

# 1 For learning morphological rules, production tasks are 2 no better than comprehension

3 Elizabeth Pankratz\*, Aislinn Keogh\*, Simon Kirby, Jennifer Culbertson

4 Centre for Language Evolution, University of Edinburgh

5 \**Joint first author*

6 e.c.pankratz@sms.ed.ac.uk

7 aislinn.keogh@ed.ac.uk

## 8 Abstract

9 This paper weaves together two strands of previous research: one which identifies that  
10 adults struggle to learn morphological rules, and another which indicates that language  
11 learning can be facilitated by language production. Here, we ask: can a production task  
12 help adults learn morphological rules? In two artificial language learning experiments,  
13 we taught participants a language that indicated thematic role with both a fixed word  
14 order (a word-level rule, which should be easier for adults to learn) and case marking  
15 (a morphological rule, which should be harder for adults to learn). We manipulated  
16 whether participants practised this artificial language using a comprehension task or a  
17 production task, and then asked whether participants who did the production task were  
18 more likely to learn the case marking rule. We also assessed how aware participants  
19 were of the morphological pattern that results from the case marking, even if they did  
20 not associate certain markers with certain thematic roles *per se*. Experiment 1 tested L1  
21 English participants, and Experiment 2 tested L1 German participants: populations that  
22 differ in their prior experience of case. In both experiments, we found that participants  
23 across the board failed to learn the case marking rule, even though the majority did de-  
24 tect the morphological pattern that was the consequence of case marking. We conclude  
25 that the production task we used in this study did not suffice to help adults learn a less  
26 accessible morphological rule.

27 *Keywords:* artificial language learning, case marking, word order, language production,  
28 segmentation

## 29 1 Introduction

30 Learners of a new language will often discover that no matter how many target-language  
31 books, films, or podcasts they absorb, their language skills do not truly blossom until they  
32 have practised producing the language themselves.

33 Language production benefits both infant learners of their first language as well as  
34 adult learners of further languages. In L1 acquisition, children who use their target lan-  
35 guage more frequently show stronger expressive abilities in that language throughout  
36 development, independent of their level of comprehension (Bohman et al., 2010; Don-  
37 nelly and Kidd, 2021; Ribot et al., 2018). And in L2 acquisition, production tasks have  
38 been shown to improve how L1 Mandarin users learn English relative clauses (Izumi,

39 2002) and how people with diverse L1s learn German grammatical gender (Keppenne  
40 et al., 2021); the way production tasks benefit L2 acquisition has been influentially re-  
41 ferred to as the Output Hypothesis (Swain, 2005). Artificial language learning studies  
42 also illustrate that production practice helps adults both to learn rules (Hopman and  
43 MacDonald, 2018) and to generalise them (Hopman, 2022).

44 A separate strand of research has shown that adult learners don't acquire all kinds of  
45 rules equally well. Particularly troublesome are morphological rules; adult learners' dif-  
46 ficulty with both nominal and verbal inflectional morphology has been well documented  
47 (see, e.g., Bentz and Winter, 2013; Holmes and Dejean De La Bâtie, 1999; Parodi et al.,  
48 2004; Rogers, 1987; Sagarra and Ellis, 2013). Case marking poses a particular challenge:  
49 L2 learners of German and Turkish struggle to learn the case-marking morphology, even  
50 if their L1 also has case (Jordens et al., 1989; Papadopoulou et al., 2011). In contrast,  
51 rules that apply to larger chunks, such as words and phrases, seem more accessible. For  
52 example, when learning noun classification systems, adults tend to rely more on class  
53 membership cues that do not require them to segment below word level (i.e., deter-  
54 miners) compared to sub-word cues that do require segmentation (i.e., suffixes; Keogh and  
55 Lupyan, in press). And typological evidence also suggests that adults may prefer word-  
56 level rules: languages with more adult L2 learners tend to be morphologically simpler  
57 (Bentz and Winter, 2013; Lupyan and Dale, 2010).

58 In this study, we ask whether a production task can cause adult learners to acquire a  
59 hard-to-learn morphological rule that requires words to be segmented, moving beyond  
60 an easier word-level rule that requires no segmentation. This question builds on intrigu-  
61 ing results from Hopman and MacDonald (2018). In their artificial language learning  
62 experiment, participants who did a production task seem to have acquired morphologi-  
63 cal rules better than word-level ones. Specifically, those participants appear to be more  
64 sensitive to errors in suffixing than errors in word order (see their Figure 5, p. 968).

65 However, this boost to morphological rule learning is just a descriptive result that  
66 the original paper does not explore further. Additionally, this finding might come not  
67 from the production task *per se*, but rather from properties of the artificial language that  
68 Hopman and MacDonald used. The language was very complex in that every sentence  
69 contained multiple modifiers and adverbial phrases, so the word order rules might have  
70 been hard to identify. On the other hand, several words in every sentence contained  
71 identical suffixes, making the morphological pattern highly salient. Here, we aim to  
72 follow up on Hopman and MacDonald's result using an artificial language designed to  
73 test whether production tasks help adults learn morphological rules over word-level  
74 rules.

75 Why would we expect language production to help adults learn morphological rules  
76 at all? One explanation for why production strengthens language learning is what Swain  
77 (2005) describes as its "noticing role". The idea is that when people produce a language,  
78 they process it more deeply and are thus more likely to notice linguistic patterns and  
79 induce possible generalisations (see also Izumi, 2002). As long as the task is more active  
80 than a recognition-based comprehension task, we would expect the noticing role of pro-  
81 duction to take effect; based on literature on the effects of different kinds of tests, any  
82 kind of test beyond passive recognition should improve learning (see, e.g., Kang et al.,  
83 2007; McDaniel et al., 2007; McDermott et al., 2014). We therefore hypothesised that a  
84 production task could draw people's attention to morphological patterns that they may

85 otherwise have missed.

86 As a testing ground for this hypothesis, we used the well-studied trade-off between  
87 case marking, an example of a morphological rule, and fixed word order, an exam-  
88 ple of a word-level rule (Bentz and Winter, 2013; Fedzechkina et al., 2011; Levshina,  
89 2020; Lupyan and Dale, 2010). The rest of this paper discusses two preregistered exper-  
90 iments ([https://osf.io/qbjda/?view\\_only=cdd6600d8e9c45e5bc153058ea97df29](https://osf.io/qbjda/?view_only=cdd6600d8e9c45e5bc153058ea97df29))  
91 that test this hypothesis on two populations that differ in their prior experience with  
92 case-marking systems: Experiment 1 tests L1 English participants, and Experiment 2  
93 tests L1 German participants.

94 To foreshadow our results: overall, participants learned the fixed word order rule but  
95 failed to acquire the case marking rule, although the majority did notice the recurring  
96 syllable pattern that was the consequence of case marking. Even participants already  
97 familiar with the concept of case (the German L1 participants in Experiment 2) showed  
98 the same clear preference to treat words as the smallest unit in the language and not to  
99 segment below this level. With respect to our main hypothesis, we found that taking  
100 part in a production task does not make participants more likely to learn the case mark-  
101 ing rule. This result suggests that, although production tasks may generally facilitate  
102 learning, they don't necessarily help adult learners to discover morphological rules.

## 103 2 Experiment 1

104 Participants were trained on a series of sentences that each described a transitive event  
105 between two human characters. These sentences were designed to be compatible with  
106 both word-level and morphological strategies for marking thematic role. Specifically,  
107 each sentence had the same fixed word order (SOV), a consistent word-level cue, *and*  
108 each noun bore a suffix corresponding to its grammatical role (nominative for the agent  
109 role and accusative for the patient role), a consistent morphological cue.

110 For example, participants might see an image of a fairy pushing a doctor and learn  
111 the corresponding sentence *fuvu zijo gix*. Then they might see a cowboy kicking a pi-  
112 rate and learn the sentence *lovu wujo kuv*. In both sentences, the word order is SOV, and  
113 in both sentences, the agent is marked with *-vu* and the patient with *-jo*. Thus partici-  
114 pants could analyse the language in two different ways: like (1), in which nouns remain  
115 unsegmented, or like (2), in which nouns are segmented into stem and case marker.

