby, F. G., & Maddox, W. T. (2011). Human category learning 2.0. *Annals of the New York*

346        *Academy of Sciences*, *1224*(1), 147–161. https://doi.org/10.1111/j.1749-6632.2010.05874.x

347    Beale, J. M., & Keil, F. C. (1995). Categorical effects in the perception of faces. *Cognition*, *57*,

348        217–239.

349    Carpenter, A. C., & Schacter, D. L. (2017). Flexible retrieval: When true inferences produce

350        false memories. *Journal of Experimental Psychology: Learning, Memory, and Cognition*,

351        *43*(3), 335–349.

352    Carpenter, A. C., & Schacter, D. L. (2018). False memories, false preferences: Flexible retrieval

353        mechanisms supporting successful inference bias novel decisions. *Journal of Experimental*

354        *Psychology: General*, *147*(7), 988–1004. https://doi.org/10.1037/xge0000391

355    Chanales, A. J. H., Oza, A., Favila, S. E., & Kuhl, B. A. (2017). Overlap among spatial

356        memories triggers repulsion of hippocampal representations. *Current Biology*, *27*, 1–11.

357        https://doi.org/10.1016/j.cub.2017.06.057

358    Etcoff, N. L., & Magee, J. J. (1992). Categorical perception of facial expressions. *Cognition*,

359        *44*(3), 227–240. https://doi.org/10.1016/0010-0277(92)90002-Y

360    Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using

361        G*Power 3.1: Tests for correlation and regression analyses. *Behavior Research Methods*,

362        *41*(4), 1149–1160. https://doi.org/10.3758/BRM.41.4.1149

363    Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G*Power: A flexible statistical power

analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, *39*(2), 175–191. https://doi.org/10.3758/BF03193146

Favila, S. E., Chanales, A. J. H., & Kuhl, B. A. (2016). Experience-dependent hippocampal pattern differentiation prevents interference during subsequent learning. *Nature Communications*, *7*, 11066. https://doi.org/10.1038/ncomms11066

Folstein, J. R., Palmeri, T. J., & Gauthier, I. (2013). Category learning increases discriminability of relevant object dimensions in visual cortex. *Cerebral Cortex*, *23*(4), 814–823. https://doi.org/10.1093/cercor/bhs067

Goldstone, R. L. (1994a). Influences of categorization on perceptual discrimination. *Journal of Experimental Psychology: General*, *123*(2), 178–200. https://doi.org/10.1037/0096-3445.123.2.178

Goldstone, R. L. (1994b). The role of similarity in categorization: Providing a groundwork. *Cognition*, *52*, 125–157.

Goldstone, R. L., & Hendrickson, A. T. (2010). Categorical perception. *Wiley Interdisciplinary Reviews: Cognitive Science*, *1*(1), 69–78. https://doi.org/10.1002/wcs.26

Goldstone, R. L., Lippa, Y., & Shiffrin, R. M. (2001). Altering object representations through category learning. *Cognition*, *78*(1), 27–43. https://doi.org/10.1016/S0010-0277(00)00099-8

Goldstone, R. L., & Steyvers, M. (2001). The sensitization and differentiation of dimensions during category learning. *Journal of Experimental Psychology: General*, *130*(1), 116–139.

Harnad, S. (2006). Categorical Perception. In L. Nadel (Ed.), *Encyclopedia of Cognitive Science* (pp. 1–5). https://doi.org/10.1002/0470018860.s00490

Hintzman, D. L. (1984). MINERVA 2: A simulation model of human memory. *Behavior Research Methods, Instruments, & Computers*, *16*(2), 96–101.

387        https://doi.org/10.3758/BF03202365

388    Hulbert, J. C., & Norman, K. A. (2015). Neural differentiation tracks improved recall of

389        competing memories following interleaved study and retrieval practice. *Cerebral Cortex*,

390        *25*(10), 3994–4008. https://doi.org/10.1093/cercor/bhu284

391    Kim, G., Norman, K. A., & Turk-Browne, N. B. (2017). Neural differentiation of incorrectly

392        predicted memories. *The Journal of Neuroscience*, *37*(8), 2022–2031.

393        https://doi.org/10.1523/JNEUROSCI.3272-16.2017

394    Knowlton, B. J., & Squire, L. R. (1993). The learning of categories: Parallel brain systems for

395        item memory and category knowledge. *Science*, *262*, 1747–1749.

396        https://doi.org/10.1126/science.8259522

397    Kruschke, J.K. (1992). ALCOVE: An exemplar-based connectionist model of category learning.

398        *Psychological Review*, *99*(1), 22–44.

399    Kruschke, John K. (1996). Dimensional relevance shifts in category learning. *Connection*

400        *Science*, *8*(2), 225–247. https://doi.org/10.1080/095400996116893

401    Levin, D. T., & Angelone, B. L. (2002). Categorical perception of race. *Perception*, *31*(5), 567–

402        578. https://doi.org/10.1068/p3315

403    Liberman, A. M., Harris, K. S., Hoffman, H. S., & Griffith, B. C. (1957). The discrimination of

404        speech sounds within and across phoneme boundaries. *Journal of Experimental Psychology*,

405        *54*(5), 358–368. https://doi.org/10.1037/h0044417

406    Livingston, K. R., Andrews, J. K., & Harnad, S. (1998). Categorical perception effects induced

407        by category learning. *Journal of Experimental Psychology: Learning Memory and*

408        *Cognition*, *24*(3), 732–753. https://doi.org/10.1037/0278-7393.24.3.732

409    Lohnas, L. J., Thesen, T., Doyle, W. K., Devinsky, O., Duncan, K., & Davachi, L. (2018). Time-

410    resolved neural reinstatement and pattern separation during memory decisions in human

411    hippocampus. *Proceedings of the National Academy of Sciences*, *115*(31), E7418–E7427.

