

1 **Chestnut-crowned babbler calls are composed of meaningless shared building blocks**

2  
3 Sabrina Engesser <sup>a,b\*</sup>, Jennifer L. Holub <sup>c</sup>, Louis G. O'Neill <sup>d,e</sup>, Andrew F. Russell <sup>e†</sup> &

4 Simon W. Townsend <sup>a,b,f†</sup>

5  
6 <sup>a</sup> Department of Comparative Linguistics, University of Zurich, 8032 Zurich, Switzerland

7 <sup>b</sup> Center for the Interdisciplinary Study of Language Evolution, University of Zurich, 8032  
8 Zurich, Switzerland

9 <sup>c</sup> Fowlers Gap Arid Zone Research Station, School of Biological, Earth & Environmental  
10 Sciences, University of New South Wales, NSW 2052, Australia

11 <sup>d</sup> Department of Biological Sciences, Macquarie University, NSW 2109, Australia

12 <sup>e</sup> Centre for Ecology & Conservation, College of Life & Environmental Sciences, University of  
13 Exeter, Penryn, Cornwall TR10 9FE, United Kingdom

14 <sup>f</sup> Department of Psychology, University of Warwick, Coventry CV4 7AL, United Kingdom

15  
16 † These authors share last authorship

17  
18 \* Corresponding author

19 Email: [sabrina.engesser@uzh.ch](mailto:sabrina.engesser@uzh.ch), phone: +41 (0)44 634 0223

20  
21 **Short title:** Building blocks of babbler multi-element calls

22  
23 **Keywords:** Language evolution, phonology, combinatoriality, vocal communication,  
24 habituation-discrimination

25 **Abstract**

26 A core component of human language is its combinatorial sound system: meaningful signals are  
27 built from different combinations of meaningless sounds. Investigating whether non-human  
28 communication systems are also combinatorial is hampered by difficulties in identifying the  
29 extent to which vocalizations are constructed from shared, meaningless building blocks. Here we  
30 present a novel approach to circumvent this difficulty and show that a pair of functionally distinct  
31 chestnut-crowned babbler (*Pomatostomus ruficeps*) vocalizations can be decomposed into  
32 perceptibly distinct, meaningless entities that are shared across the two calls. Specifically, by  
33 focusing on the acoustic distinctiveness of sound elements using a habituation-discrimination  
34 paradigm on wild-caught babblers under standardized aviary conditions, we show that two multi-  
35 element calls are composed of perceptibly distinct sounds that are reused in different  
36 arrangements across the two calls. Furthermore, and critically, we show that none of the five  
37 constituent elements elicits functionally relevant responses in receivers, indicating that the  
38 constituent sounds do not carry the meaning of the call; so are contextually meaningless. Our  
39 work, which allows combinatorial systems in animals to be more easily identified, suggests that  
40 animals can produce functionally distinct calls that are built in a way superficially reminiscent of  
41 the way that humans produce morphemes and words. The results reported lend credence to the  
42 recent idea that language's combinatorial system may have been preceded by a superficial stage  
43 where signalers neither needed to be cognitively aware of the combinatorial strategy in place, nor  
44 of its building blocks.

45 **Significance statement**

46 Word generation in human language is fundamentally based on the ability to use a finite set of  
47 meaningless sounds in different combinations across contexts. Investigating whether animals  
48 share this basic capacity has been hampered by difficulties in identifying the extent to which  
49 animal vocalizations can be decomposed into smaller meaningless, yet shared sounds. Using a  
50 novel implementation of habituation-discrimination experiments, we show for the first time that a  
51 pair of functionally distinct chestnut-crowned babbler (*Pomatostomus ruficeps*) vocalizations are  
52 composed of perceptibly distinct, contextually meaningless sounds that are shared across the  
53 different calls. We conclude that the individual sounds represent building blocks that generate  
54 meaning when combined in a particular way, akin to word formation in human language.

55 \body

## 56 **Introduction**

57 A universal feature of human language is its combinatorial structure: a finite set of perceptibly  
58 distinct, meaningless sounds (building blocks) can be productively recombined to create a  
59 theoretically limitless set of meaningful signals [1]. One way to elucidate candidate origins  
60 and/or early forms of the combinatorial feature of language is to test for analogues in the basic  
61 process that underpins combinatoriality in the vocalizations of non-human animals [2]. While  
62 animals are clearly able to communicate using combinatorial vocal signals [3-8], whether they  
63 use meaningless sound elements in different arrangements to generate new meaning is  
64 contentious [9, 10]. This contention stems from two sources. First, from ambiguous associations  
65 between sound arrangements and meaning: for example, although animal songs are often  
66 composed of smaller sound units in different arrangements, precise arrangements are not known  
67 to underpin context-specific, or ‘propositional’, meaning [10-12]. Second, it also stems from  
68 difficulties of identifying whether functionally distinct vocalizations can be comprised of a  
69 recombinatorial system of shared meaningless sounds (i.e. building blocks) [13-16].

70         The traditional approach used to deconstruct the building blocks of the combinatorial  
71 sound system of human language is through the analysis of minimal pairs: pairs of semantically  
72 distinct words that differ in a *single* meaningless sound element, for example ‘lap’ versus ‘tap’  
73 [9, 17]. The elements that differ in minimal pairs, in this case /t/ and /l/, are semantically  
74 meaningless, but are what serve to differentiate the meaning encoded in the two words. By  
75 extension, /t/ and /l/ must each represent distinct, meaning-contrasting sounds. This minimal pairs  
76 approach is feasible in human language because its sound elements are present in a plethora of  
77 permutations, such that each one used, and the role it plays in differentiating meaning, can be  
78 contrasted systematically with others in the repertoire [18]. However, this approach becomes  
79 unfeasible for communication systems where different sounds are not productively recombined

80 and occur in prohibitively few combinations to allow direct contrasts of the impacts of single  
81 sounds on meaning to be made. Given that the productive usage of different sounds is likely a  
82 derived language-specific trait and is not a known feature of animal communication [13], an  
83 alternative method is required to test whether functionally distinct vocal signals are built from  
84 recombinations of shared sounds that are meaningless in isolation – the hallmark of  
85 combinatoriality in human language.