- 116 (1) a. fuvu zijo gix  
117       fairy doctor push  
118       b. lov u wujo kuv  
119       cowboy pirate kick
- 118 (2) a. fu-vu      zi-jo      gix  
119       fairy-NOM doctor-ACC push  
119       b. lo-vu       wu-jo      kuv  
119       cowboy-NOM pirate-ACC kick

120 Note that the recurring syllables at the end of each noun could also be analysed in  
121 terms of their linear order: participants could arrive at an analysis like “the first noun  
122 always ends in *vu*, and the second noun always ends in *jo*”. This is not a case marking

123 analysis *per se*, since it's not based on thematic roles. But it is still of interest to us,  
124 because we're concerned with how well participants can identify patterns below word  
125 level. For this reason, in what follows, we refer to the two possible analyses not as  
126 "fixed word order" and "case marking", but rather as "unsegmented" and "segmented",  
127 respectively.

128 A crucial aspect of the training phase's design is that participants received no di-  
129 rect evidence that nouns have morphological structure, because none of the characters  
130 appeared as both agent and patient. Thus, the language's grammar is ambiguous. To  
131 illustrate concretely: a participant would only ever see the fairy character as an agent,  
132 only ever labelled as *fuvu*. They receive no information about what form this word would  
133 take if the fairy were a patient. The word might become *fujo*, following the segmented  
134 analysis, or remain *fuvu*, following the unsegmented analysis. Thus it was possible for  
135 participants to successfully learn the training data without segmenting the words.

136 After training, we split participants into two groups to introduce the manipulation by  
137 task. Half of the participants practised the sentences they had learned using a more ac-  
138 tive production task, while the other half practised using a more passive comprehension  
139 task. The production task involved constructing sentences by clicking on the component  
140 syllables in the correct order, while the comprehension task simply involved choosing  
141 the correct image from an array of two.

142 Next, in the testing phase, we showed participants the same scenes they saw in train-  
143 ing, but with the characters' roles reversed. For example, where in training they saw a  
144 fairy pushing a doctor, now they saw the doctor pushing the fairy. We then asked them  
145 to judge two different sentences that might describe this scene. The first kind of sen-  
146 tence was formed using the unsegmented analysis: the full words for the agent and  
147 patient were rearranged. The second kind of sentence was formed using the segmented  
148 analysis: only the stems were rearranged, and the case markers stayed in place.

149 If participants learned the nouns as unsegmented, holistic chunks, they should ac-  
150 cept the first kind of sentence. If they segmented the nouns into stem and suffix, they  
151 should accept the second kind of sentence. Given previous findings that adults struggle  
152 to learn case morphology, we expected our participants to prefer sentences formed using  
153 the unsegmented analysis. However, crucially, here we test whether this preference is  
154 affected by the type of practice task they did: comprehension or production.

155 Finally, participants completed a one-shot cloze task with a novel character, i.e., a  
156 character held out from the set encountered in training. The goal here was to assess  
157 whether participants were aware of the language's morphological patterns (i.e., that the  
158 first noun always ends in a particular syllable, and that the second noun always ends in  
159 another), whether or not they actually analysed these syllables as case markers.

## 160 2.1 Materials

161 The artificial language contained transitive sentences made up of three words: one for  
162 the agent, one for the patient, and one for the action, in that order (i.e., SOV). All verbs  
163 were monosyllabic CVC nonsense words, and all nouns were disyllabic CVCV nonsense  
164 words. Verbs were randomly selected from a set of 28: *gax, gix, gox, hix, jeg, jix, juf, juz,*  
165 *kex, kez, kuv, kux, nuz, puv, pux, vaf, vof, wez, wox, zax, zok, zox, zud, zuf, zug, zup, zuv,*  
166 *and zux.* Nouns were randomly assembled from nine possible stem syllables (*bu, fu, gu,*

167 *ki, lo, ru, wu, ze, and zi*) and two suffix syllables (*vu* and *jo*) such that all agent nouns  
168 took one suffix and all patient nouns took the other.

169 Each sentence accompanied an image, a line drawing of two human characters inter-  
170 acting. A few examples are shown in Figure 1. The nine possible characters were: a chef,  
171 a cowboy, a doctor, a fairy, a footballer, a nun, a pirate, a princess, and a wizard. Each  
172 scene showed the agent character engaging in a reversible transitive action toward the  
173 patient character. The nine possible actions were: admiring, greeting, kicking, kissing,  
174 patting, poking, pushing, seeing, and yelling. Each image had two mirrored versions:  
175 one with the agent on the left, and one with the agent on the right.

176 To keep the artificial lexicon easily learnable, we randomly selected only six char-  
177 acters and two actions for each participant. The characters and actions were randomly  
178 associated with nonsense noun stems and verbs from the sets listed above. Then, each  
179 character was mapped to the thematic role they would appear in during the training  
180 phase. The mapping between characters and roles was random, with one constraint:  
181 we disallowed any permutations in which all agents were female and all patients were  
182 male (or vice versa), to forestall analyses of the suffixes as gender markers. All in all,  
183 participants saw 18 unique scenes during training: 3 agents  $\times$  3 patients  $\times$  2 verbs.

## 184 2.2 Procedure

185 We wrote the experiment in JavaScript using the jsPsych library (de Leeuw et al., 2023).  
186 It contains four phases, detailed below and illustrated in Figure 1.

### 187 2.2.1 Training

188 In each training trial, participants saw an image alone for 1000 ms. Then the correspond-  
189 ing sentence in the artificial language appeared below it. 2500 ms later, a ‘next’ button  
190 appeared below the sentence. Clicking on it advanced participants to the next trial.

191 The whole training phase consisted of three blocks of 18 trials each, one trial per  
192 scene. Participants could optionally take a short break between each block.

### 193 2.2.2 Practice

194 After training, participants were divided into two groups: one group completed a pro-  
195 duction practice task (the PRODUCTION condition), and the other completed a compre-  
196 hension practice task (the COMPREHENSION condition). Both practice tasks involved fa-  
197 miliar scenes and sentences that participants had encountered during training.

198 Participants in the PRODUCTION condition saw a familiar scene and had to build the  
199 correct sentence for this scene out of its component syllables. Below the image were five  
200 gaps, and below the gaps was one button per syllable in the sentence, shown in a random  
201 order. Although this task is less active than, say, speaking the artificial language sentence  
202 aloud, it still involves reproducing the linguistic signal that participants received. This  
203 reproduction places additional demands on participants that the comprehension task, as  
204 a simple recognition task, does not (more detail on the comprehension task below).

205 Clicking one of the buttons added that syllable into the leftmost gap in the sentence,  
206 so the sentence was filled in left to right as each syllable was clicked. An ‘undo’ button

