

1 The emergence of words from vocal imitations

2 Pierce Edmiston¹, Marcus Perlman², & Gary Lupyan¹

3 ¹ University of Wisconsin-Madison

4 ² Max Planck Institute for Psycholinguistics

5 Author Note

6 Pierce Edmiston and Gary Lupyan, Department of Psychology, University of
7 Wisconsin-Madison, Madison, Wisconsin. Marcus Perlman, Max Planck Institute for
8 Psycholinguistics, Nijmegen, Netherlands. The work was supported in part by the National
9 Science Foundation (INSPIRE Grant 1344279 to G.L.).

10 Correspondence concerning this article should be addressed to Pierce Edmiston, 1202
11 W. Johnson St., Madison, WI, 53703. E-mail: pedmiston@wisc.edu

12

Abstract

13 People have long pondered the origins of language, especially the words that compose them.
14 Here, we report a series of experiments investigating how conventional spoken words might
15 emerge from imitations of environmental sounds. Does the repeated imitation of an
16 environmental sound gradually give rise to novel word forms? In what ways do these words
17 resemble the original sounds that motivated them? Participants played a version of the
18 children's game "Telephone". The first generation of participants imitated recognizable
19 environmental sounds (e.g., glass breaking, water splashing). Subsequent generations
20 imitated the imitations for a maximum of 8 generations. The results showed that the
21 imitations became more stable and word-like, and later imitations were easier to learn as
22 category labels. At the same time, even after 8 generations, both spoken imitations and their
23 written transcriptions could be matched above chance to the category of environmental
24 sound that motivated them. These results show how repeated imitation can create
25 progressively more word-like forms while continuing to retain a resemblance to the original
26 sound that motivated them, and speak to the possible role of human vocal imitation in
27 explaining the origins of at least some spoken words.

28 *Keywords:* language evolution, iconicity, vocal imitation, transmission chain

29 Word count: 5552

30 The emergence of words from vocal imitations

31 The importance of imitation and depiction in the origin of signs is clearly observable in
32 signed languages (Goldin-Meadow, 2016; Kendon, 2014; Klima & Bellugi, 1980), but in
33 considering the idea that imitation in the vocal modality may be key to understanding the
34 origin of spoken words, many have argued that the human capacity for vocal imitation is far
35 too limited to play a significant role (Arbib, 2012; Armstrong & Wilcox, 2007; Corballis,
36 2003; Hewes, 1973; Hockett, 1978; Tomasello, 2010). For example, Pinker and Jackendoff
37 (2005) argued that, “most humans lack the ability... to convincingly reproduce
38 environmental sounds... Thus ‘capacity for vocal imitation’ in humans might be better
39 described as a capacity to learn to produce speech” (p. 209). Consequently, it is still widely
40 assumed that vocal imitation — or more broadly, the use of any sort of resemblance between
41 form and meaning — cannot be important to understanding the origin of spoken words. We
42 challenge this view by demonstrating that spoken words can emerge from vocal imitations
43 even without the intention to communicate. We find that repeating vocal imitations of
44 environmental sounds over generations of unique speakers is sufficient to create more
45 word-like vocalizations both in form and function.

46 Although most words of contemporary spoken languages are not clearly imitative in
47 origin, there has been a growing recognition of the importance of imitative words in spoken
48 languages (Dingemanse, Blasi, Lupyan, Christiansen, & Monaghan, 2015; Perniss,
49 Thompson, & Vigliocco, 2010) and the frequent use of vocal imitation and depiction in
50 spoken discourse (Clark & Gerrig, 1990; Lewis, 2009). This has led some to argue for the
51 importance of imitation for understanding the origin of spoken words (e.g., Brown, Black, &
52 Horowitz, 1955; Dingemanse, 2014; Donald, 2016; Imai & Kita, 2014; Perlman, Dale, &
53 Lupyan, 2015). In addition, counter to previous assumptions, people are highly effective at
54 using vocal imitations to refer to environmental sounds such as coins dropping in a jar or
55 mechanical events such as scraping — in some cases, even more effective than when using
56 conventional words (Lemaitre & Rocchesso, 2014). Recent work has also shown that people

57 are able to create novel imitative vocalizations for more abstract meanings (e.g. “slow”,
58 “rough”, “good”, “many”) that are understandable to naïve listeners (Perlman et al., 2015).
59 These imitations are effective not because people can mimic environmental sounds with high
60 fidelity, but because people are able to produce imitations that capture the salient features of
61 sounds in ways that are understandable to listeners (Lemaitre, Houix, Voisin, Misdariis, &
62 Susini, 2016). Similarly, the features of onomatopoeic words might highlight distinctive
63 aspects of the sounds they represent. For example, the initial voiced, plosive /b/ in “boom”
64 represents an abrupt, loud onset, the back vowel /u/ a low pitch, and the nasalized /m/ a
65 slow, muffled decay (Rhodes, 1994).

66 Thus, converging evidence suggests that people can use vocal imitation as an effective
67 means of communication. But can vocal imitations ever give rise to words that can be
68 integrated into the vocabulary of a language? And if so, by what means might this happen?
69 To answer these questions, we recruited participants to play an online version of the
70 children’s game of “Telephone”. In the children’s game, a spoken message is whispered from
71 one person to the next. In our version, the original message or “seed sound” was a recording
72 of an environmental sound. The initial group of participants (first generation) imitated these
73 seed sounds, the next generation imitated the previous imitators, and so on for up to 8
74 generations.

75 We then conducted a series of analyses and additional experiments to systematically
76 answer the following questions: First, do imitations stabilize in form and become more
77 word-like as they are repeated? Second, do the imitations retain a resemblance to the original
78 environmental sound that inspired them? If so, it should be possible for naïve participants to
79 match the emergent words back to the original seed sounds. Third, do the imitations become
80 more suitable as labels for the category of sounds that motivated them? For example, does
81 the imitation of a particular water-splashing sound become, over generations of repeated
82 imitation, a better label for the more general category of water-splashing sounds?

83

Experiment 1: Stabilization of imitations through repetition

84

85

86

87

88

89

90

91

92

93

94

In the first experiment, we collected the vocal imitations, and assessed the extent to which repeating imitations of environmental sounds over generations of unique speakers results in progressive stabilization toward more word-like forms. After collecting the imitations, we measured changes in the stability of the imitations in three ways. First, we measured changes in the perception of acoustic similarity between subsequent generations of imitations along contiguous transmission chains. Second, we used algorithmic measures of acoustic similarity to assess the similarity of imitations sampled within and between transmission chains. Third, we obtained transcriptions of imitations, and measured the extent to which later generation imitations were transcribed with greater consistency and agreement. The results show that repeated imitation results in vocalizations that are easier to repeat with high fidelity and easier to transcribe into English orthography.

95

Methods

96

97

98

99

100

101

102

103

104

Selecting seed sounds. To avoid sounds having lexicalized or conventionalized onomatopoeic forms in English, we used inanimate categories of environmental sounds. Using an odd-one-out norming procedure ($N=105$ participants), an initial set of 36 sounds in 6 categories was reduced to a final set of 16 “seed” sounds: 4 sounds in each of 4 categories. The purpose of this norming procedure was to reach a set of approximately equally distinguishable sounds within each category by systematically removing the sounds that stood out in each category. The results of the norming procedure are shown in Fig. S1. The four final categories were: water, glass, tear, zipper. The final 16 seed sounds can be downloaded from here: osf.io/n6g7d/download.