86         We propose that testing whether individuals perceive sound elements within and across  
87 functionally distinct calls as acoustically different or equivalent can also serve to decompose the  
88 potential building blocks of an animal’s vocal system. Further, this approach can be implemented  
89 using established habituation-discrimination paradigms previously applied for speech-sound  
90 perception in human infants [19] and to assess the information content of whole calls in animals  
91 [20-23]. The utility of this habituation-discrimination approach to unpacking the characteristics  
92 of elements within calls is based on recent simulations on the emergence of combinatorial signals  
93 that define combinatorial structures using trajectories through acoustic and perceptual space [13,  
94 18, 24]. In such simulations, the distance between points along trajectories of acoustic space  
95 reflect confusion probabilities, and hence the perceptual discreteness of sound elements.  
96 Accordingly, sound elements that are so close in acoustic parameter space so as to be easily  
97 confused are in essence perceptibly equivalent, while those that are more distant and seldom  
98 confused are essentially distinct. The advantage of this approach is that by focusing on sound  
99 discrimination and sharing within and across functionally distinct calls, comparative work  
100 investigating whether animal signals are composed of meaningless, recombinatorial entities (or  
101 building blocks) becomes feasible; with the potential to shed important light on the origins of  
102 combinatoriality.

103         Our overall aim is to use this new approach to test whether a pair of structurally similar

104 but functionally distinct vocalizations of the chestnut-crowned babbler (*Pomatostomus ruficeps*;  
105 Fig 1A) can be decomposed into perceptibly distinct, contextually meaningless entities that are  
106 shared across the two calls – the defining feature of combinatoriality. The two calls of this highly  
107 social passerine bird from inland southeastern Australia [25] in question are: bi-element flight  
108 calls which are uttered when a bird flies off and which function to coordinate group movement  
109 (composed of the elements F<sub>1</sub>F<sub>2</sub>; Fig 1B); and tri-element prompt calls which are produced by an  
110 individual when entering the breeding nest in order to stimulate nestling begging during food  
111 provisioning (composed of the elements P<sub>1</sub>P<sub>2</sub>P<sub>3</sub>; Fig 1B) [26, 27]. The functional distinction  
112 between the two calls is confirmed in playbacks on wild birds in on-site aviaries: flight calls  
113 induce greater movement and looking outside the aviary, presumably in response to an  
114 anticipated incoming bird, while prompt calls induce an 8-fold increase in the amount of time  
115 spent looking at a nest placed inside the aviary, presumably because of the natural association  
116 between nests and prompt calls [14]. Further, none of the five elements in the two calls is known  
117 to be used as stand-alone calls despite >1000 h of recordings in all known socio-ecological  
118 contexts, and all differ significantly from uni-element short-distance contact calls used to  
119 maintain contact and spacing during feeding [26]. Finally, previous aviary playback experiments  
120 also suggested that the distinct meaning encoded in these two multi-element calls is generated by  
121 the specific arrangement of the constituent sound elements [14]. However, what is not known is  
122 whether or not the constituent elements within these multi-element calls are: (a) perceptibly  
123 distinct within calls; (b) perceptibly equivalent across calls; and (c) contextually meaningless.  
124 Each of these three facets is required to resolve whether functionally distinct calls are built from  
125 smaller, perceptibly distinct and shared, meaningless sounds.

126 To test these core components of combinatoriality, we used standardized aviary playbacks  
127 on wild-caught chestnut-crowned babblers: (i) to identify which of the five sound elements

128 constituting flight and prompt calls (i.e. F<sub>1</sub>, F<sub>2</sub>, P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>) are perceptibly distinct; (ii) to identify  
129 which, if any, are shared across the two calls; and (iii) to investigate whether contextually  
130 relevant information is encoded in the individual sound elements. To test element distinction  
131 versus equivalence, birds were exposed individually to a habituation-discrimination paradigm  
132 (Fig 1C). If two elements (e.g. F<sub>1</sub> & F<sub>2</sub>) represent perceptibly distinct sounds, we would expect  
133 that, after habituating subjects to a series of repetitions of one element (e.g. F<sub>1</sub>), switching to the  
134 other element (e.g. F<sub>2</sub>) would result in a renewed response, measured by investigating changes in  
135 the time subjects spent looking into the direction from which the sounds were broadcast – as is  
136 customary in habituation-discrimination approaches [20-23]. On the other hand, a lack of  
137 response renewal following the habituation sequence would indicate that the contrasted elements  
138 are not discriminated and therefore are perceptibly equivalent sounds. Further, to test whether the  
139 five elements constituting flight and prompt calls carry contextually relevant meaning, we  
140 analyzed functionally relevant behavioral responses, including vocal responses, during the initial  
141 habituation phase of each playback. If elements carry relevant meaning, playbacks of flight call  
142 elements would be expected to result in babblers looking outside the aviary more and/or moving  
143 around the aviary more (see above [14]), whilst for prompt call elements we would expect an  
144 increase in time spent looking at the nest provided (see above [14]).

145

## 146 **Results**

147 *(a) Are calls built from perceptibly distinct sounds?*

148 We first tested whether flight and prompt calls are each comprised of distinct sounds by playing  
149 back habituation-discrimination sequences of F<sub>1</sub>-F<sub>2</sub> elements from flight calls, and P<sub>1</sub>-P<sub>2</sub>, P<sub>2</sub>-P<sub>3</sub>  
150 and P<sub>1</sub>-P<sub>3</sub> elements from prompt calls to up to 12 birds individually (see Methods). In this  
151 experiment, habituation-discrimination sequences were played in natural order to avoid

152 expectancy violation (i.e. discrimination performance being inflated through playing back  
153 elements in an unnatural order). Receivers habituated to habituation sequences (each composed  
154 of 20 element repetitions played back at three-second time intervals): subjects spent a median of  
155 19% ( $IQR = 12,29$ ) of their time looking at the speakers during playbacks of the first two  
156 elements in habituation sequences but only 1% ( $IQR = 0,6$ ) of their time doing so during the last  
157 two elements of habituation sequences. One-sample Wilcoxon-tests were then used to investigate  
158 whether any changes in the proportion of time birds spent looking at the loudspeaker during the  
159 end of the habituation phase (last two habituation elements) and beginning of the discrimination  
160 phase (first two discrimination elements) significantly deviated from zero. Values significantly  
161 greater than zero indicate that habituation and discrimination elements were perceptibly distinct,  
162 while values not significantly different from zero indicate elements were not discriminated (i.e.  
163 perceived as equivalent sounds).