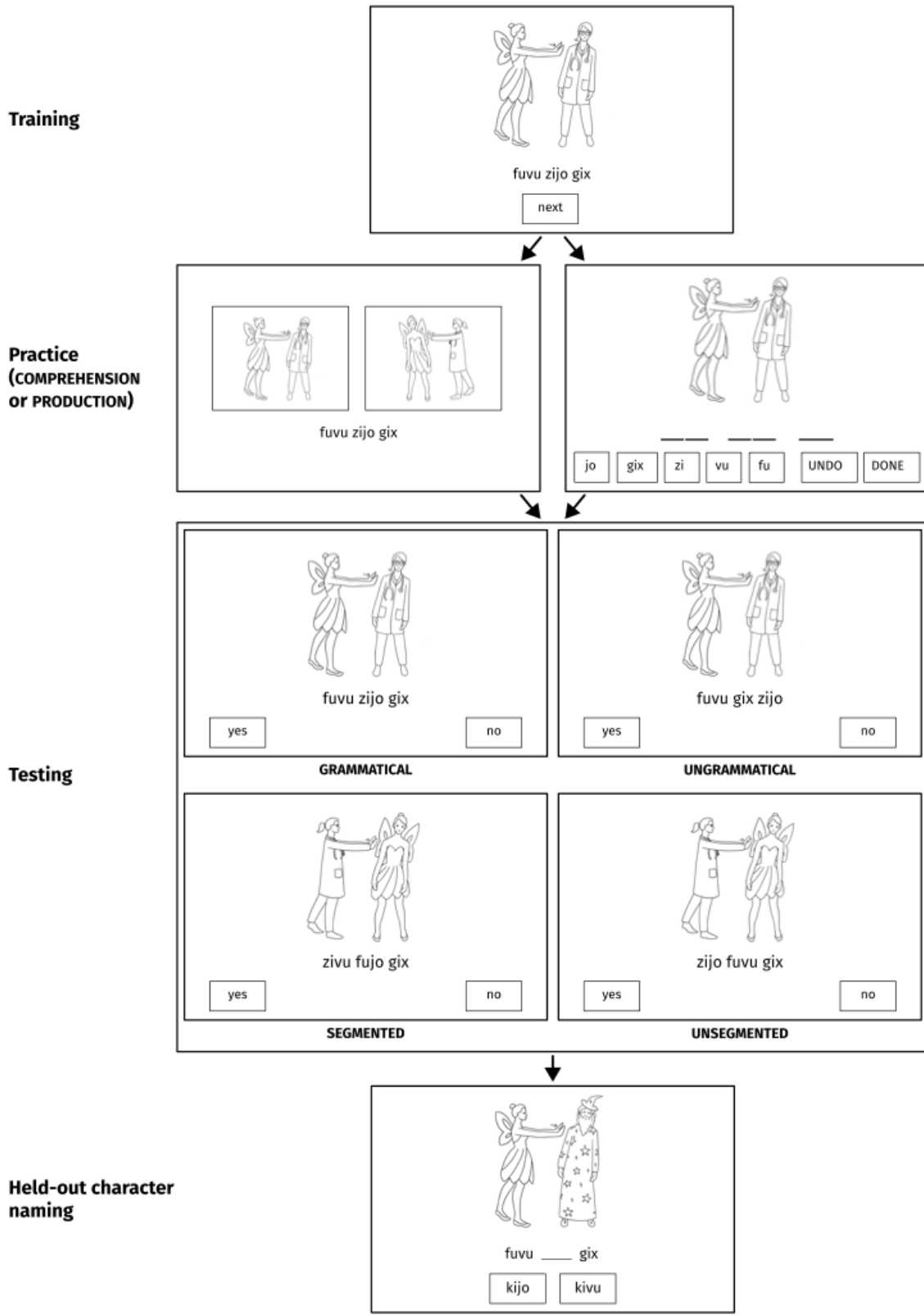


Figure 1: A schematic overview of the four phases of the experiment. All participants do the same training, then do either the comprehension or the production practice task. Then all participants complete the same testing and character naming phases. In other words, the two conditions differ *only* in the practice task.

207 emptied the most recently filled gap. Participants could submit their sentence with the  
208 ‘done’ button as long as the sentence included every syllable once.

209 After submitting the sentence, participants received feedback on their response and  
210 were shown the correct sentence. The feedback stayed on-screen until participants  
211 clicked ‘next’ to continue. Each participant did 18 production trials, one per familiar  
212 scene, shown in a random order.

213 Participants in the COMPREHENSION condition were shown a familiar sentence and  
214 had to select the corresponding scene from an array of two. The target scene was a  
215 familiar one encountered during training; the foil image contained the same characters  
216 but with the thematic roles reversed (that is, if the target showed the fairy pushing the  
217 doctor, then the foil would show the doctor pushing the fairy). The order of target and  
218 foil was randomised on each trial. The agent appeared on the left in one image and on  
219 the right in the other, so that the characters themselves remained in the same position  
220 in each image.

221 So that we do not confound our results by giving production participants a segmen-  
222 tation advantage (in that they see each syllable individually on its own button), we made  
223 the sentence in the comprehension task appear on screen one syllable at a time, with a  
224 new syllable appearing every 500 ms. Once the full sentence was visible, participants  
225 could click on one of the two scenes. They received feedback on their response which  
226 stayed on-screen until they clicked ‘next’ to move to the next trial. Each participant did  
227 18 comprehension trials, one per familiar scene, shown in a random order.

### 228 **2.2.3 Testing**

229 After the practice phase, all participants were asked to judge a number of sentences,  
230 some familiar and some novel. In each trial, participants saw a scene and a sentence,  
231 along with the prompt “Could someone who speaks this language describe this scene  
232 using the sentence below?”. We used the  $\pounds$  and  $\jmath$  keys for “yes” and “no”, with the  
233 mapping randomly determined for each participant (but kept the same for each trial).  
234 Participants received no feedback during this phase: pressing either  $\pounds$  or  $\jmath$  immediately  
235 moved them on to the next trial.

236 The testing phase contained four kinds of trials. First, there were GRAMMATICAL  
237 trials: nine of the familiar scenes and sentences from training, randomly sampled. If  
238 participants learned the language, they should always accept these sentences. Second,  
239 there were UNGRAMMATICAL trials: the other nine familiar scenes from training, but with  
240 sentences rearranged into a different word order (SVO, rather than the SOV participants  
241 were trained on). If participants learned the word order rule in the language, we reasoned  
242 that they should always reject these sentences.

243 We preregistered a particular exclusion criterion for these sentences which allowed  
244 participants to make up to and including four mistakes across these 18 GRAMMATICAL  
245 and UNGRAMMATICAL trials. In other words, the minimum accuracy permitted was 77.7%.  
246 However, it is worth noting that by excluding participants who accepted a different word  
247 order, we might be rejecting exactly those participants who had adopted a case marking  
248 analysis, since case marking languages generally permit freer word order (Fedzechkina  
249 et al., 2011; Lupyan and Dale, 2010). In an exploratory analysis reported in Appendix  
250 A, we removed the ungrammaticality criterion and re-ran the analysis we describe be-  
251 low, this time including participants who accepted any number of “ungrammatical” sen-

252 tences. The overall pattern of results remains the same regardless of whether we use  
253 this criterion. This suggests that participants who accept the “ungrammatical” SVO sen-  
254 tences do so not because they have learned a free word order along with a case marking  
255 rule, but because they haven’t learned the language reliably.

256 The final two trial types in the testing phase provide the critical data for our research  
257 question. In both trial types, the scenes contained familiar characters, but their thematic  
258 roles are reversed from the ones participants saw them in during training. Reversing  
259 the thematic roles causes the segmented analysis to yield a different sentence than the  
260 unsegmented analysis.

261 To illustrate: if a participant learned that the sentence *fuvu zijo gix* goes along with  
262 the fairy pushing the doctor, then in the testing phase, they would encounter two trials  
263 with a scene of the doctor pushing the fairy. In the SEGMENTED trial, they would see the  
264 doctor pushing the fairy along with the sentence in (3), which was formed by swapping  
265 just the CV stems. In the UNSEGMENTED trial, they would see this same scene along with  
266 the sentence in (4), which was formed by swapping the entire nouns. Participants were  
267 asked to judge novel sentences formed according to these two rules for all 18 reversed-  
268 role scenes.

269 (3) *zi-vu fu-jo gix*  
doctor-NOM fairy-ACC push

270 (4) *zijo fuvu gix*  
doctor fairy push

271 All in all, the testing phase contained 54 trials (9 GRAMMATICAL + 9 UNGRAMMATICAL  
272 + 18 SEGMENTED + 18 UNSEGMENTED). The order of these trials was randomised for each  
273 participant.

#### 274 2.2.4 Held-out character naming

275 The final phase of the experiment involved a one-shot trial in which participants saw a  
276 scene with one familiar character, one held-out character that had not been previously  
277 seen in the experiment, and a familiar action happening between them. These elements  
278 were all randomly chosen. The familiar character always appeared in the same thematic  
279 role from training, so the label for that character was also familiar. The held-out char-  
280 acter assumed the other role.