105

106

107

108

Collecting vocal imitations. Participants ($N=94$) recruited from Amazon Mechanical Turk were paid to participate in an online version of the children’s game of “Telephone”. Participants were instructed that they would hear some sound and their task is to reproduce it as accurately as possible using their computer microphone. Full instructions

109 are provided in the Supplemental Materials.

110 Each participant listened to and imitated four sounds: one from each of the four
111 categories of environmental sounds. Sounds were assigned at random such that participants
112 were unlikely to imitate the same person more than once. Participants were allowed to listen
113 to each target sound multiple times, but were only allowed a single recording in response.
114 Recordings that were too quiet (less than -30 dBFS) were not accepted.

115 Imitations were monitored by an experimenter to catch any gross errors in recording
116 before they were heard by the next generation of imitators. For example, recordings with
117 loud sounds in the background were removed, and recordings were trimmed to the length of
118 the imitation prior to the next generation. The experimenter also removed sounds that
119 violated the rules of the experiment, e.g., by saying something in English. A total of 115
120 (24%) imitations were removed prior to subsequent analysis. The final sample contained 365
121 imitations along 105 contiguous transmission chains (Fig. 1).

122 **Measuring acoustic similarity.**

123 ***Acoustic similarity judgments.*** Acoustic similarity judgments were gathered
124 from five research assistants who listened to pairs of sounds (approx. 300) and rated their
125 subjective similarity. On each trial, raters heard two sounds from subsequent generations
126 played in random order. They then indicated the similarity between the sounds on a 7-point
127 Likert scale from *Entirely different and would never be confused* to *Nearly identical*. Raters
128 were encouraged to use as much of the scale as they could while maximizing the likelihood
129 that, if they did this procedure again, they would reach the same judgments. Full
130 instructions are provided in the Supplemental Materials. Inter-rater reliability was calculated
131 as the intra-class coefficient treating the group as the unit of analysis (Gamer, Lemon,
132 Fellows, & Singh, 2012; Shrout & Fleiss, 1979): $ICC = 0.76$, 95% CI [0.70, 0.81], $F(170, 680)$
133 $= 4.18$, $p < 0.001$. Ratings were normalized for each rater (z-scored) prior to analysis.

134 ***Algorithmic acoustic similarity.*** To obtain algorithmic measures of acoustic
135 similarity, we used the acoustic distance functions included in Phonological Corpus Tools

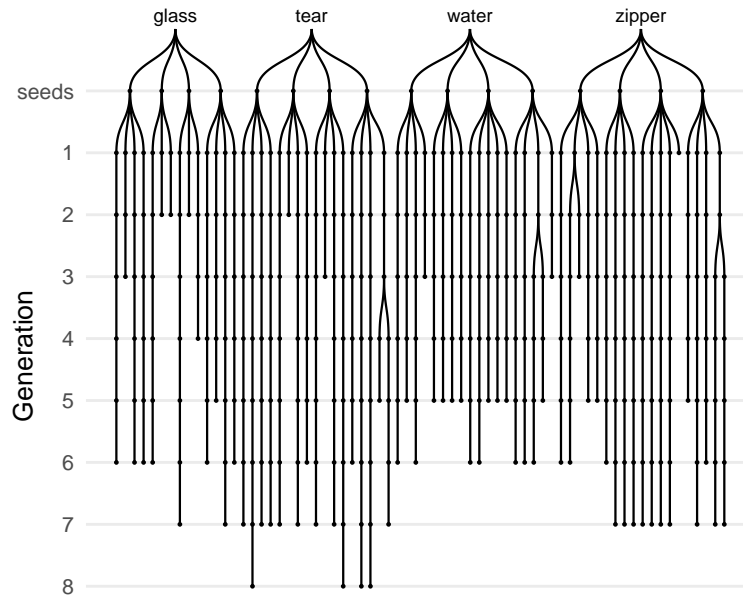


Figure 1. Vocal imitations collected in the transmission chain experiment. Seed sounds (16) were sampled from four categories of environmental sounds: glass, tear, water, zipper. Participants imitated each seed sound, and then the next generation of participants imitated the imitations, and so on, for up to 8 generations. Chains are unbalanced due to random assignment and the exclusion of some low quality recordings.

136 (Hall, Allen, Fry, Mackie, & McAuliffe, 2016). We computed Mel-frequency cepstral
137 coefficients (MFCCs) between pairs of imitations using 12 coefficients in order to obtain
138 speaker-independent estimates.

139 **Collecting transcriptions of imitations.** Participants ($N=216$) recruited from
140 Amazon Mechanical Turk were paid to transcribe vocalizations using English orthography,
141 being instructed to write down what they heard as a single “word” so that the written word
142 would sound as much like the sound as possible. Participants were instructed that this was a
143 word creation task and so to avoid transcribing the vocalizations into existing English words.
144 Each participant completed 10 transcriptions. Transcriptions were gathered for the first and
145 the last three generations of imitations collected in the transmission chain experiment.
146 Participants also provided “transcriptions” of the original environmental seed sounds.
147 Analyses of these transcriptions are reported in the Supplementary Materials (Fig. S5).

148 To measure similarity among transcriptions of the same imitation, we used the
149 `SequenceMatcher` functions in the `difflib` package of the Python standard library, which
150 implements a version of Ratcliff and Obershelp’s “gestalt pattern matching” algorithm.
151 Alternative measures of transcription agreement including exact string matching and the
152 length of the longest substring match were also collected.

153 **Analyses.** Statistical analyses were conducted in R using linear mixed-effects models
154 provided by the `lme4` package (Bates, Mächler, Bolker, & Walker, 2015). Degrees of freedom
155 and corresponding significance tests for linear mixed-effects models were estimated using the
156 Satterthwaite approximation via the `lmerTest` package (Kuznetsova, Bruun Brockhoff, &
157 Haubo Bojesen Christensen, 2016). Random effects (intercepts and slopes) for subjects and
158 for items were included wherever appropriate, and are described below.

159 **Data availability.** Our data along with all methods, materials, and analysis scripts,
160 are available in public repositories described on the Open Science Framework page for this
161 research here: osf.io/3navm.

162 Results

163 **Acoustic similarity increased through iteration.** Imitations of environmental
164 sounds became more stable over the course of being repeated as revealed by increasing
165 acoustic similarity judgments along individual transmission chains. Acoustic similarity
166 ratings were fit with a linear mixed-effects model predicting perceived acoustic similarity
167 from generation with random effects (intercepts and slopes) for raters. To test whether the
168 hypothesized increase in acoustic similarity was true across all seed sounds and categories, we
169 added random effects (intercepts and slopes) for seed sounds nested within categories. The
170 results showed that, across raters and seeds, imitations from later generations were rated as
171 sounding more similar to one another than imitations from earlier generations, $b = 0.10$ (SE
172 = 0.03), $t(11.9) = 3.03$, $p = 0.011$ (Fig. 2). This result suggests that imitations became
173 more stable (i.e., easier to imitate with high fidelity) with each generation of repetition.