164 For the two flight call elements, the proportion of time receivers looked at the speaker  
165 increased 6-fold during the discrimination phase, indicating that birds discriminated  $F_2$  from  $F_1$   
166 ( $V = 36$ ,  $P = 0.008$ ,  $N = 11$ ; Fig 2A). As a consequence, we can conclude that the two elements in  
167 bi-element flight calls are perceptibly distinct (i.e.  $F_1 \neq F_2$ ). By contrast, tri-element prompt calls  
168 do not appear to be composed of three distinct elements. Within prompt calls, significant 2 to 4-  
169 fold increases in the time spent looking at the speaker during the discrimination phase were found  
170 when  $P_2$  followed  $P_1$  ( $V = 28$ ,  $P = 0.016$ ,  $N = 9$ ; Fig 2A) and when  $P_3$  followed  $P_2$  ( $V = 55$ ,  
171  $P = 0.002$ ,  $N = 10$ ; Fig 2A). However, there was no significant change in the proportion of time  
172 spent looking at the speaker between the end of the habituation phase and the discrimination  
173 phase when  $P_3$  followed  $P_1$  ( $V = 11$ ,  $P = 0.69$ ,  $N = 10$ ; Fig 2A). These results suggest that the first  
174 and third prompt call elements are perceptibly equivalent, and that both are distinct from the  
175 second prompt call element.



176 To confirm the precise make-up of prompt calls, we conducted two further analyses. First,  
177 a Friedman test confirmed that there was a significant difference between the extent to which  
178 birds discriminated the three contrasted elements in prompt calls ( $\chi^2_2 = 10.6$ ,  $P = 0.005$ ,  $N = 7$ ).  
179 Second, post-hoc two-sample Wilcoxon tests were used to compare the differences in the changes  
180 in the proportion of time birds spent looking at the speaker during the last two habituation stimuli  
181 versus the first two discrimination stimuli across each of the three sets of contrasted elements.  
182 These analyses confirmed: (a) that birds did not significantly differ in the extent to which they  
183 distinguished  $P_1$  from  $P_2$  versus  $P_2$  from  $P_3$  ( $V = 10$ , adjusted  $P = 0.16$ ,  $N = 9$ ;  $P$  value adjusted for  
184 multiple post-hoc testing; Fig 2A); but (b) that responses to  $P_2$  following  $P_1$  and to  $P_3$  following  
185  $P_2$  were both greater than responses to  $P_3$  following  $P_1$  ( $P_1$ - $P_2$  vs.  $P_1$ - $P_3$ :  $V = 28$ , adjusted  
186  $P = 0.031$ ,  $N = 7$ ;  $P_2$ - $P_3$  vs.  $P_1$ - $P_3$ :  $V = 36$ , adjusted  $P = 0.023$ ,  $N = 8$ ; Fig 2A). Thus, we are  
187 confident that the tri-element prompt call is composed of two perceptibly distinct sound types,  
188 with  $P_1 = P_3$ , but  $P_1$  and  $P_3$  to an equal extent  $\neq P_2$ .

189

190 *(b) Are perceptibly equivalent sounds shared across calls?*

191 Critical to elucidating whether multi-element calls ostensibly comprise building-blocks is to test  
192 whether elements are shared across functionally distinct calls. To investigate whether this is the  
193 case for flight and prompt calls, a different set of up to 13 birds received habituation-  
194 discrimination sequences comprising combinations of the two flight and three prompt call  
195 elements (see Methods). These were  $F_1$  and  $P_2$ ,  $F_2$  and  $P_1$ ,  $F_2$  and  $P_3$ ,  $P_{1/3}$  and  $F_1$  – with the  
196 elements used as habituation and discrimination stimuli, in this case, alternated because we  
197 wished to ensure that any expectancy violation was comparable across contrasts. Again, evidence  
198 for habituation during habituation phases was shown, with birds decreasing the percentage of  
199 time spent looking at the loudspeaker from a median of 17% ( $IQR = 10,30$ ) to a median of 3%

200 ( $IQR = 0,8$ ) between the beginning and end of the habituation sequences.

201           Subsequent one-sample Wilcoxon-tests, comparing the change in the proportion of time  
202 looking at the speaker between the last two elements of habituation phases and the first two  
203 elements of discrimination phases against a null expectation of zero, revealed that the two distinct  
204 flight call elements were each perceptually equivalent to at least one of the prompt call elements.  
205 In three of the four comparisons, the proportion of time spent looking at the loudspeaker did not  
206 significantly increase between the last two stimuli of the habituation phase and the discrimination  
207 phase. Specifically, we found  $F_1$  to be perceptually equivalent to  $P_2$  ( $V = 18, P = 0.58, N = 12$ ;  
208 Fig 2B), and  $F_2$  to be perceptibly equivalent to both  $P_1$  ( $V = 2, P = 0.19, N = 10$ ; Fig 2B) and  $P_3$   
209 ( $V = 27, P = 0.65, N = 9$ ; Fig 2B). In contrast, the proportion of time birds spent looking at the  
210 loudspeaker increased by 4-fold when the prompt call element  $P_1$  or  $P_3$  (which are equivalent, see  
211 above) was contrasted with the flight call element  $F_1$ ; meaning that  $P_1/P_3$  are distinct from  $F_1$   
212 ( $V = 55, P = 0.002, N = 11$ ; Fig 2B). Thus, these results indicate that bi-element flight calls and  
213 tri-element prompt calls both consist of the same two sound types: the first flight and second  
214 prompt call elements are perceptibly equivalent (i.e.  $F_1 = P_2$ ), as are the second flight and both  
215 first and third prompt call elements (i.e.  $F_2 = P_1 = P_3$ ). In other words, flight and prompt calls  
216 comprise the same two building blocks in different combinations.