281 Along with the scene, participants saw a sentence with a gap where the word for the  
282 new character would be. They were asked “What seems like the most plausible word for  
283 the new character in this scene?”. Two alternatives were provided, formed by combining  
284 a random held-out stem with *-vu* and with *-jo*. For example, if the scene was the fairy  
285 (familiar noun) pushing the wizard (unfamiliar noun), and the sentence was *fuvu \_\_\_\_*  
286 *gix*, participants would be asked to choose between *kivu* and *kijo* as the label for the  
287 wizard.

### 288 2.3 Participants and exclusions

289 We used Prolific to recruit 183 adults resident in the UK who self-reported that their  
290 first language was English and that they had no known language disorders. They all  
291 gave informed consent to participate in the experiment.



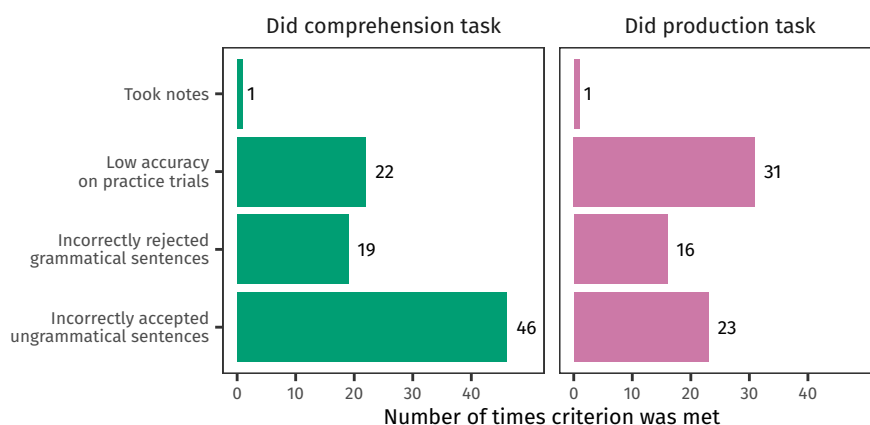


Figure 2: How many times each preregistered exclusion criterion was met in Experiment 1 (participants caught by more than one criterion contribute to each criterion’s count). On the whole, exclusion criteria were met more often in the COMPREHENSION condition than in the PRODUCTION condition.

292 The experiment took around 20 minutes to complete (median time = 17:38), and par-  
 293 ticipants were paid £3.50 (above UK National Minimum Wage at the time of running the  
 294 experiment). Participants were randomly assigned to either the COMPREHENSION condi-  
 295 tion (100 people) or the PRODUCTION condition (83 people). We excluded 103 participants  
 296 for the following preregistered reasons: self-reporting the use of written notes in an exit  
 297 questionnaire contrary to instructions (2); low accuracy (< 77.7%, i.e., 14/18) on practice  
 298 trials (16), testing trials (49) or both (36).

299 Figure 2 illustrates how many times each exclusion criterion was met in each condi-  
 300 tion. This plot does not reflect how the criteria may overlap, so participants caught by  
 301 multiple criteria contribute to multiple counts; see Appendix B for a full breakdown of  
 302 how many participants were caught by each combination of criteria.

303 We had to exclude many more participants who had been originally recruited into  
 304 the COMPREHENSION condition, and fewer who were recruited into the PRODUCTION con-  
 305 dition. This asymmetry indicates at least anecdotally that the production task does seem  
 306 to have helped participants learn the sentences that they were exposed to—in line with  
 307 previous evidence that production is good for learning.

308 After exclusions, we were left with analysable data from 40 participants in each con-  
 309 dition. (Appendix C contains the same analysis that we report below run on all 183  
 310 participants.) The remaining participants’ accuracy on the grammatical and ungram-  
 311 matical sentences was all close to ceiling (naturally, since these are the participants who  
 312 were not excluded for low accuracy), and there were no substantial differences between  
 313 conditions. For the COMPREHENSION group, grammatical sentences were correctly ac-  
 314 cepted 96% of the time, and ungrammatical sentences were correctly rejected 98% of the  
 315 time. And for the PRODUCTION group, grammatical sentences were also correctly ac-  
 316 cepted 96% of the time, and ungrammatical sentences were correctly rejected 97% of the  
 317 time.

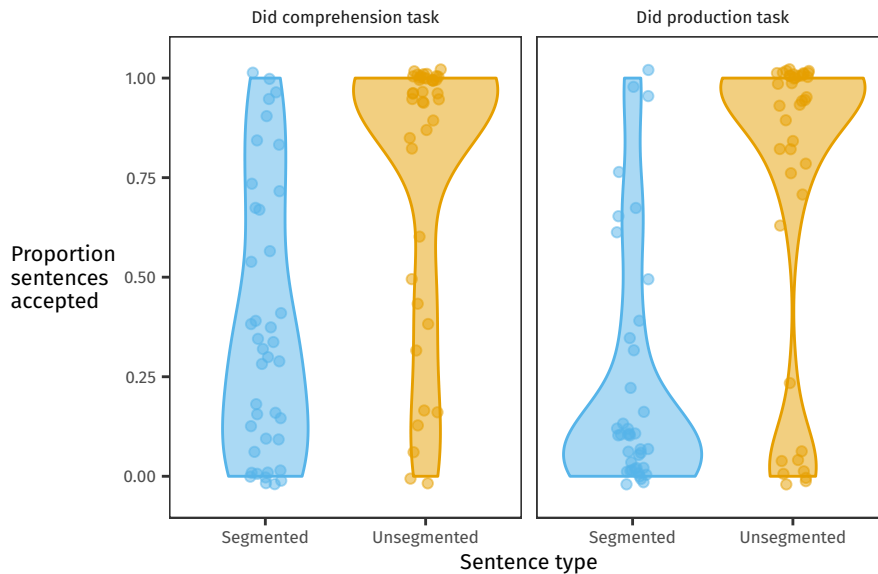


Figure 3: Participants in both the COMPREHENSION and PRODUCTION conditions of Experiment 1 accepted novel sentences that followed the unsegmented analysis more frequently than sentences that followed the segmented analysis. Each dot represents one participant’s proportion of accepted sentences of each type.

## 318 2.4 Results

### 319 2.4.1 Judgement

320 Participants in both conditions tended to accept novel sentences formed using the un-  
 321 segmented analysis, and they were more ambivalent about novel sentences formed using  
 322 the segmented analysis. Figure 3 shows the proportion of each kind of novel sentence  
 323 that participants accepted.

324 Following our preregistered analysis plan, we used brms (Bürkner, 2017) in R (R  
 325 Core Team, 2024) to fit a Bayesian linear model with a Bernoulli likelihood to this data.  
 326 This model predicts sentence acceptance as a function of condition (COMPREHENSION  
 327 versus PRODUCTION), sentence type (SEGMENTED versus UNSEGMENTED), and their inter-  
 328 action. The group-level effects in the model included varying intercepts by participant  
 329 and varying slopes over sentence type by participant. We selected the model’s weakly  
 330 regularising priors using prior predictive checks. The model converged, as indicated by  
 331 all Rhats = 1.00. Appendix D contains the full model specification.

332 We sum-coded condition (COMPREHENSION as  $-0.5$ , PRODUCTION as  $+0.5$ ) and sen-  
 333 tence type (SEGMENTED as  $-0.5$ , UNSEGMENTED as  $+0.5$ ). The interaction term was also  
 334 scaled to  $\pm 0.5$  so that we could use the same weakly regularising prior for all three pre-  
 335 dictors.

336 We hypothesised that, if a production task helps participants learn morphological  
 337 rules, then participants in the PRODUCTION condition would be more likely to accept  
 338 sentences generated by the segmented analysis than participants in the COMPREHENSION  
 339 condition. We would see this in the model as an interaction between condition and  
 340 sentence type.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	0.54	0.27	0.03	1.10
Condition	-0.81	0.51	-1.79	0.20
Sentence type	4.10	0.70	2.76	5.51
Condition:Sent. type	0.37	0.66	-0.88	1.70

Table 1: The posterior probability distributions estimated by the model for the English participants' sentence acceptance data in Experiment 1. Values are on the log-odds scale.