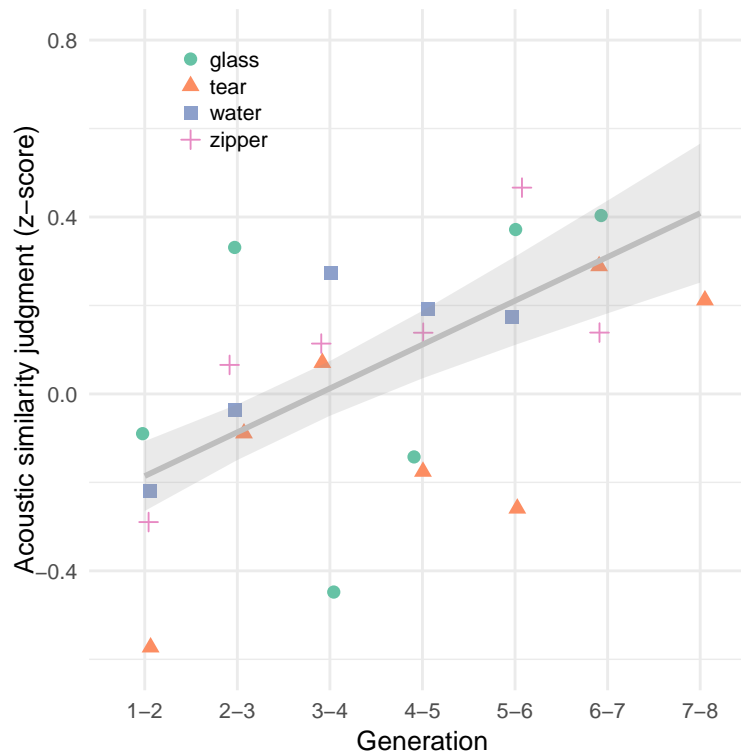


Figure 2. Change in perception of acoustic similarity over generations of iterated imitation. Points depict mean acoustic similarity ratings for pairs of imitations in each category. The predictions of the linear mixed-effects model are shown with ± 1 SE. Acoustic similarity increased over generations, indicating that repetition made the vocalizations easier to imitate with high fidelity.

174 **Acoustic similarity was highest within transmission chains.** Increasing
175 similarity along transmission chains could also reflect the continuous degradation of the
176 signal due to repeated imitation, in which case we would expect acoustic similarity to
177 increase both within as well as between transmission chains as a function of generation of
178 imitation. To rule out this alternative explanation, we calculated MFCCs for pairs of sounds
179 sampled from within and between different transmission chains from consecutive generations
180 across categories. To analyze the results, we fit a linear model predicting normalized acoustic
181 similarity scores (z-scores) from the generation of sounds. A hierarchical model was not
182 appropriate for this analysis because the between-chain pairs of sounds were sampled from

183 different categories, preventing any random effects due to category or seed from being
184 included in the model. We found that acoustic similarity increased within chains more than
185 it increased between chains, $b = -0.07$ (SE = 0.03), $t(6674.0) = -2.13$, $p = 0.033$ (Fig. S2).
186 This result supports the conclusion that transmission chains were stabilizing on divergent
187 acoustic forms as opposed to all chains converging on similar forms through continuous
188 degradation.

189 **Later generation imitations were transcribed more consistently.** An
190 additional test of stabilization and word-likeness was to measure whether later generation
191 imitations were transcribed more consistently than first generation imitations. We collected
192 a total of 2163 transcriptions — approximately 20 transcriptions per sound. Of these, 179
193 transcriptions (8%) were removed because they contained English words. Some examples of
194 the final transcriptions are presented in Table 1.

195 To measure the similarity among transcriptions, we calculated the orthographic
196 distance between the most frequent transcription and all other transcriptions of a given
197 imitation. The orthographic distance measure was a ratio based on longest contiguous
198 matching subsequences between pairs of transcriptions. We then fit a hierarchical linear
199 model predicting orthographic distance from the generation of the imitation (First
200 generation, Last generation) with random effects (intercepts and slopes) for seed sound
201 nested within category¹. The results showed that transcriptions of last generation imitations
202 were more similar to one another than transcriptions of first generation imitations, $b = -0.12$
203 (SE = 0.03), $t(3.0) = -3.62$, $p = 0.035$ (Fig. 3). The same result is reached through
204 alternative measures of orthographic distance, such as the percentage of exact transcription
205 matches for each imitation, $b = 0.10$ (SE = 0.03), $t(90.0) = 2.84$, $p = 0.006$, and the length
206 of the longest matching substring, $b = 0.98$ (SE = 0.24), $t(15.1) = 4.14$, $p < 0.001$ (Fig. S3).

¹Random effects for subject were not appropriate because the distance measure was derived from pairwise comparisons of transcriptions generated by different transcribers. As a result, the degrees of freedom for the significance tests for the parameters of this model reflect the Satterthwaite approximation based on the number of seed sounds (16) nested within categories (4), not the number of unique transcribers ($N=216$).

207 Differences between transcriptions of human vocalizations and transcriptions directly of
208 environmental sounds are presented in the Supplementary Materials (Fig. S5).

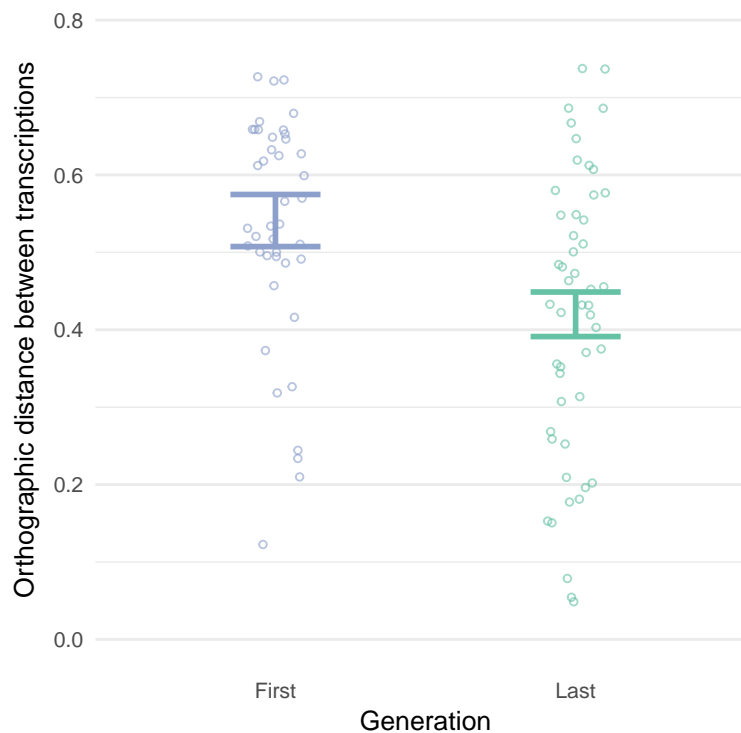


Figure 3. Orthographic agreement among transcriptions of first and last generation imitations. Points depict the mean orthographic distance between the most frequent transcription and all other transcriptions of a given imitation, with error bars denoting ± 1 SE of the hierarchical linear model predictions. Transcriptions of later generation imitations were more similar to one another than transcriptions of first generation imitations, suggesting that repeating imitations made them easier to transcribe into English orthography than direct imitations of environmental sounds.