217

218 *(c) Do sound elements carry contextual meaning?*

219 In human languages, meaningful signals are built from recombinations of meaningless sounds.  
220 To test whether or not the constituent elements of flight and prompt calls carry context-specific  
221 meaning, we measured the vocal responses and activity budgets of birds during the first two  
222 habituation stimuli of each playback (i.e. H-start, Fig 1C). First, we found no evidence to suggest  
223 that playbacks induce birds to respond with either flight or prompt calls: the median number of

224 each call given during the 6 s period of the 82 playbacks included, was zero ( $IQR = 0,0$ ). Second,  
225 we found no evidence to suggest that birds modify key behaviors in response to the playbacks.  
226 For example, we have previously shown that playbacks of flight calls on lone individuals in the  
227 aviary environment cause individuals to move around the aviary and to look outside more, while  
228 prompt call playbacks cause birds to look more at a nest in an upper corner of the aviary [14].  
229 Here, by contrast, individuals spent little time engaging in behaviors of relevance during the 6 s  
230 of each playback analyzed, spending on average: 1.3 s ( $SD = 1.1$ ) of their time in-movement; 1.3  
231 s ( $SD = 1.2$ ) looking outside the aviary; and 0.07 s ( $SD = 0.3$ ) of their time looking at the nest. In  
232 addition, the amount of time individuals spent engaged in each of these behaviors was  
233 independent of the precise element played ( $F_1, F_2, P_1, P_2, P_3$ ) (Linear Mixed Model: behavior \*  
234 element interaction,  $\chi^2 = 9.48, DF = 8, P = 0.30$ ; Fig 3A) as well as whether or not the elements  
235 played were from a flight call (F elements) or a prompt call (P elements) (LMM: behavior \*  
236 element interaction  $\chi^2 = 1.93, DF = 2, P = 0.38$ ; Fig 3B). Thus, babblers do not seem to extract  
237 contextually meaningful information from the sound elements of the two calls when played back  
238 in isolation.

239

## 240 **Discussion**

241 Using a novel application of the established habituation-discrimination paradigm, we here  
242 demonstrate that a pair of functionally distinct, multi-element calls produced by chestnut-  
243 crowned babblers are composed of two perceptibly distinct, contextually meaningless sounds,  
244 which are shared across the two vocalizations. Specifically, we show that the first element from  
245 bi-element flight calls is distinct from its second element but equivalent to the second element  
246 from tri-element prompt calls. Further, the second flight call element is equivalent to the first and  
247 third prompt call elements. In addition, none of the individual elements that make up these two

248 calls elicits differential vocal or behavioral responses of relevance in receivers. For example,  
249 subjects rarely responded to playbacks with flight or prompt calls, with a total of just nine such  
250 calls recorded across the 82 x 6 s playbacks. Moreover, babblers spent little time engaged in  
251 behaviors of relevance and the amount of time they did so was not modified by the element  
252 played; which would otherwise be expected if the elements encoded flight or prompt call-related  
253 information [14]. Together, these results suggest for the first time, that a non-human animal uses  
254 meaningless (shared) building blocks in different arrangements to encode distinct meaning.

255         A core feature of human language is that perceptibly discrete, meaningless sounds are  
256 combined in various ways to generate distinct meaning. Testing whether animals use this basic  
257 process has been hampered by a focus on minimal pairs as a way to decompose the sound system  
258 of a language - that is, identifying building-blocks through a sound's role in differentiating  
259 meaning [9, 17] . This approach necessarily requires sounds to occur across a sufficient number  
260 of vocalizations to permit meaningful comparisons, which is problematic for largely non-  
261 productive communication systems such as those utilized by animals. We demonstrate here that  
262 one can identify elements that, in essence, function like building blocks, by rather focusing on the  
263 individual perceptibility of sounds used within and across functionally distinct animal calls. We  
264 suggest that this novel approach opens up new opportunities to investigate any parallels between  
265 animal vocalizations and combinatoriality in human language.

266         We caution, of course, that any parallel between the combinatorial constructs of animal  
267 communication and word generation in human language must be tempered. First, in contrast to  
268 the combinatorial structures found in animal communication systems, combinatoriality in human  
269 language is hypothetically open-ended, with finite numbers of phonemes used in myriad  
270 combinations to generate potentially limitless information. Second, while we have shown  
271 previously that at least one element ( $P_1$ ) appears to be meaning-contrasting [14] and we have

272 shown here that elements across babbler calls (including P<sub>1</sub>) can function like building blocks,  
273 confirming that shared elements are meaning-differentiating will always be challenging in  
274 animals. To mitigate this problem, investigations into whether or not animals use building blocks  
275 in their communication systems should limit their comparisons to functionally distinct calls. This  
276 will ensure that constituent elements that are shared also play a potential role in generating  
277 meaning. Third, the building blocks of babbler calls are separated by silence, whereas in human  
278 language, they are not. Whether this is a significant distinction or a likely precursor is yet to be  
279 determined.