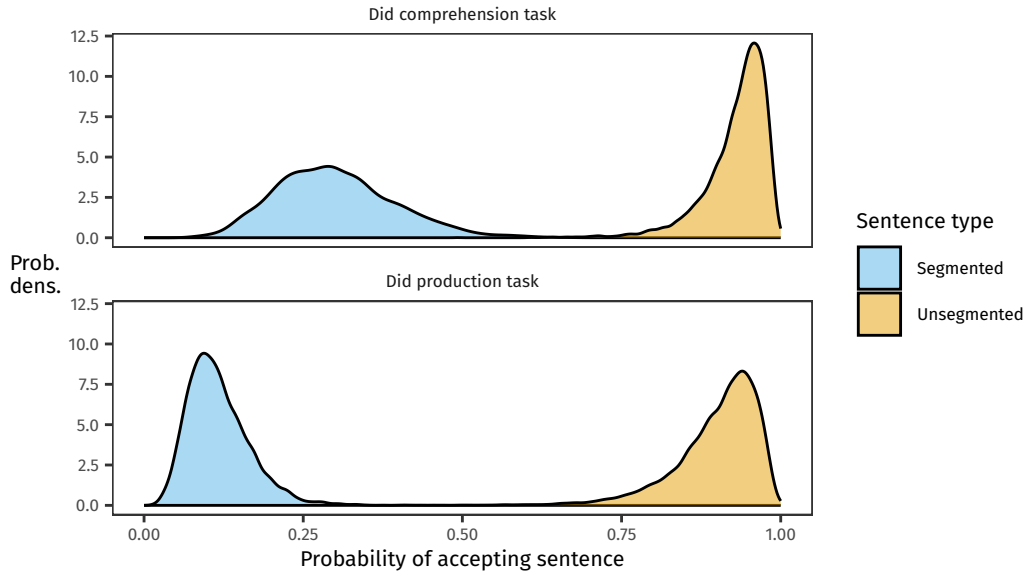


Figure 4: Conditional posterior probability distributions of the probability of accepting a sentence in Experiment 1. UNSEGMENTED sentences are more likely to be accepted than SEGMENTED sentences, regardless of whether participants did a comprehension or production task.

341 The model's posterior estimates for the population-level effects are summarised in  
342 Table 1. Figure 4 shows the conditional posterior probability distributions—that is, the  
343 posterior distributions over the probabilities of accepting a sentence for all combinations  
344 of condition and sentence type.

345 Overall, the model indicates with high certainty that participants are more likely to  
346 accept a novel sentence formed with the unsegmented analysis compared to a novel sen-  
347 tence formed with the segmented analysis. Concerning condition, the model's estimates  
348 indicate uncertainty about a difference in sentence acceptance probabilities between the  
349 PRODUCTION condition and the COMPREHENSION condition, as well as uncertainty about  
350 the interaction that our hypothesis targeted. Our prediction that participants who did  
351 the production task would be more likely to accept the novel segmented sentences was  
352 not borne out, and in fact, the results tend slightly in the opposite direction.

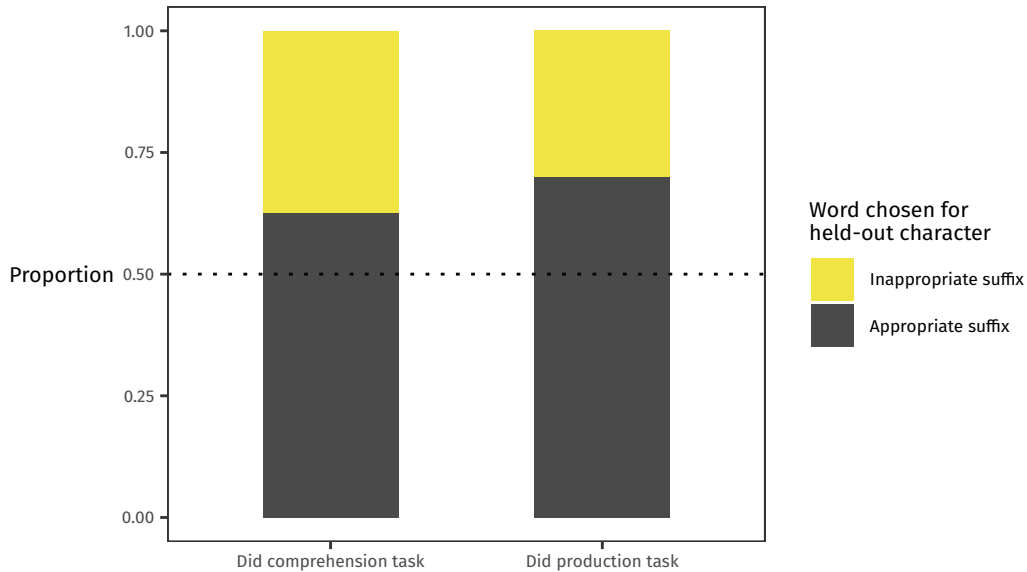


Figure 5: In the held-out character naming task of Experiment 1, more than half of participants in each condition selected the word with the appropriate suffix. Slightly more participants in the PRODUCTION condition selected the appropriate suffix.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	0.68	0.24	0.23	1.16
Condition	0.33	0.47	-0.57	1.22

Table 2: The posterior probability distributions estimated by the model for the English participants' held-out character naming data in Experiment 1. Values are on the log-odds scale.

### 353 2.4.2 Held-out character naming

354 Figure 5 shows the proportion of participants who chose a label for the held-out char-  
 355 acter that contained the appropriate suffix. Even if they didn't arrive at a fully-fledged  
 356 case marking analysis, more than half of the participants in each condition seem to have  
 357 noticed that each noun reliably ends in a particular syllable.

358 We preregistered the analysis of this data as exploratory. To see whether participants  
 359 in the PRODUCTION condition showed greater awareness of these morphological patterns  
 360 (even if they did not analyse them as case markers *per se*), we fit a Bayesian linear model  
 361 with a Bernoulli likelihood to this data, predicting appropriate suffix choice as a function  
 362 of condition (COMPREHENSION coded as -0.5, PRODUCTION as +0.5). Every participant  
 363 gave only one data point, so no group-level effects were needed. We used the same  
 364 weakly regularising priors as in other Bernoulli models reported in this paper. The model  
 365 converged, as indicated by all Rhats = 1.00.

366 Table 2 summarises the posterior distributions of the population-level effects esti-  
 367 mated by this model, and Figure 6 shows the conditional posterior probability distribu-  
 368 tions over the probabilities of selecting the appropriate suffix.

369 The model indicates that participants in both groups chose the label containing the

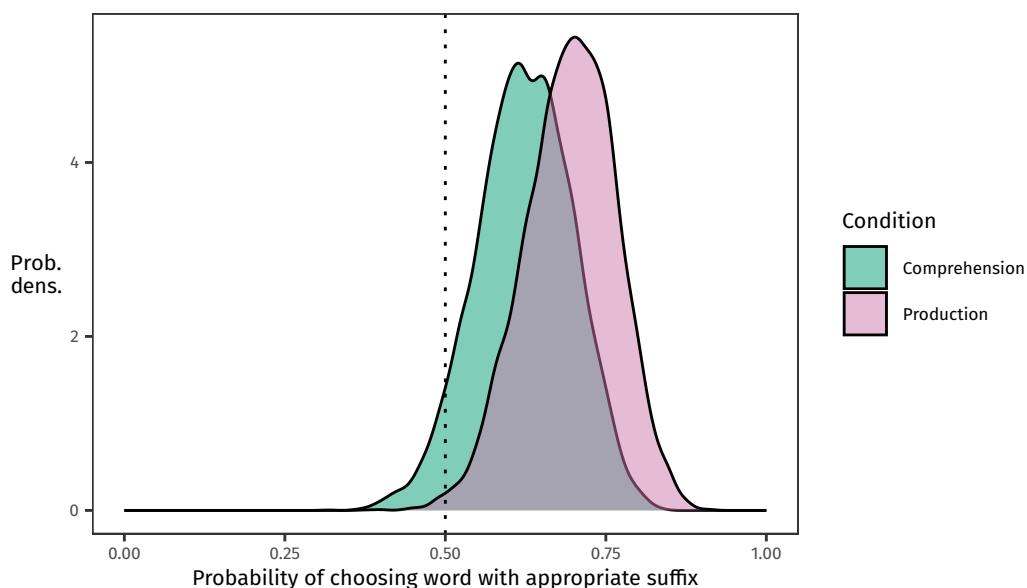


Figure 6: Conditional posterior probability distributions over the probability of selecting a word that contains the appropriate suffix in Experiment 1. The overlap of these posteriors suggests uncertainty about whether and how much the groups might differ.