209 Discussion

210 Repeating imitations of environmental sounds over generations of unique speakers was
211 sufficient to create more wordlike forms even without any instruction to do so. We defined
212 wordlike-ness in terms of acoustic stability and orthographic agreement. With additional
213 repetitions, the acoustic forms of the imitations became more similar to one another,

214 indicating they became easier to repeat with high fidelity. The possibility that this similarity
215 was due to uniform degradation across all transmission chains was ruled out by algorithmic
216 analyses of acoustic similarity within and between chains demonstrating that acoustic
217 similarity increased within chains but not between them. Additionally, later generation
218 imitations were transcribed more consistently into English orthography, further supporting
219 our hypothesis that repeating imitations makes them more word-like.

220 The results of Experiment 1 demonstrate the ease with which iterated imitation gives
221 rise to unique word forms. However, the results do not address how these emergent words
222 relate to the original sounds that were being imitated. As the imitations became more
223 word-like, were they stabilizing on arbitrary acoustic and orthographic forms, or did they
224 maintain some resemblance to the environmental sounds that motivated them? The purpose
225 of Experiment 2 was to assess the extent to which repeated imitations and their
226 transcriptions maintained a resemblance to the original set of seed sounds.

227 **Experiment 2: Resemblance of imitations to original seed sounds**

228 To assess the resemblance of repeated imitations to the original seed sounds, we
229 measured the ability of participants naïve to the design of the experiment to match
230 imitations and their transcriptions back to their original sound source relative to other seed
231 sounds from either the same category or from different categories (Fig. 4). We used match
232 accuracies to answer two questions concerning the effect of iterated imitation on resemblance
233 to the original seed sounds. First, we asked whether and for how many generations the
234 imitations and their transcriptions could be matched back to the original sounds. Second, we
235 asked whether repeated imitation resulted in a uniform degradation of the signal in each
236 imitation, or if repeated imitation resulted in some kinds of information degrading more
237 rapidly than others. Specifically, we tested the hypothesis that if imitations were becoming
238 more word-like, then they should also be interpreted more categorically, and thus we
239 predicted that the imitations might lose individuating information that identifies the specific

240 source of an imitation more rapidly than category information that identifies the general
241 category of environmental sound being imitated.

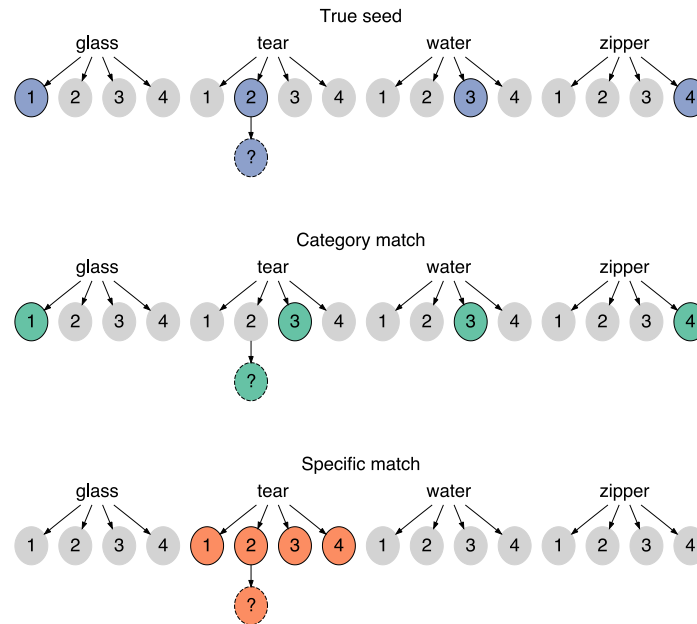


Figure 4. Three types of matching questions used to assess the resemblance between the imitation (and transcriptions of imitations) and the original seed sounds. For each question, participants listened an imitation (dashed circles) or read a transcription of one, and had to guess which of 4 sound choices (solid circles) they thought the person was trying to indicate. True seed questions contained the specific sound that generated the imitation as one of the choices (the correct response). The remaining sound choices were sampled from different categories. Category match questions replaced the original seed sound with another sound from the same category. Specific match questions pitted the actual seed against the other seeds within the same category.

242 Methods

243 **Matching imitations to seed sounds.** Participants ($N=751$) recruited from
244 Amazon Mechanical Turk were paid to listen to imitations, one at a time, and for each one,
245 choose one of four possible sounds they thought the person was trying to imitate. The task

246 was unspeeded and no feedback was provided. Participants completed 10 questions at a time.

247 All 365 imitations were tested in each of the three question types depicted in Fig. 4.
248 These questions differed in the relationship between the imitation and the four seed sounds
249 provided as the choices in the question. Question types (True seed, Category match, Specific
250 match) were assigned between-subject. Participants in the True seed and Category match
251 conditions were provided four seed sounds from different categories as choices in each
252 question. Participants in the Specific match condition were provided four seed sounds from
253 the same category.

254 **Matching transcriptions to seed sounds.** Participants ($N=468$) recruited from
255 Amazon Mechanical Turk completed a modified version of the matching survey described
256 above. Instead of listening to imitations, participants now read a word (a transcription of an
257 imitation), which they were told was an invented word. They were instructed that the word
258 was invented to describe one of the four presented sounds, and they had to guess which one.
259 The distractors for all questions were between-category, i.e. true seed and category match.
260 Specific match questions were omitted.

261 Of the unique transcriptions that were generated for each sound (imitations and seed
262 sounds), only the top four most frequent transcriptions were used in the matching
263 experiment. Participants who failed a catch trial ($N=6$) were excluded, leaving 461
264 participants in the final sample.

265 Results

266 **Imitations retained category information more than individuating**
267 **information.** Response accuracies in matching imitations to seed sounds were fit by a
268 generalized linear mixed-effects model predicting match accuracy as different from chance
269 (25%) based on the type of question being answered (True seed, Category match, Specific
270 match) and the generation of the imitation. Question types were contrast coded using
271 Category match questions as the baseline condition in comparison to the other two question

272 types each containing the actual seed that generated the imitation as one of the choices. The
273 model included random intercepts for participant², and random slopes and intercepts for
274 seed sounds nested within categories.

275 Accuracy in matching imitations to seed sounds was above chance for all question
276 types for the first generation of imitations, $b = 1.65$ (SE = 0.14) log-odds, odds = 0.50, $z =$
277 11.58, $p < 0.001$, and decreased steadily over generations, $b = -0.16$ (SE = 0.04) log-odds, $z =$
278 -3.72 , $p < 0.001$. We then tested whether this increase in difficulty was constant across the
279 three types of questions or if some question types became more difficult than others. The
280 results are shown in Fig. 5A. Performance decreased over generations more rapidly for
281 questions requiring a within-category distinction than for between-category questions, $b =$
282 -0.08 (SE = 0.03) log-odds, $z = -2.68$, $p = 0.007$, suggesting that between-category
283 information was more resistant to loss through repeated imitation.

284 An alternative explanation for this result is that the within-category match questions
285 are simply more difficult because the sounds provided as choices are more acoustically
286 similar to one another than the between-category questions, and therefore, performance
287 might be expected to drop off more rapidly with repeated imitation for these more difficult
288 questions³. However, performance also decreased for the easiest type of question where the
289 correct answer was the actual seed generating the imitation (True seed questions; see Fig. 4);
290 the advantage of having the true seed among between-category distractors decreased over
291 generations, $b = -0.07$ (SE = 0.02) log-odds, $z = -2.77$, $p = 0.006$. The observed increase in
292 the “category advantage” (i.e., the advantage of having between-category distractors)

²Random slopes for generation were not appropriate in the by-subject random effects because data collection was batched by generation of imitation, and therefore each participant did not sample across the range of generations.