280         The acknowledged distinctions between babbler and human combinatoriality  
281 notwithstanding, the complexities of human language likely evolved from more rudimentary  
282 beginnings. Indeed, recent theoretical work suggests that language's productive combinatorial  
283 system was preceded by a superficial stage where the sound elements of signals overlap in their  
284 acoustic and perceptual space, but neither needed to be recognized as recombinatorial units nor  
285 utilized in a productive way by the system's users [13, 18, 24]. Subsequently, once signalers  
286 became aware of their recombinatorial system (i.e. recognize signals as being composed of  
287 smaller building blocks), they could evolve strategies (e.g. learning mechanisms) to exploit the  
288 combinatorial mechanism productively [13, 18, 24]. We propose that our study provides evidence  
289 for such a superficial vocal system by demonstrating bounded, unproductive combinatoriality  
290 (i.e. two sounds build only two signals) in babbler vocalizations. Although simple in its structure,  
291 this data supports recent hypotheses on human combinatorial systems transitioning from a more  
292 rudimentary evolutionary stage (i.e. 'superficial' combinatorial layer) before it fledged into a  
293 fully productive combinatorial system. Further experiments are now needed to clarify whether  
294 similar more superficial combinatorial structures exist in the communication systems of other  
295 species and the precise forms they take.

296 To conclude, our work provides new insights into the potential similarities between  
297 animal communication systems and the combinatorial structures of human language, with  
298 chestnut-crowned babblers reusing perceptibly distinct elements that are meaningless in isolation,  
299 but when used in different arrangements generate distinct meaning. Our study has at least three  
300 important implications. First, although we provide novel evidence for ‘superficial’  
301 combinatoriality in non-human animals, we deem it highly improbable that chestnut-crowned  
302 babblers are unique amongst animals in their ability to recombine perceptibly distinct and  
303 equivalent sounds to generate context-specific calls. Indeed, we are confident that by shifting the  
304 empirical focus to an approach that allows combinatorial systems in animals to be more easily  
305 identified, additional data in other species will undoubtedly accumulate. Second, whilst species  
306 with clearly identifiable internally structured calls, as is the case with chestnut-crowned babblers,  
307 represent intuitively more straight-forward test systems, we advocate a more general search for  
308 analogues incorporating vocalizations without clear temporal separation as happens to be the case  
309 in human language [10]. Either way, further cases are required to provide a coherent  
310 understanding of the form of early combinatorial systems, as well as their eco-evolutionary  
311 correlates. Finally, using the approach outlined, we believe that comparative work on  
312 combinatorial communication in animals will become a significant compliment to game-theoretic  
313 modelling [13, 28]; multi-agent simulations [24]; emerging sign language [29]; and  
314 communication game work [30] that aim to unpack the evolutionary origins and forms of  
315 combinatorial structures and capacities in humans and other animals.

316

## 317 **Material and methods**

### 318 **Study species and housing**

319 The study was conducted from July to September 2017 on 25 individuals from 13 different

320 groups of a free-living, color-ringed population of chestnut-crowned babblers, at the Fowlers Gap  
321 Arid Zone Research Station in New South Wales (141°42'E, 31°06'S; for details on the study  
322 population and habitat see [25]). Chestnut-crowned babblers are 50 g, group-living, cooperatively  
323 breeding passerine birds endemic to inland south-eastern Australia [25], with a known vocal  
324 repertoire of at least 18 functionally distinct calls [26]. For experimental procedures, birds were  
325 captured and housed in standardized aviaries, and were released back into their original groups  
326 after a maximum time of 48 hours (for details on capturing and aviary set-up see [14]). We have  
327 confirmed previously that birds are accepted back into their groups without retribution following  
328 their temporary absence [31], and in this study measurements of mass following their period in  
329 the aviary indicated that birds gained an average of 0.1 g (SD = 2.0) in the aviary. Birds for  
330 testing were selected randomly with respect to age and sex, although we never removed the  
331 group's breeding female or individuals with any juvenile plumage (indicating all removed  
332 individuals were nutritionally independent and > six months old).

333         During and between tests, single birds were kept in one of six compartments of a larger  
334 aviary (dimensions of each compartment: 2 x 2 x 2.5 m). Each compartment consisted of a  
335 babbler nest, perches and natural substrate. The back side of the aviary comprised a metal-mesh  
336 of 1 cm<sup>2</sup> allowing the birds a view to the outside, while the sides were opaque metal and the front  
337 consisted of one-way Perspex. During daylight, birds were fed 20 mealworms every two to three  
338 hours, and water was provided throughout (see also [14] for details on housing conditions). If two  
339 birds were removed from a group at the same time, birds were kept in different compartments,  
340 but joined into one compartment overnight. During playback experiments, only one test subject  
341 remained in the aviary, while any other birds were removed to an accommodation block out of  
342 earshot, to prevent interference with the playback.

343

344 **Playback stimuli and procedure**

345 Flight and prompt calls used for the creation of playback sequences were recorded using Electret  
346 EM-400 condenser tie-clip microphones in combination with a Sony IC-UX533 recorder  
347 (sampling frequency 44.1 kHz, 24-bit accuracy). Only high-quality vocalizations were chosen,  
348 and flight and prompt call elements were extracted and normalized using Adobe Audition CC  
349 2015. Each playback sequence consisted of 20 habituation stimuli (of one element type) and two  
350 subsequent discrimination stimuli (of another element type) broadcast at three-second intervals  
351 (Fig 1C). All test subjects were only ever exposed to stimuli originating from unfamiliar  
352 individuals. Additionally, to account for pseudo-replication and inevitable among-individual  
353 variation in element characteristics owing to, for example body size, the 20 elements used in each  
354 habituation sequence always originated from at least eight different individuals (average = 12),  
355 while the two discrimination stimuli within a sequence always originated from different  
356 individuals. Flight and prompt calls are often given by different individuals in quick succession,  
357 so babblers are accustomed to hearing flight and prompt call elements from different individuals  
358 in the field. Finally, the 20 elements within the habituation sequences and the two element within  
359 the discrimination sequences were randomly ordered, and each playback sequence/track was only  
360 used once, resulting in each test subject receiving unique playback sequences.