370 appropriate suffix for the missing word with a probability slightly greater than chance.  
 371 Although being in the PRODUCTION condition is associated with a slightly higher prob-  
 372 ability of choosing the appropriate label, there is a great deal of overlap between condi-  
 373 tions and thus a great deal of uncertainty about whether participants in either condition  
 374 are more likely to select the appropriate label.

## 375 2.5 Interim discussion

376 Experiment 1 has shown that participants in both groups overwhelmingly preferred  
 377 novel sentences formed using the unsegmented analysis over sentences formed using  
 378 the segmented analysis. This preference was unaffected by whether participants com-  
 379 pleted a production or comprehension practice task, counter to our hypothesis.

380 Interestingly, the preference for the unsegmented analysis was resounding, even  
 381 though the held-out character naming task indicated that many participants were aware  
 382 of a morphological pattern in the language—namely that the first noun always ends in  
 383 a particular syllable (the nominative marker) while the second noun always ends in an-  
 384 other one (the accusative marker).

385 One straightforward explanation for this result is that the L1 English participants in  
 386 this experiment might not have arrived at a case marking analysis because case is not  
 387 morphologically marked outside of the pronominal system in English. In other words,  
 388 English uses word order alone to indicate grammatical roles, and thus our participants  
 389 may have been particularly unlikely to look beyond word order to notice that the case-  
 390 marking suffixes also indicated these roles. We collected data about further languages  
 391 that participants know or understand, and, in an exploratory analysis, compared the  
 392 performance of participants who do know a case-marking language (15 people) to those

393 who do not (65). The pattern of results remains the same; see Appendix E for details.

394 Nonetheless, it is possible that a population whose L1 includes more widespread use  
395 of case would be more likely to access the case marking analysis. We therefore ran a  
396 follow-up experiment with L1 speakers of German, a language with a productive case  
397 marking system featuring (among other cases) nominative and accusative differentially  
398 marked on nominal dependents like determiners.

## 399 **3 Experiment 2**

### 400 **3.1 Materials**

401 We used largely the same materials as in Experiment 1, described above in Section 2.1.  
402 Only a handful of changes were made for German-speaking participants.

403 First, we removed any forms from the language that resembled German words: *zug*  
404 is like German *Zug* ‘train’, *kex* might be read as *Keks* ‘cookie’, and so on.

405 Second, to ensure that the full set of stimuli was grammatically equivalent in German,  
406 we removed all images containing the pirate character. The German word *Pirat* is a so-  
407 called “strong masculine” noun: a noun that itself inflects for case, in addition to the usual  
408 inflection on the determiner (cf. nominative *der Pirat* ‘the pirate’, accusative *den Piraten*).  
409 All other characters correspond to German nouns that are grammatically “weak”, that  
410 is, the nouns don’t inflect for case.

411 Third, we changed the default word order from SOV to VSO, because SOV is the basic  
412 word order of German (Haftka, 1996; Haider, 2020). This means that the Experiment 1  
413 sentence *fuvu zijo gix* would become *gix fuvu zijo* in Experiment 2. The “ungrammatical”  
414 word order in the judgement phase remained SVO, akin to German’s V2 (though we did  
415 not use rejection of SVO sentences as a criterion for excluding participants in Experiment  
416 2; we will discuss this further in Section 3.3).

### 417 **3.2 Procedure**

418 Experiment 2 followed the same procedure as Experiment 1 (see Section 2.2), with one  
419 modification. For English participants, we had randomly mapped the keys *f* and *j* to  
420 ‘yes’ and ‘no’. Since German *ja* ‘yes’ begins with *J*, we instead used *p* and *q* as the  
421 decision keys for the sentence judgement task.

### 422 **3.3 Participants and exclusions**

423 We used Prolific to recruit 135 participants who self-reported that their first language  
424 was German and that they had no known language disorders. They all gave informed  
425 consent to participate in the experiment.

426 The experiment took around 20 minutes to complete (median time = 17:39), and par-  
427 ticipants were paid £3.85 (approx. €4.50), above UK National Minimum Wage at the time  
428 of running the experiment. As in Experiment 1, participants were randomly assigned to  
429 either the COMPREHENSION condition (68 people) or the PRODUCTION condition (67 peo-  
430 ple). We excluded 43 participants for the following preregistered reasons: low accuracy  
431 on practice trials (17), GRAMMATICAL testing trials (12), or both (14). Figure 7 illustrates

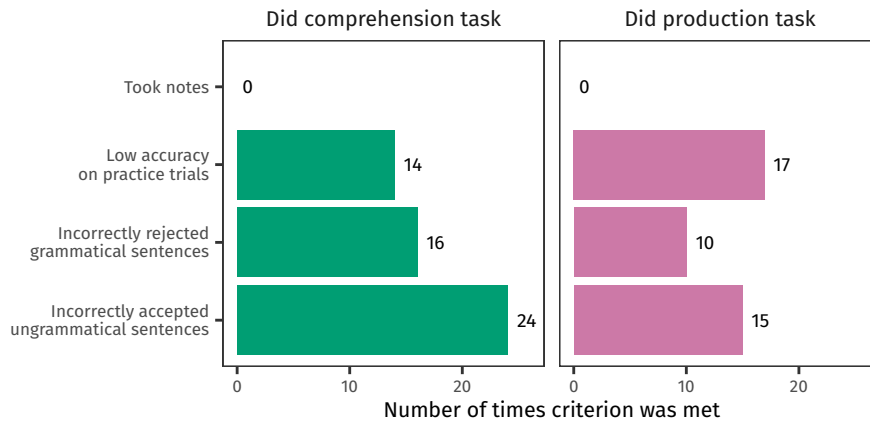


Figure 7: How many times each preregistered exclusion criterion was met in Experiment 2 (participants caught by more than one criterion contribute to each criterion’s count). Exclusions are more balanced between conditions in Experiment 2 compared to Experiment 1, though still, more participants in the COMPREHENSION group compared to the PRODUCTION group incorrectly rejected sentences that were grammatical. (The ungrammatical sentences criterion is included here only for completeness; in Experiment 2 it was not used to exclude participants.)

432 how many times each exclusion criterion was met in each condition (note that this plot  
 433 does not reflect how criteria may overlap, so participants caught by multiple criteria  
 434 contribute to multiple counts).

435 This figure includes German participants’ performance on the so-called “ungram-  
 436 matical” sentences, the ones with word order that differs from training, though we did  
 437 not use this criterion to exclude participants from the analysis. We ignored this criterion  
 438 for German speakers because German permits a relatively free word order, so partici-  
 439 pants may not have had the expectation that word order should be fixed, particularly if  
 440 they accessed the segmented (case marking) analysis. Recall that removing this criterion  
 441 for the English-speaking participants in Experiment 1 did not affect the pattern of results  
 442 (Appendix A).

443 After exclusions, we were left with data from 46 participants in each condition. The  
 444 remaining participants’ accuracy on the grammatical and ungrammatical sentences was  
 445 fairly high, with no substantial differences between conditions. For the COMPREHENSION  
 446 group, grammatical sentences were correctly accepted 96% of the time, and ungrammat-  
 447 ical sentences were correctly rejected 78% of the time. And for the PRODUCTION group,  
 448 grammatical sentences were also correctly accepted 96% of the time, and ungrammatical  
 449 sentences were correctly rejected 82% of the time.