³We observed that performance on some Specific match questions dropped below chance for later generations indicating participants had an apparent aversion to the nominally correct answer. Additional analyses showed that participants were not converging on a single incorrect response. The reason for this pattern is at present unclear. Removing these trials from the analysis does not substantively change the conclusions.

293 combined with a decrease in the “true seed advantage” (the advantage of having the actual
294 seed among the choices), shows that the changes induced by repeated imitation caused the
295 imitations to lose some of properties that linked the earlier imitations to the specific sound
296 that motivated them, while nevertheless preserving a more abstract category-based
297 resemblance.

298 **Transcriptions retained information about seed sources.** We next report the
299 results of matching the written transcriptions of the auditory sounds back to the original
300 environmental sounds. Remarkably, participants were able to guess the correct meaning of a
301 word that was transcribed from an imitation that had been repeated up to 8 times, $b = 0.83$
302 (SE = 0.13) log-odds, odds = -0.18, $z = 6.46$, $p < 0.001$ (Fig. 5B). This was true for True
303 seed questions containing the actual seed generating the transcribed imitation, $b = 0.75$ (SE
304 = 0.15) log-odds, $z = 4.87$, $p < 0.001$, and for Category match questions where participants
305 had to associate transcriptions with a particular category of environmental sounds, $b = 1.02$
306 (SE = 0.16) log-odds, $z = 6.39$, $p < 0.001$. The effect of generation did not vary across these
307 question types, $b = 0.05$ (SE = 0.10) log-odds, $z = 0.47$, $p = 0.638$. The results of matching
308 “transcriptions” directly of the environmental sounds are shown in Fig. S5.

309 Discussion

310 Even after being repeated up to 8 times, imitations retained a resemblance to the
311 environmental sound that motivated them, even after being transcribed into orthographic
312 forms. For imitations, but not for transcriptions, this resemblance was stronger for the
313 category of environmental sound than the actual seed sound, suggesting that through
314 repetition, the imitations were becoming more categorical. This result supports the results of
315 Experiment 1 in demonstrating another aspect of wordlike-ness achieved through repeated
316 imitation: Words, in addition to being stable in acoustic and orthographic forms, are also
317 categorical, denoting all members of a category equally as opposed to identifying individual
318 category members. Repeating imitations of environmental sounds is sufficient to remove

319 some of the individuating characteristics of the imitation while retaining a category-based
320 resemblance.

321 The reason the same effect was not observed in matching accuracy for transcriptions is
322 unknown. One possible reason is that the process of transcribing a non-linguistic
323 vocalization into a written word encourages transcribers to emphasize individuating
324 information about the vocalization. However, the fact that transcriptions of imitations can
325 be matched back to other category members (Category match questions) suggests that
326 transcriptions are still carrying some category information. Another possible reason is that
327 by subsetting the most frequent transcriptions, we unintentionally excluded less frequent
328 transcriptions that were more diagnostic of category information.

329 Experiments 1 and 2 document a process of gradual change from an imitation of an
330 environmental sound to a more wordlike form. But do these emergent words function like
331 other words in the language? In Experiment 3, we test the suitability of words taken from
332 the beginning and end of transmission chains in serving as category labels in a category
333 learning task.

334 **Experiment 3: Suitability of created words as category labels**

335 One consequence of imitations becoming more word-like is that they may make for
336 better category labels. For example, an imitation from a later generation, by virtue of having
337 a more word-like form, may be easier to learn as a label for the category of sounds that
338 motivated it than an earlier imitation, which is more closely yoked to a particular
339 environmental sound. To the extent that repeating imitations abstracts away the
340 idiosyncrasies of a particular category member (Edmiston & Lupyan, 2015; Lupyan &
341 Thompson-Schill, 2012), it may also be easier to generalize to new category members. We
342 tested these predictions using a category learning task in which participants learned novel
343 labels as category labels of the seed environmental sounds. The novel labels were
344 transcriptions of either first or last generation imitations gathered in Experiment 1.

345 **Methods**

346 **Selecting words to learn as category labels.** Our transmission chain design and
347 subsequent transcription procedure created 1814 unique words. From these, we sampled
348 words transcribed from first and last generation imitations, as well as transcriptions of the
349 original seed sounds. Our procedure for sampling transcriptions to use as category labels was
350 as follows: First, we removed transcriptions that contained less than 3 unique characters and
351 transcriptions that were over 10 characters long. Of the remaining transcriptions, a sample
352 of 56 were selected that were approximately equally associated with the target category. To
353 measure the association between each imitation and its target category (the category of the
354 seed sound), we used the match accuracy scores reported in Experiment 2. The reason for
355 using this measure of association strength as a control for selecting words to learn as
356 category labels was to be able to select words that were initially equally associated with the
357 target categories. Equating along this dimension allowed for a more focused test of
358 differences between the words in terms of generalization to new category members. The final
359 sample of transcriptions were selected using a bootstrapping procedure which involved
360 selecting a desired mean (the average association strength for eligible transcriptions of last
361 generation imitations) and sampling transcriptions from first generation imitations and from
362 seed sounds until the match accuracy of those imitations matched the desired mean within 1
363 standard deviation.

364 **Procedure.** Participants ($N=67$) were University of Wisconsin undergraduates who
365 received course credit for participation. Participants were randomly assigned four novel
366 labels to learn for four categories of environmental sounds. Full instructions are provided in
367 the Supplementary Materials. Participants were assigned between-subject to learn labels
368 (transcriptions) of either first or last generation imitations. Some participants learned labels
369 from transcriptions of seed sounds (Fig. S6). On each trial, participants heard one of the 16
370 seed sounds. After a 1s delay, participants saw a label (one of the transcribed imitations)
371 and responded yes or no using a gamepad controller depending on whether the sound and

372 the word went together. Participants received accuracy feedback (a bell sound and a green
373 checkmark if correct; a buzzing sound and a red “X” if incorrect). Four outlier participants
374 were excluded from the final sample due to high error rates and slow RTs.

375 Participants categorized all 16 seed sounds over the course of the experiment, but they
376 learned them in blocks of 4 sounds at a time. Within each block of 24 trials, participants
377 heard the same four sounds and the same four words multiple times, with a 50% probability
378 of the sound matching the word on any given trial. At the start of a new block of trials,
379 participants heard four new sounds they had not heard before, and had to learn to associate
380 these new sounds with the words they had learned in the previous blocks.

381 Results

382 Later generation transcriptions yielded more efficient responding.

383 Participants began by learning through trial-and-error to associate four written labels with
384 four categories of environmental sounds. The small number of categories made this an easy
385 task (mean accuracy after the first block of 24 trials was 81%; Fig. S4). Participants
386 learning transcriptions of first or last generation imitations did not differ in overall accuracy,
387 $p = 0.887$, or reaction time, $p = 0.616$. After this initial learning phase (i.e. after the first
388 block of trials), accuracy performance quickly reached ceiling and did not differ between
389 groups $p = 0.775$. However, the response times of participants learning last generation
390 transcriptions declined more rapidly with practice than participants learning first generation
391 transcriptions, $b = -114.13$ (SE = 52.06), $t(39.9) = -2.19$, $p = 0.034$ (Fig. 6A). These faster
392 responses suggest that, in addition to becoming more stable both in terms of acoustic and
393 orthographic properties, repeating imitations makes them easier to process as category labels.
394 We predict that given a harder task (i.e., more than four categories and 16 exemplars) would
395 yield differences in initial learning rates as well.