361         Each bird was exposed to 4 unique habituation-discrimination sequences with a break of  
362 at least 10 minutes between treatments, leading to a maximum of 100 trials across the 25 birds  
363 (but see below). Ten minutes was decided as a minimum because we wished to minimize the  
364 amount of time that any co-inhabitant of the aviary was removed for during the playback (with a  
365 minimum of 10 mins between treatments, this could be reduced to ca. 40 mins) and pilot work  
366 suggested that 10 min intervals did not confound habituation effects. In line with this pilot work,  
367 we found here that the change in looking response between H-end and H-start was equivalent for



368 the first and last habitation trials both in the within-call element comparisons (paired, two-sample  
369 Wilcoxon test:  $V = 32$ ,  $P = 0.62$ ,  $N = 12$  individuals) and among-call element comparisons  
370 ( $V = 42$ ,  $P = 0.85$ ,  $N = 12$  individuals). Playbacks were broadcast with a natural flight and prompt  
371 call amplitude of 50 dB at two meters (measured with a Castle GA206 sound level meter, C-  
372 weighted) and using a Braven BRV-X loudspeaker. The loudspeaker was placed outside 1 m  
373 away and 1 m shifted towards the side of the open, mesh-enclosed part of the aviary  
374 compartment, and was concealed by vegetation. This position was chosen because it facilitated  
375 our judgment of gaze direction towards the speaker, which is the key data of interest resulting  
376 from habituation-discrimination experiments [20-23]. In order to assess the time subjects looked  
377 into the direction of the loudspeaker (and engaged in other relevant behaviors), playbacks were  
378 video-taped using a Sony HDR-CX240.

379

#### 380 **Video coding and trial inclusion criteria**

381 Videos were analyzed frame-by-frame and blindly with respect to playback type using Adobe  
382 Audition CC 2015, with the following data extracted from each subject: number of flight and  
383 prompt calls given; number of hops/flights; and the amount of time spent looking outside, at the  
384 nest in the upper corner and at the loudspeaker. Vocalizations, movement and looking outside  
385 were easily coded, but quantifying gaze direction towards specific objects is more challenging  
386 because birds have relatively laterally-set eyes compared with humans. Nevertheless, all birds  
387 have binocular overlap in their vision to allow them to avoid obstacles during flight, interact with  
388 conspecifics, obtain food and pinpoint predators. For passerines, binocular overlaps range from  
389  $35-51^\circ$  ( $N=13$  species, including 6 non-tool-using corvids) {Troscianko, 2012 #615}. Given that  
390 babblers are passerines in the same super family as corvids (Corvoidea), suggests that they will  
391 have binocular overlap of at least  $30^\circ$  and probably closer to the  $40^\circ$  characteristic of corvids.

392 Further, for one such corvid, the common raven (*Corvus corax*, binocular overlap = 43°) looking  
393 direction towards specific objects during habituation-discrimination experiments has been  
394 assessed previously using bill orientation {Reber, 2016 #614}. In line with previous work, we  
395 here measure looking at the speaker or the nest by assessing the orientation of the test bird's bill  
396 which had to directly point towards the object in question ( $\pm 30^\circ$ , well within the expected field  
397 of binocular overlap). Babblers routinely turn their head in order to pinpoint food, conspecifics  
398 and predators, and we have substantial experience with gaze direction for each of these stimuli in  
399 the aviary setting. Through double-blind scoring of time spent looking at the speaker during the  
400 end of habituation (H-end) and discrimination phases of 41 trials (50% of the 82 included), we  
401 found substantial inter-scorer agreement (Interclass Correlation Coefficient for two-way model  
402 based on absolute agreement and single rater scores ICC = 0.83,  $P < 0.001$ , 95% CI = 0.75-0.89)  
403 {Hallgren, 2012 #616}.

404 Out of the 100 potential trials, 82 were included in the analyses. Two trials were not  
405 obtained because we released a bird early due to concerns over a loss of appetite and failed to  
406 capture H-start of another trial in the camera. Further, in 5 trials, birds failed to look in the  
407 direction of the speaker during the habituation phase, a prerequisite of the habituation-  
408 discrimination paradigm, and likewise, a further 11 had to be excluded as they looked at the  
409 speaker at least as much during H-end as H-start. There was no systematic bias in the habituation  
410 stimuli that were excluded, with each of the 5 habituation elements being removed at least twice.

411

## 412 **Statistical analyses**

### 413 *Element discrimination*

414 Testing whether elements are perceived as dissimilar or equivalent was primarily investigated  
415 using a series of one-sample Wilcoxon tests. Specifically, the change in the proportion of time

416 individuals spent looking at the speaker between the discrimination phase (D) and the end of the  
417 habituation phase (H-end) was contrasted against a null expectation of zero change (Figs. 1C, 2).  
418 The only exception was to further clarify the form of prompt calls. In this case, we additionally  
419 used Friedman combined with post-hoc two-sample Wilcoxon tests to test the *differences* in the  
420 changes of responses between H-end and D for contrasted pairs of elements (i.e. P<sub>1</sub>-P<sub>2</sub> vs. P<sub>2</sub>-P<sub>3</sub>  
421 vs. P<sub>1</sub>-P<sub>3</sub>) - post-hoc *P*-values were adjusted using the Bonferroni-holm method [32]. For all  
422 analyses of element discrimination, we used the proportion of time looking at the speaker (rather  
423 than absolute time) since the birds were not always in camera view for the entire 6 s H-end and D  
424 phases (H-end: mean time in view = 5.9; SD = 0.2, range = 4.8-6.2; D: mean = 6.0, SD = 0.1  
425 range = 5.3-6.4). All statistical analyses were conducted in R (version 3.4.2) - Wilcoxon tests  
426 using the “exactRankTest”-package [33], and Friedman tests using the “stats”-package [34].