412    https://doi.org/10.1073/pnas.1717088115

413  Love, B. C. (2002). Comparing supervised and unsupervised category learning. *Psychonomic*

414    *Bulletin & Review*, *9*(4), 829–835.

415  Nosofsky, R. M. (1988). Exemplar-Based Accounts of Relations Between Classification,

416    Recognition, and Typicality. *Journal of Experimental Psychology: Learning, Memory, and*

417    *Cognition*, *14*(4), 700–708. https://doi.org/10.1037/0278-7393.14.4.700

418  Nosofsky, R. M. (1991). Tests of an exemplar model for relating perceptual classification and

419    recognition memory. *Journal of Experimental Psychology: Human Perception and*

420    *Performance*, *17*(1), 3–27. https://doi.org/10.1037/0096-1523.17.1.3

421  Nosofsky, R. M., & Zaki, S. R. (1998). Dissociations between categorization and recognition in

422    amnesic and normal individuals. *Psychological Science*, *9*(4), 247–255.

423  Poldrack, R., Clark, J., Paré-Blagoev, E. J., Shohamy, D., Creso Moyano, J., Myers, C., &

424    Gluck, M. A. (2001). Interactive memory systems in the human brain. *Nature*, *414*, 546–

425    550. https://doi.org/10.1038/35107080

426  Reber, P. J., Gitelman, D. R., Parrish, T. B., & Mesulam, M. M. (2003). Dissociating explicit and

427    implicit category knowledge with fMRI. *Journal of Cognitive Neuroscience*, *15*(4), 574–

428    583. https://doi.org/10.1162/089892903321662958

429  Reber, P. J., Stark, C. E. L., & Squire, L. R. (1998). Contrasting cortical activity associated with

430    category memory and recognition memory. *Learning & Memory*, *5*, 420–428.

431    https://doi.org/10.1101/lm.5.6.420

432  Rosch, E., & Mervis, C. B. (1975). Family resemblances. *Cognitive Psychology*, *7*, 573–605.

433          https://doi.org/10.1186/gb-2002-3-12-reports0063

434    Schapiro, A. C., Turk-Browne, N. B., Botvinick, M. M., & Norman, K. A. (2017).

435          Complementary learning systems within the hippocampus: A neural network modelling

436          approach to reconciling episodic memory with statistical learning. *Philosophical*

437          *Transactions of the Royal Society B*, *372*(1711), 20160049.

438          https://doi.org/10.1098/rstb.2016.0049

439    Schlichting, M. L., Mumford, J. A., & Preston, A. R. (2015). Learning-related representational

440          changes reveal dissociable integration and separation signatures in the hippocampus and

441          prefrontal cortex. *Nature Communications*, *6*, 1–10. https://doi.org/10.1038/ncomms9151

442    Schlichting, M. L., Zeithamova, D., & Preston, A. R. (2014). CA1 subfield contributions to

443          memory integration and inference. *Hippocampus*, *24*(10), 1248–1260.

444          https://doi.org/10.1002/hipo.22310

445    Seger, C. A. (2008). How do the basal ganglia contribute to categorization? Their roles in

446          generalization, response selection, and learning via feedback. *Neuroscience and*

447          *Biobehavioral Reviews*, *32*(2), 265–278. https://doi.org/10.1016/j.neubiorev.2007.07.010

448    Shohamy, D., & Wagner, A. D. (2008). Integrating memories in the human brain: Hippocampal-

449          midbrain encoding of overlapping events. *Neuron*, *60*, 378–389.

450          https://doi.org/10.1016/j.neuron.2008.09.023

451    Squire, L. R. (1992). Memory and the hippocampus: A synthesis from findings with rats,

452          monkeys, and humans. *Psychological Review*, *99*(2), 195–231.

453          https://doi.org/10.1037/0033-295X.99.3.582

454    Teyler, T. J., & DiScenna, P. (1986). The hippocampal memory indexing theory. *Behavioral*

455          *Neuroscience*, *100*(2), 147–154. https://doi.org/10.1037/0735-7044.100.2.147

456    Winocur, G., Moscovitch, M., & Sekeres, M. (2007). Memory consolidation or transformation:

457        Context manipulation and hippocampal representations of memory. *Nature Neuroscience*,

458        *10*(5), 555–557. https://doi.org/10.1038/nn1880

459    Zeithamova, D., Dominick, A. L., & Preston, A. R. (2012). Hippocampal and ventral medial

460        prefrontal activation during retrieval-mediated learning supports novel inference. *Neuron*,

461        *75*(1), 168–179. https://doi.org/10.1016/j.neuron.2012.05.010

462    Zeithamova, D., Schlichting, M. L., & Preston, A. R. (2012). The hippocampus and inferential

463        reasoning: Building memories to navigate future decisions. *Frontiers in Human*

464        *Neuroscience*, *6*, 1–14. https://doi.org/10.3389/fnhum.2012.00070