396 **Later generation transcriptions were better generalized.** Next, we examined
397 whether transcriptions from last generation imitations were easier to generalize to novel

398 category exemplars. To test this hypothesis, we compared RTs on trials immediately prior to
399 the introduction of novel sounds (new category members) and the first trials after the block
400 transition (± 6 trials). The results revealed a reliable interaction between the generation of
401 the transcribed imitation and the block transition, $b = -110.77$ (SE = 52.84), $t(39.7) = -2.10$,
402 $p = 0.042$ (Fig. 6B). This result suggests that transcriptions from later generation imitations
403 were easier to generalize to new category members.

404 Discussion

405 The results of a simple category learning experiment demonstrate a possible benefit to
406 the stabilization of repeated imitations on more wordlike forms. As a consequence of being
407 more wordlike, repeated imitations were responded to more quickly, and generalized to new
408 category members more easily. These results suggest an advantage to repeating imitations
409 from the perspective of the language learner in that they afford better category
410 generalization.

411 General Discussion

412 Imitative words are found across the spoken languages of the world (Dingemanse et al.,
413 2015; Imai & Kita, 2014; Perniss et al., 2010). Counter to past assumptions about the
414 limitations of human vocal imitation, people are surprisingly effective at using vocal
415 imitation to represent and communicate about the sounds in their environment (Lemaitre et
416 al., 2016) and more abstract meanings (Perlman et al., 2015), making the hypothesis that
417 early spoken words originated from imitations a plausible one. We examined whether simply
418 repeating an imitation of an environmental sound—with no intention to create a new word
419 or even to communicate—produces more word-like forms.

420 Our results show that through simple repetition, imitative vocalizations became more
421 word-like both in form and function. In form, the vocalizations gradually stabilized over
422 generations, becoming more similar from imitation to imitation. They also became
423 increasingly standardized in accordance with English orthography, as later generations were

424 more consistently transcribed into English words, providing converging evidence of
425 stabilization. In function, the increasingly word-like forms became more effective as category
426 labels. In a category learning experiment, naïve participants were faster at matching
427 category labels derived from later-generation imitations than those derived directly from
428 imitations of environmental sounds. This fits with previous research showing that the
429 relatively arbitrary forms that are typical of words (e.g. “dog”) makes them better suited to
430 function as category labels compared to direct auditory cues (Boutonnet & Lupyan, 2015;
431 Edmiston & Lupyan, 2015; e.g. the sound of a dog bark; Lupyan & Thompson-Schill, 2012).

432 Even as the vocalizations became more word-like, they nevertheless maintained an
433 imitative quality. After eight generations they could no longer be matched to the particular
434 sound from which they originated any more accurately than they could be matched to the
435 general category of environmental sound. Thus, information that distinguished an imitation
436 from other sound categories was more resilient to transmission decay than exemplar
437 information within a category. Remarkably, even after the vocalizations were transcribed
438 into English orthography, participants were able to guess their original sound category from
439 the written “words”. In contrast to the vocalizations, participants continued to be more
440 accurate at matching late generation transcriptions back to their particular source sound
441 relative to other exemplars from the same category.

442 Although the number of imitative words in contemporary languages may appear to be
443 very small (Crystal, 1987; Newmeyer, 1992), increasing evidence from disparate languages
444 shows that vocal imitation is, in fact, a widespread source of vocabulary. Cross-linguistic
445 surveys indicate that onomatopoeia—imitative words used to represent sounds—are a
446 universal lexical category found across the world’s languages (Dingemanse, 2012). Even
447 English, a language that has been characterized as relatively limited in iconic vocabulary
448 (Vigliocco, Perniss, & Vinson, 2014), is documented as having hundreds of clearly imitative
449 words including words for human and animal vocalizations as well as various types of
450 environmental sounds (Rhodes, 1994; Sobkowiak, 1990). Besides words that are directly

451 imitative of sounds—the focus of the present study — many languages contain semantically
452 broader inventories of ideophones. These words comprise a grammatically and phonologically
453 distinct class of words that are used to express various sensory-rich meanings, such as
454 qualities related to manner of motion, visual properties, textures and touch, inner feelings
455 and cognitive states (Dingemanse, 2012; Nuckolls, 1999; Voeltz & Kilian-Hatz, 2001). As
456 with onomatopoeia, ideophones are often recognized by naïve speakers as bearing a degree of
457 resemblance to their meaning (Dingemanse, Schuerman, & Reinisch, 2016).

458 Our study focused on imitations of environmental sounds and more work remains to be
459 done to determine the extent to which vocal imitation can ground de novo vocabulary
460 creation in other semantic domains (Lupyan & Perlman, 2015; e.g., Perlman et al., 2015).
461 What the present results make clear is that the transition from imitation to word can be a
462 rapid and simple process: the mere act of iterated imitation can drive vocalizations to
463 become more word-like in both form and function. Notably, just as onomatopoeia and
464 ideophones of natural languages maintain a resemblance to the quality they represent, the
465 present vocal imitations transitioned to words while retaining a resemblance to the original
466 sound that motivated them.

467 **References**

- 468 Arbib, M. A. (2012). *How the brain got language: The mirror system hypothesis* (Vol. 16).
469 Oxford University Press.
- 470 Armstrong, D. F., & Wilcox, S. (2007). *The gestural origin of language*. Oxford University
471 Press.
- 472 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects
473 Models Using lme4. *Journal of Statistical Software*, *67*(1), 1–48.
- 474 Boutonnet, B., & Lupyan, G. (2015). Words Jump-Start Vision: A Label Advantage in
475 Object Recognition. *Journal of Neuroscience*, *35*(25), 9329–9335.
- 476 Brown, R. W., Black, A. H., & Horowitz, A. E. (1955). Phonetic symbolism in natural

- 477 languages. *Journal of Abnormal Psychology*, 50(3), 388–393.
- 478 Clark, H. H., & Gerrig, R. J. (1990). Quotations as demonstrations. *Language*, 66, 764–805.
- 479 Corballis, M. C. (2003). *From hand to mouth: The origins of language*. Princeton University
480 Press.
- 481 Crystal, D. (1987). *The Cambridge Encyclopedia of Language* (Vol. 2). Cambridge Univ
482 Press.
- 483 Dingemanse, M. (2012). Advances in the Cross-Linguistic Study of Ideophones. *Language
484 and Linguistics Compass*, 6(10), 654–672.
- 485 Dingemanse, M. (2014). Making new ideophones in Siwu: Creative depiction in conversation.
486 *Pragmatics and Society*.
- 487 Dingemanse, M., Blasi, D. E., Lupyan, G., Christiansen, M. H., & Monaghan, P. (2015).
488 Arbitrariness, Iconicity, and Systematicity in Language. *Trends in Cognitive Sciences*,
489 19(10), 603–615.
- 490 Dingemanse, M., Schuerman, W., & Reinisch, E. (2016). What sound symbolism can and
491 cannot do: Testing the iconicity of ideophones from five languages. *Language*, 92.
- 492 Donald, M. (2016). Key cognitive preconditions for the evolution of language. *Psychonomic
493 Bulletin & Review*, 1–5.
- 494 Edmiston, P., & Lupyan, G. (2015). What makes words special? Words as unmotivated cues.
495 *Cognition*, 143(C), 93–100.
- 496 Gamer, M., Lemon, J., Fellows, I., & Singh, P. (2012). *irr: Various Coefficients of Interrater
497 Reliability and Agreement*.
- 498 Goldin-Meadow, S. (2016). What the hands can tell us about language emergence.
499 *Psychonomic Bulletin & Review*, 24(1), 1–6.
- 500 Hall, K. C., Allen, B., Fry, M., Mackie, S., & McAuliffe, M. (2016). Phonological
501 CorpusTools. *14th Conference for Laboratory Phonology*.
- 502 Hewes, G. W. (1973). Primate Communication and the Gestural Origin of Language.