427

#### 428 *Element meaning*

429 To investigate whether the five constituent elements of flight and prompt call elements carry  
430 contextual meaning, we performed two Linear Mixed effects Models (LMM). In both models, the  
431 response term was the amount of time (during the 6 s of H-start for each element, square-root  
432 transformed) that individuals were observed: looking outside (not at the speaker); looking at the  
433 nest in an upper corner of the aviary; and in-movement (mainly hopping among perches). These  
434 behaviors were chosen because we have previously shown in the same aviary set-up that babblers  
435 change the duration of each behavior in response to playbacks of flight and prompt calls [14]. It  
436 is important to note that the sum percentage of time that individuals engaged in these 3 behaviors  
437 averaged just 44%, meaning that individuals could respond to each behavior independently. The  
438 term of interest in the first model was the interaction between element type (F<sub>1</sub>, F<sub>2</sub>, P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>) and  
439 behavioral response (in-movement, looking-out, looking-nest); while in the second model, we

440 interacted whether or not the element in question was from a flight call (F elements) or a prompt  
441 call (P elements) with behavioral response. In both models, time in view was fitted as a covariate  
442 and trial identity nested within individual identity were fitted as random intercepts to account for  
443 the fact that trials had 3 behavioral responses and that multiple elements were played to the same  
444 individual. Model reduction were not performed for either model as in both cases the key result is  
445 the interaction between element and behavior. The above two models were fitted in R using the  
446 “lme4” package, and the full model with and without the interaction of interest were compared  
447 using log-likelihood ratio tests to determine the significance of the interaction term [34, 35].

448  
449 **Acknowledgements:** We thank Simon Griffith, Keith Leggett and the Dowling family for  
450 logistical support at Fowlers Gap; Kiara L’Herpinier, Joseph England and Jennifer Page for help  
451 with fieldwork; Steven Moran, Volker Dellwo and Stuart Watson for discussions, and three  
452 anonymous reviewers for their constructive feedback. The research was approved by the ethics  
453 committee of the University of Exeter (Application number 2018/2301).

454 **Data accessibility:** All data to reproduce the work is provided as supplementary material.

455

## 456 **References**

- 457 1. Hockett CF (1960) The Origin of Speech. *Sci Am* 203:88-111.
- 458 2. Hauser MD, Chomsky N, Fitch WT (2002) The Faculty of Language: What Is It, Who Has  
459 It, and How Did It Evolve? *Science* 298(5598):1569-1579.
- 460 3. Hurford J (2007) *The origins of meaning* (Oxford University Press, Oxford).
- 461 4. Arnold K, Zuberbühler K (2006) Language evolution: Semantic combinations in primate  
462 calls. *Nature* 441(7091):303.
- 463 5. Zuberbühler K (2018) Combinatorial capacities in primates. *Curr Opin Behav Sci* 21:164-  
464 169.

- 465 6. Ouattara K, Lemasson A, Zuberbühler K (2009) Campbell's monkeys concatenate  
466 vocalizations into context-specific call sequences. *Proc Natl Acad Sci USA* 106(51):22026-  
467 22031.
- 468 7. Engesser S, Ridley AR, Townsend SW (2016) Meaningful call combinations and  
469 compositional processing in the southern pied babbler. *Proc Natl Acad Sci USA*  
470 113(21):5976-5981.
- 471 8. Suzuki TN, Wheatcroft D, Griesser M (2016) Experimental evidence for compositional  
472 syntax in bird calls. *Nat Commun* 7:10986.
- 473 9. Yip MJ (2006) The search for phonology in other species. *Trends Cogn Sci* 10(10):442-446.
- 474 10. Bowling DL, Fitch WT (2015) Do Animal Communication Systems Have Phonemes?  
475 *Trends Cogn Sci* 19(10):555-557.
- 476 11. Berwick RC, Okanoya K, Beckers GJL, Bolhuis JJ (2011) Songs to syntax: the linguistics of  
477 birdsong. *Trends Cogn Sci* 15(3):113-121.
- 478 12. Engesser S, Townsend SW (2019) Combinatoricity in the vocal systems of non-human  
479 animals. *WIREs Cogn Sci* e1493.
- 480 13. Zuidema W, de Boer B (2009) The evolution of combinatorial phonology. *J Phon* 37(2):125-  
481 144.
- 482 14. Engesser S, Crane JM, Savage JL, Russell AF, Townsend SW (2015) Experimental Evidence  
483 for Phonemic Contrasts in a Nonhuman Vocal System. *PLoS Biol* 13(6):e1002171.
- 484 15. Hailman JP, Ficken MS, Ficken RW (1985) The Chick-a-Dee Calls of *Parus atricapillus* - a  
485 Recombinant System of Animal Communication Compared with Written-English. *Semiotica*  
486 56(3-4):191-224.
- 487 16. Suzuki TN (2013) Communication about predator type by a bird using discrete, graded and  
488 combinatorial variation in alarm calls. *Anim Behav* 87:59-65.
- 489 17. Chomsky N, Halle M (1968) *The Sound Pattern of English* (Harper & Row, New York).
- 490 18. Zuidema W, de Boer B (2018) The evolution of combinatorial structure in language. *Curr*  
491 *Opin Behav Sci* 21:138-144.
- 492 19. Eimas PD, Siqueland ER, Jusczyk P, Vigorito J (1971) Speech Perception in Infants. *Science*  
493 171(3968):303-306.
- 494 20. Charlton BD, Ellis WA, McKinnon AJ, Brumm J, Nilsson K, Fitch WT (2011) Perception of  
495 male caller identity in Koalas (*Phascolarctos cinereus*): acoustic analysis and playback  
496 experiments. *PLoS ONE* 6(5):e20329.
- 497 21. Cheney DL, Seyfarth RM (1988) Assessment of meaning and the detection of unreliable  
498 signals by vervet monkeys. *Anim Behav* 36(2):447-486.