### 450 3.4 Results

451 Overall, the results from the German participants in Experiment 2 are similar to the  
 452 results from the English participants in Experiment 1.

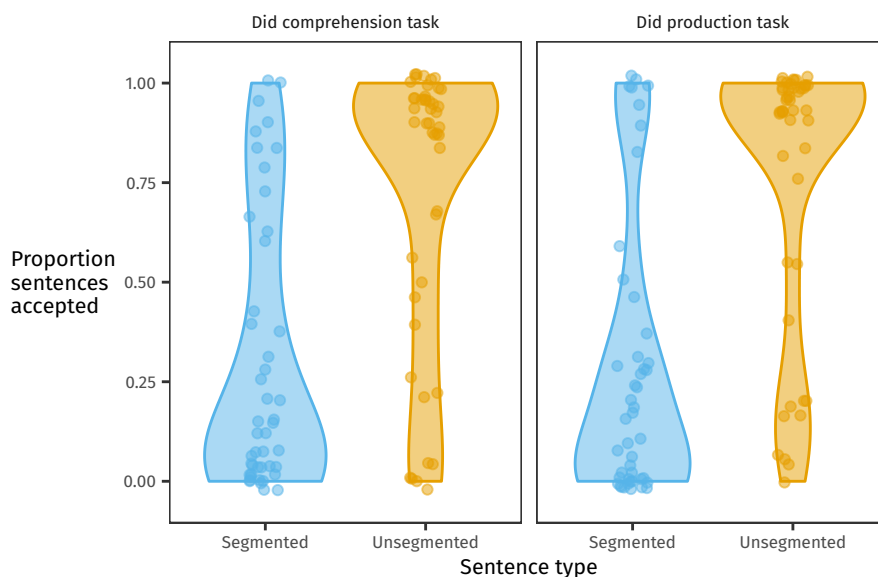


Figure 8: In Experiment 2, participants in both the COMPREHENSION and PRODUCTION conditions again accepted novel sentences that followed the unsegmented analysis more frequently than sentences that followed the segmented analysis.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	0.17	0.21	-0.24	0.59
Condition	0.33	0.40	-0.45	1.13
Sentence type	3.84	0.59	2.68	5.01
Condition:Sent. type	0.52	0.58	-0.61	1.68

Table 3: The posterior probability distributions estimated by the model for the German participants' sentence acceptance data in Experiment 2. Values are on the log-odds scale.

### 453 3.4.1 Judgement

454 Like the English-speaking participants, the German participants showed a clear prefer-  
 455 ence for the unsegmented analysis (see Figure 8), even though German has productive  
 456 morphological case marking. We fit the same Bayesian linear model as described above  
 457 in Section 2.2.3 to the data from the German participants. The posterior distributions for  
 458 the population-level effects estimated by the model are given in Table 3, and the condi-  
 459 tional posterior probability distributions are shown in Figure 9. Again, the interaction  
 460 that would support our hypothesis about a production task enabling participants to learn  
 461 the segmented analysis was not borne out.

### 462 3.4.2 Held-out character naming

463 About three-quarters of German participants appear to have noticed that one noun al-  
 464 ways ends in *-vu* and the other always ends in *-jo*; see Figure 10.

465 We fit the same model as described in Section 2.2.4 to this data. Table 4 summarises  
 466 the posterior distributions of the population-level effects, and Figure 11 shows the con-  
 467 ditional posterior probability distributions over the probabilities of selecting the appro-



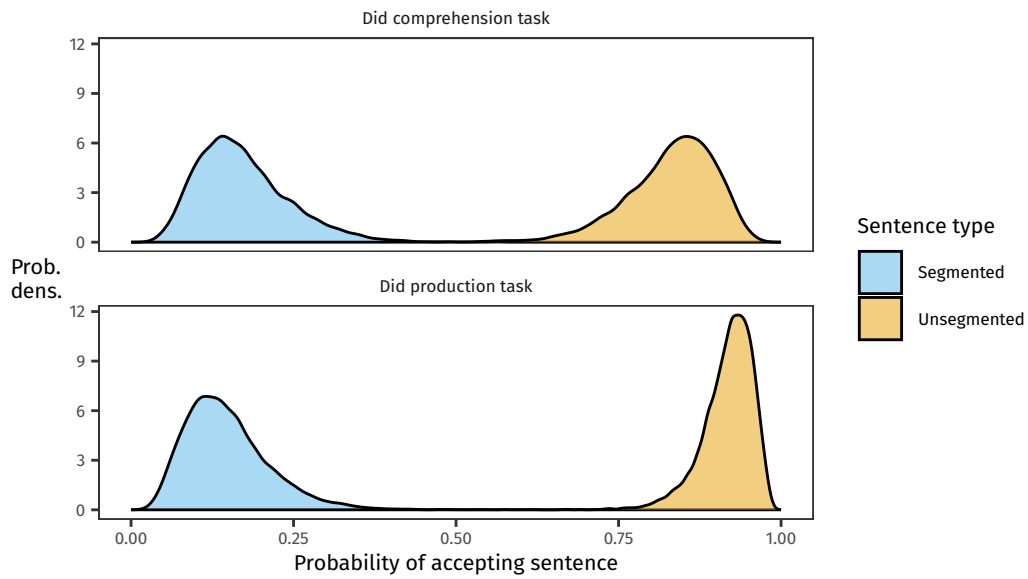


Figure 9: Conditional posterior probability distributions of the probability of accepting a sentence for the participants in Experiment 2. As in Experiment 1, UNSEGMENTED sentences are more likely to be accepted than SEGMENTED sentences, regardless of whether participants did a comprehension or production task.

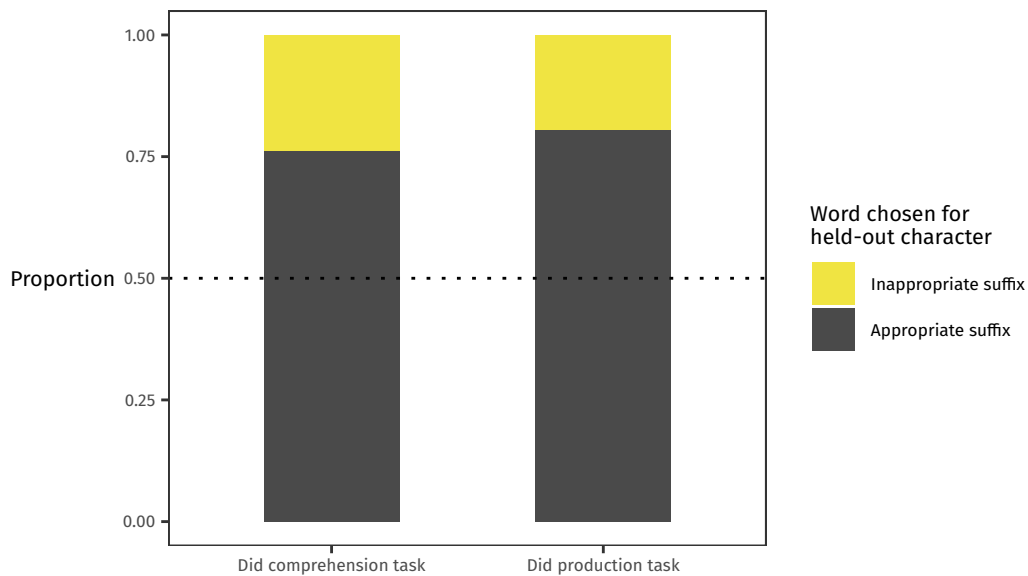


Figure 10: In the held-out character naming task of Experiment 2, around three-quarters of German participants selected the form in which the word ended in the appropriate suffix; the proportion of appropriate choices is slightly higher for the production group.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	1.28	0.25	0.80	1.79
Condition	0.24	0.50	-0.73	1.22

Table 4: The posterior probability distributions estimated by the model for the German participants' held-out character naming data in Experiment 2. Values are on the log-odds scale.