- 503 *Current Anthropology*, 14(1/2), 5–24.
- 504 Hockett, C. F. (1978). In search of Jove's brow. *American Speech*, 53(4), 243–313.
- 505 Imai, M., & Kita, S. (2014). The sound symbolism bootstrapping hypothesis for language
506 acquisition and language evolution. *Philosophical Transactions of the Royal Society*
507 *B: Biological Sciences*, 369(1651).
- 508 Kendon, A. (2014). Semiotic diversity in utterance production and the concept of 'language'.
509 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 369(1651),
510 20130293–20130293.
- 511 Klima, E. S., & Bellugi, U. (1980). *The signs of language*. Harvard University Press.
- 512 Kuznetsova, A., Bruun Brockhoff, P., & Haubo Bojesen Christensen, R. (2016). *lmerTest:*
513 *Tests in Linear Mixed Effects Models*.
- 514 Lemaitre, G., & Rocchesso, D. (2014). On the effectiveness of vocal imitations and verbal
515 descriptions of sounds. *The Journal of the Acoustical Society of America*, 135(2),
516 862–873.
- 517 Lemaitre, G., Houix, O., Voisin, F., Misdariis, N., & Susini, P. (2016). Vocal Imitations of
518 Non-Vocal Sounds. *PloS One*, 11(12), e0168167–28.
- 519 Lewis, J. (2009). As well as words: Congo Pygmy hunting, mimicry, and play. In *The cradle*
520 *of language*. The cradle of language.
- 521 Lupyan, G., & Perlman, M. (2015). The vocal iconicity challenge! In *The th biennial*
522 *protolanguage conference*. Rome, Italy.
- 523 Lupyan, G., & Thompson-Schill, S. L. (2012). The evocative power of words: Activation of
524 concepts by verbal and nonverbal means. *Journal of Experimental Psychology:*
525 *General*, 141(1), 170–186.
- 526 Newmeyer, F. J. (1992). Iconicity and generative grammar. *Language*.
- 527 Nuckolls, J. B. (1999). The case for sound symbolism. *Annual Review of Anthropology*,
528 28(1), 225–252.
- 529 Perlman, M., Dale, R., & Lupyan, G. (2015). Iconicity can ground the creation of vocal

- 530 symbols. *Royal Society Open Science*, 2(8), 150152–16.
- 531 Perniss, P., Thompson, R. L., & Vigliocco, G. (2010). Iconicity as a General Property of
532 Language: Evidence from Spoken and Signed Languages. *Frontiers in Psychology*, 1.
- 533 Pinker, S., & Jackendoff, R. (2005). The faculty of language: what's special about it?
534 *Cognition*, 95(2), 201–236.
- 535 Rhodes, R. (1994). Aural images. *Sound Symbolism*, 276–292.
- 536 Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: uses in assessing rater reliability.
537 *Psychological Bulletin*, 86(2), 420–428.
- 538 Sobkowiak, W. (1990). On the phonostatistics of English onomatopoeia. *Studia Anglica*
539 *Posnaniensia*, 23, 15–30.
- 540 Tomasello, M. (2010). *Origins of human communication*. MIT press.
- 541 Vigliocco, G., Perniss, P., & Vinson, D. (2014). Language as a multimodal phenomenon:
542 implications for language learning, processing and evolution. *Philosophical*
543 *Transactions of the Royal Society B: Biological Sciences*, 369(1651),
544 20130292–20130292.
- 545 Voeltz, F. E., & Kilian-Hatz, C. (2001). *Ideophones* (Vol. 44). John Benjamins Publishing.

Table 1

Examples of words transcribed from imitations.

Category	Seed	First generation	Last generation
glass	1	tingtingting	deetdedededeet
glass	2	chirck	correcto
glass	3	dirrng	wayew
glass	4	boonk	baroke
tear	1	scheeept	cheecheea
tear	2	feeshefee	cheeoooo
tear	3	hhhweerrr	chhhhhhewwwe
tear	4	ccccchhhhyeahh	shhhhh
water	1	boococucuwich	eeverlusha
water	2	chwoochwoochwooo	cheiopshpshcheiopsh
water	3	atoadelchoo	mowah
water	4	awakawush	galonggalong
zipper	1	euah	izoo
zipper	2	zoop	veeeep
zipper	3	arrgt	owww
zipper	4	bzzzzup	izzip

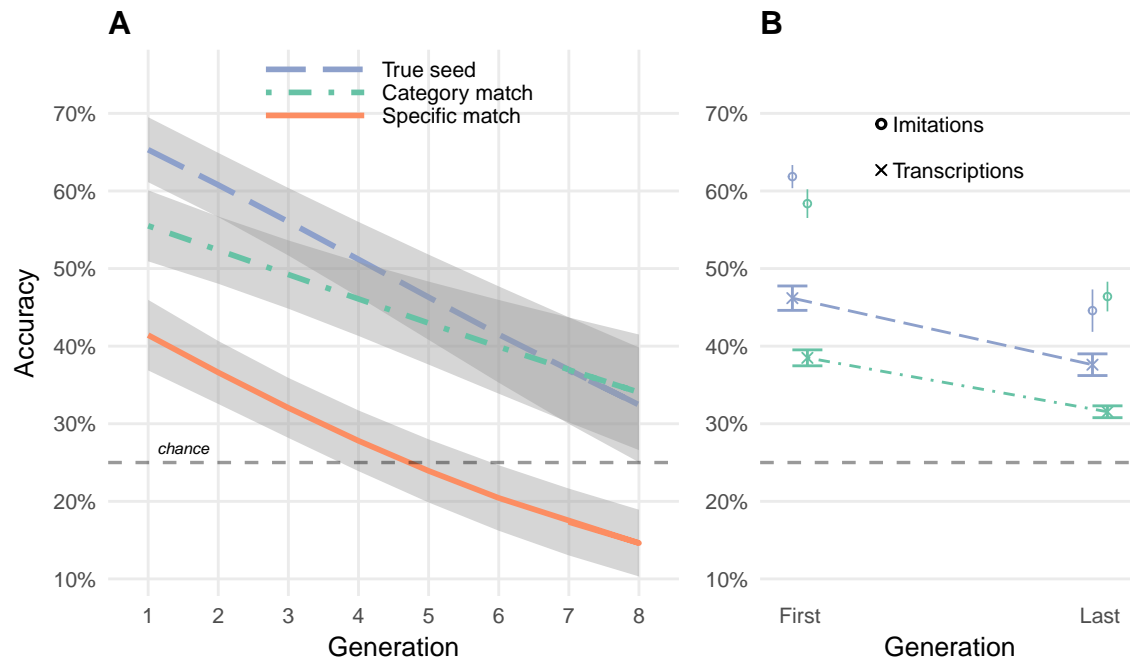


Figure 5. Repeated imitations retained category resemblance. A. Accuracy of matching vocal imitations to original seed sounds as a function of the generation during which the imitation was produced. Curves show predictions of the generalized linear mixed effects models with ± 1 SE of the model predictions. The “category advantage” (Category match vs. Specific match) increased over generations, while the “true seed advantage” (True seed v. Category match) decreased (see main text), suggesting that imitations lose within-category information more rapidly than between-category information. B. Accuracy of matching transcriptions of the imitations to original seed sounds (e.g., “boococucuwich” to a water splashing sound). Transcriptions of imitations could still be matched back to the category of sound that motivated the original imitation even after 8 generations. Circles show mean matching accuracy for the corresponding vocal imitations for comparison.