- 499 22. Fitch WT (2006) Rhesus macaques spontaneously perceive formants in conspecific  
500 vocalizations. *J Acoust Soc Am* 120(4):2132-2141.
- 501 23. Reby D, Hewison M, Izquierdo M, Pépin D (2008) Red Deer (*Cervus elaphus*) Hinds  
502 Discriminate Between the Roars of Their Current Harem-Holder Stag and Those of  
503 Neighbouring Stags. *Ethology* 107:954-959.
- 504 24. de Boer B, Zuidema W (2010) Multi-Agent Simulations of the Evolution of Combinatorial  
505 Phonology. *Adapt Behav* 18(2):141-154.
- 506 25. Russell AF (2016) Chestnut-crowned babblers: Dealing with climatic adversity and  
507 uncertainty in the Australian arid zone. *Cooperative breeding in vertebrates: studies in*  
508 *ecology, evolution and behavior*, eds Koenig WD, Dickinson JL (Cambridge University  
509 Press, Cambridge, MA), pp 150-164.
- 510 26. Crane JMS, Savage JL, Russell AF (2016) Diversity and function of vocalisations in the  
511 cooperatively breeding Chestnut-crowned Babbler. *Emu* 116(3):241.
- 512 27. Young CM, Browning LE, Savage JL, Griffith SC, Russell AF (2013) No evidence for  
513 deception over allocation to brood care in a cooperative bird. *Behav Ecol* 24(1):70-81.
- 514 28. Nowak MA, Krakauer DC, Dress A (1999) An error limit for the evolution of language. *Proc*  
515 *R Soc B* 266(1433):2131-2136.
- 516 29. Sandler W, Aronoff M, Meir I, Padden C (2011) The gradual emergence of phonological  
517 form in a new language. *Nat Lang Linguist Th* 29(2):503-543.
- 518 30. Verhoef T, Kirby S, de Boer B (2014) Emergence of combinatorial structure and economy  
519 through iterated learning with continuous acoustic signals. *J Phon* 43:57-68.
- 520 31. Nomano FY, Browning LE, Savage JL, Rollins LA, Griffith SC, Russell AF (2015)  
521 Unrelated helpers neither signal contributions nor suffer retribution in chestnut-crowed  
522 babblers. *Behav Ecol* 26(4):986-995.
- 523 32. Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and  
524 powerful approach to multiple testing. *J R Statist Soc B* 57:289-300.
- 525 33. Hothorn T, Hornik K (2017) exactRankTests: Exact Distributions for Rank and Permutation  
526 Tests. R package version 0.8-29. Accessed 01 June 2017.
- 527 34. R-Core-Team (2014) R: A language and environment for statistical computing. R  
528 Foundation for Statistical Computing. Vienna, Austria. Accessed 01 June 2017.
- 529 35. Bates D, Maechler M, Bolker B, Walker S (2014) Fitting linear mixed-effects models using  
530 lme4. *J Stat Softw* 67(1):1-48.  
531

532 **Figure legends**

533 **Fig 1. Study design.** (A) Chestnut-crowned babblers (credit AF Russell). (B) Spectrogram of a  
534 flight and a prompt call, with the flight call being composed of F<sub>1</sub>F<sub>2</sub> elements and prompt calls of  
535 P<sub>1</sub>P<sub>2</sub>P<sub>3</sub> elements. (C) Schematic overview of the habituation-discrimination experiment. During  
536 the habituation phase subjects were accustomed to one element type (from at least 8 different  
537 unfamiliar individuals) constituting the habituation stimuli (H<sub>1</sub> – H<sub>20</sub>, e.g. F<sub>1</sub>), which was  
538 repeated 20 times at three-second intervals. Subsequently, two repetitions of another element type  
539 (both from different unfamiliar individuals) constituting the discrimination stimuli (D<sub>1</sub> – D<sub>2</sub>, e.g.  
540 F<sub>2</sub>) were broadcast. To assess the discrimination between contrasted elements, the change  
541 between the proportion of time subjects looked toward the loudspeaker during the first two  
542 discrimination (D) and the last two habituation stimuli (H-end) was analyzed.

543

544 **Fig 2. Element discriminations.** Results of the habituation-discrimination experiments when  
545 contrasting flight and prompt call elements: (A) within flight or prompt calls; and (B) between  
546 flight and prompt calls. Figures show the changes in the proportion of time subjects looked at the  
547 loudspeaker during the discrimination phase (D) and the end of the habituation phase (H-end) for  
548 each element comparison. The vertical (red) line represents the null expectation of no-change.  
549 Boxes represent the 25%, 50% and 75% quartiles of the raw data, whiskers extend to 1.5 x inter-  
550 quartile ranges, while dots show outliers. Significant changes in the proportion of time spent  
551 looking at the loudspeaker between H-end and D are shown with asterisks (\* p < 0.05, \*\*  
552 p < 0.01). In Figure A elements were presented in natural order (as shown), while in B element  
553 orders were randomized since no natural order exists in between-call comparisons (‡ denotes that  
554 P<sub>1</sub> was alternated with the equivalent sound P<sub>3</sub>).

555

556 **Fig 3. Element meaning.** The amount of time individuals spent engaged in behaviors of  
557 relevance during H-start when: (A) behavioral responses were considered for each of the 5  
558 element types individually (F<sub>1</sub>, F<sub>2</sub>, P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>); and (B) behavioral responses were considered for  
559 flight call (F) elements versus prompt call (P) elements. Shown are the raw data with point sizes  
560 indicating the frequency of occurrence at given time values. In Figure A dot shapes (circular or  
561 triangular) illustrate the two discriminated sound types (i.e. circular F<sub>1</sub> & P<sub>2</sub>; triangular F<sub>2</sub>,  
562 P<sub>1</sub> & P<sub>3</sub>). In Figures A & B red shaded dots illustrate flight call elements and blue shaded dots  
563 prompt call elements. Note there is no obvious tendency for different elements to elicit  
564 differential behavioral responses. Analyses in each case are based on 246 behavioral responses  
565 during the 82 playbacks. In each model, the variance component of the random term ‘trial  
566 identity’ was 0, indicating that the variation in activity budgets within and among trials were  
567 equivalent. By contrast, individual identity explained a significant 15% of the residual variance in  
568 each model (variance component = 0.04, P < 0.001), indicating that some individuals were more  
569 active than others. Finally, inclusion of the interaction term of interest in each model raised the  
570 AIC by 7 points (Model 1) and 2 points (Model 2), indicating that power of the models were  
571 reduced when the interaction terms were included (see text for statistics).