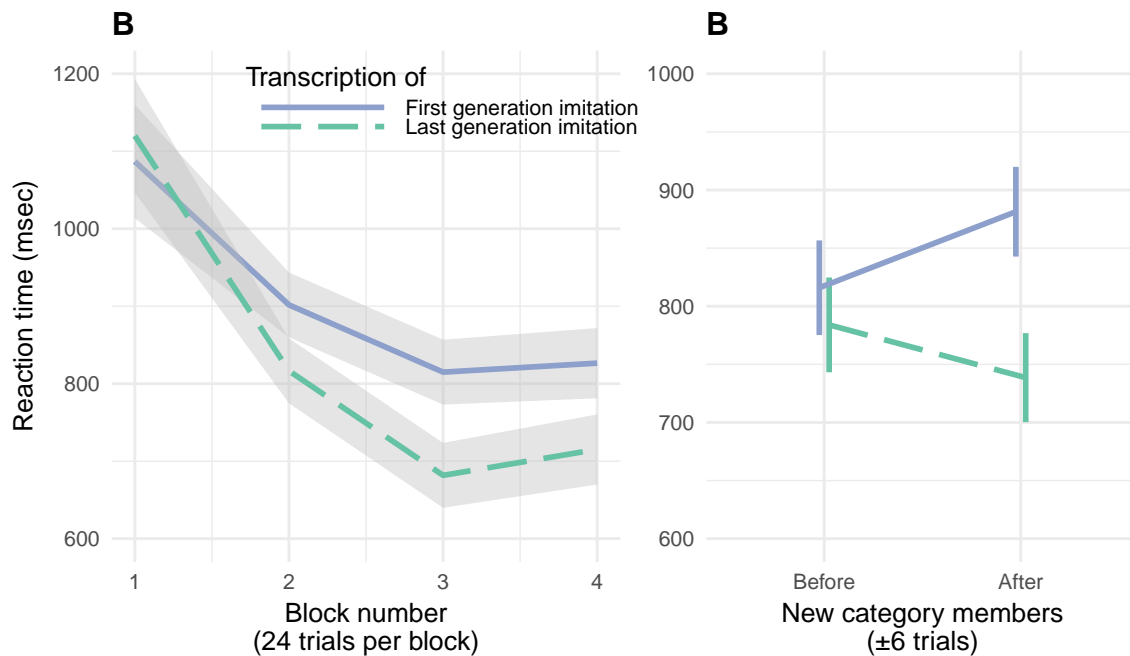


Figure 6. Repeated imitations made for better category labels. Participants learned novel labels (transcriptions of first or last generation imitations) for categories of environmental sounds. A. Mean RTs for correct responses in the category learning experiment with ± 1 SE. Participants achieved faster RTs in matching transcribed labels to environmental sounds for labels transcribed from later compared to earlier generation imitations. B. Cost of generalizing to new category members with ± 1 SE. After each block of trials, new environmental sounds were introduced, requiring participants to generalize the previously learned category labels to new category members. There was a generalization cost for the first generation labels, but not the last generation labels.

546

Table captions

547 *Table 1.* Examples of words transcribed from imitations.

548

Figure captions

549 *Figure 1.* Vocal imitations collected in the transmission chain experiment. Seed
550 sounds (16) were sampled from four categories of environmental sounds:
551 glass, tear, water, zipper. Participants imitated each seed sound, and
552 then the next generation of participants imitated the imitations, and
553 so on, for up to 8 generations. Chains are unbalanced due to random
554 assignment and the exclusion of some low quality recordings.

555 *Figure 2.* Change in perception of acoustic similarity over generations of iterated
556 imitation. Points depict mean acoustic similarity ratings for pairs of
557 imitations in each category. The predictions of the linear mixed-effects
558 model are shown with ± 1 SE. Acoustic similarity increased over genera-
559 tions, indicating that repetition made the vocalizations easier to imitate
560 with high fidelity.

561 *Figure 3.* Orthographic agreement among transcriptions of first and last generation
562 imitations. Points depict the mean orthographic distance between the
563 most frequent transcription and all other transcriptions of a given imi-
564 tation, with error bars denoting ± 1 SE of the hierarchical linear model
565 predictions. Transcriptions of later generation imitations were more
566 similar to one another than transcriptions of first generation imitations,
567 suggesting that repeating imitations made them easier to transcribe into
568 English orthography than direct imitations of environmental sounds.

569 *Figure 4.* Three types of matching questions used to assess the resemblance between
570 the imitation (and transcriptions of imitations) and the original seed
571 sounds. For each question, participants listened an imitation (dashed
572 circles) or read a transcription of one, and had to guess which of 4 sound
573 choices (solid circles) they thought the person was trying to indicate.
574 True seed questions contained the specific sound that generated the
575 imitation as one of the choices (the correct response). The remaining
576 sound choices were sampled from different categories. Category match
577 questions replaced the original seed sound with another sound from the
578 same category. Specific match questions pitted the actual seed against
579 the other seeds within the same category.

580 *Figure 5.* Repeated imitations retained category resemblance. A. Accuracy of
581 matching vocal imitations to original seed sounds as a function of the
582 generation during which the imitation was produced. Curves show pre-
583 dictions of the generalized linear mixed effects models with ± 1 SE of the
584 model predictions. The “category advantage” (Category match vs. Spe-
585 cific match) increased over generations, while the “true seed advantage”
586 (True seed v. Category match) decreased (see main text), suggesting that
587 imitations lose within-category information more rapidly than between-
588 category information. B. Accuracy of matching transcriptions of the
589 imitations to original seed sounds (e.g., “boococuwich” to a water
590 splashing sound). Transcriptions of imitations could still be matched
591 back to the category of sound that motivated the original imitation
592 even after 8 generations. Circles show mean matching accuracy for the
593 corresponding vocal imitations for comparison.

594 *Figure 6.* Repeated imitations made for better category labels. Participants learned
595 novel labels (transcriptions of first or last generation imitations) for
596 categories of environmental sounds. A. Mean RTs for correct responses
597 in the category learning experiment with ± 1 SE. Participants achieved
598 faster RTs in matching transcribed labels to environmental sounds for
599 labels transcribed from later compared to earlier generation imitations.
600 B. Cost of generalizing to new category members with ± 1 SE. After each
601 block of trials, new environmental sounds were introduced, requiring
602 participants to generalize the previously learned category labels to new
603 category members. There was a generalization cost for the first generation
604 labels, but not the last generation labels.