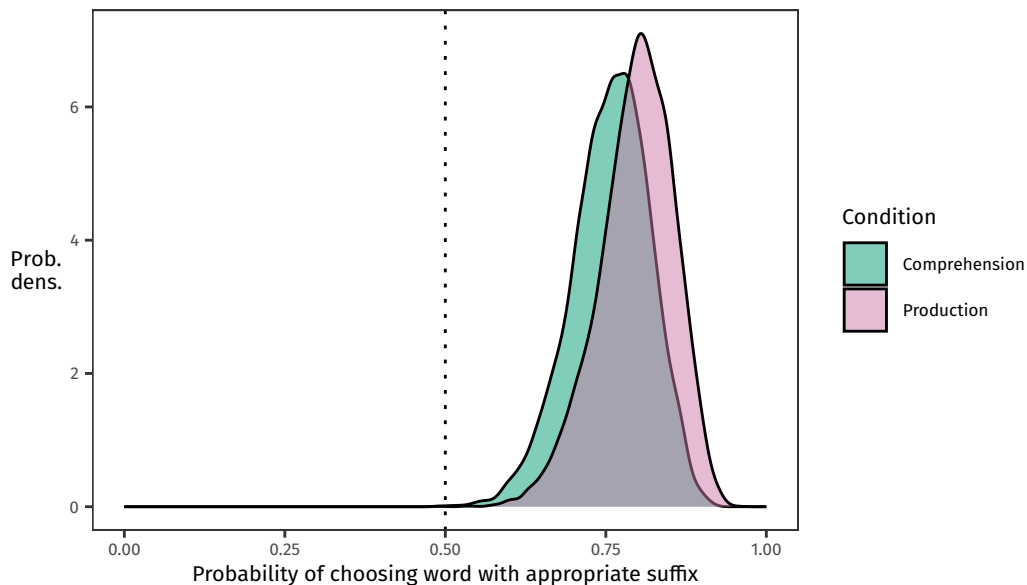


Figure 11: Conditional posterior probability distributions over the probability of German participants selecting a word that contains the appropriate suffix in Experiment 2. These posteriors overlap, so we are not certain whether and how much the groups might differ.

468 priate suffix. The model suggests that, much like the English participants, the German  
 469 group is likely to have labelled the held-out character following the morphological pat-  
 470 tern, and there is no clear association between participants' choice of label for the held-  
 471 out character and experimental condition.

## 472 4 General discussion

473 In two artificial language learning experiments, we tested whether a production task—  
 474 known to improve rule learning in a number of contexts—could also draw adult learners'  
 475 attention to rules that adults typically disprefer. Specifically, we focused on morpholog-  
 476 ical marking of thematic roles using case suffixes.

477 We trained participants on a language with fixed word order in which agent nouns  
 478 and patient nouns were always marked with distinct suffixes (e.g., *-vu* for agents and *-jo*  
 479 for patients). However, the nouns that each participant saw only ever occurred as either  
 480 agents or patients, never in both roles. Thus the suffixes could be analysed as part of the  
 481 nouns themselves (an unsegmented analysis) or as productive endings, part of a wider  
 482 case system (a segmented analysis).

483 We found that regardless of whether participants did a production or a comprehen-  
484 sion practice task, they favoured novel sentences which were formed using an unseg-  
485 mented, word-level analysis, and they tended to reject sentences formed using a seg-  
486 mented, case-marking analysis. In other words, when shown novel scenes in which  
487 familiar characters featured in a novel grammatical role (e.g., where the fairy, which ap-  
488 peared only as an agent in training, appeared as the patient), they tended to reject sen-  
489 tences in which the noun suffixes were adjusted to reflect these new grammatical roles.  
490 Nevertheless, most participants detected the morphological pattern that resulted from  
491 the case marking (i.e., that one noun in every sentence ended in *vu* and the other in *jo*),  
492 even if they did not necessarily develop this observation into a productive case marking  
493 grammar. Perhaps surprisingly, we found that the same pattern of results—sensitivity  
494 to the morphological patterns but failure to accept sentences formed according to the  
495 segmented analysis, regardless of practice condition—also held for participants whose  
496 first language, German, has extensive case marking.

497 In a sense, reanalysing an unsegmented word-order-based grammar into a case mark-  
498 ing grammar is not a trivial task, since it means overriding the chunks that have already  
499 been learned. But it is something that learners of genuine case marking languages are  
500 likely to need to do—many nouns are more likely to occur in a particular grammatical  
501 role, e.g., humans and other animate beings are more commonly found as agents than  
502 as patients (Croft, 2003; Meir et al., 2017; Silverstein, 1976). So it is not unreasonable to  
503 expect our participants to be able to break down the chunks they have learned.

504 In short, then, we have found no evidence that production tasks have an advantage  
505 over comprehension tasks for helping adult learners acquire a more difficult morpholog-  
506 ical rule over a more available word-level one.

507 Of course, as with any comparison of comprehension and production, it is difficult to  
508 be sure that we have isolated the relevant mechanism that makes production help learn-  
509 ers. For example, one potential criticism of our study is that our production task required  
510 participants to click on buttons to build up a sentence syllable-by-syllable, rather than  
511 to produce the sentence aloud themselves. Perhaps this was not active enough to elicit  
512 the benefits of language production that previous research describes. However, we find  
513 this explanation unlikely. As mentioned above, studies on the effect of different kinds  
514 of tests (e.g., short answer, multiple choice) have found that any kind of testing can im-  
515 prove learning over a passive rereading or recognition task (Kang et al., 2007; McDaniel  
516 et al., 2007; McDermott et al., 2014). These studies do suggest, however, that the degree  
517 of improvement may not be the same between all kinds of test. So perhaps the kind of  
518 production task we did only got participants part of the way, not as far as they may have  
519 come with from-scratch verbal production.

520 In our view, a more plausible explanation for why we failed to find an improve-  
521 ment with production is that participants were not required to produce the language  
522 early enough in the learning process: the critical practice phase came after an initial  
523 training phase. We designed the task in this way because we were concerned about  
524 disproportional attrition of participants in the PRODUCTION condition compared to the  
525 COMPREHENSION condition. Constructing sentences in an unfamiliar language is a much  
526 more challenging task than choosing between pictures, and we didn't want PRODUCTION  
527 participants to be discouraged (and potentially stop the study at disproportionate rates)  
528 by introducing this extra level of difficulty too early. However, introducing the different

529 practice tasks too late also has a downside: participants may have already discovered the  
530 fixed word order rule during the training phase, and since that rule perfectly explains  
531 all the data they encountered, there was no need to search for further explanations (in  
532 classical conditioning terms, an *overshadowing* effect; Pavlov, 1927). This pattern of be-  
533 haviour is characteristic of adults in non-linguistic tasks too: adults tend to identify a  
534 reliable cue and then exploit it, while children continue to explore (Liquin and Gopnik,  
535 2022; Sumner et al., 2019). A possible prediction of this account, then, is that children  
536 might be more likely than adults to accept the case marking analysis.

537 The late start of the production/comprehension tasks could also be part of why our  
538 results differ from those of Hopman and MacDonald (2018), who observe that a produc-  
539 tion task leads to slightly better learning of morphological rules than word order rules.  
540 In their design, passive exposure trials were interleaved with blocks of active production  
541 trials. And their experiment seems to have been conducted in person, a factor likely to  
542 prevent participants from withdrawing from the experiment early, compared to experi-  
543 ments run online.

#### 544 **4.1 Starting big, but noticing small**

545 Our results, and previous findings showing that adults struggle to learn morphological  
546 rules, align with the observation that as language is learned, linguistic information tends  
547 to be stored first as holistic chunks (Christiansen and Chater, 2016). The morphosyn-  
548 tactic rules that might underlie parts of those chunks are only induced when learners  
549 receive sufficient evidence from the input. This idea has been referred to as “needs-  
550 only analysis” (Wray, 2002, 2006), and it is closely related to the “starting big” approach  
551 described by Havron and Arnon (2021) and Siegelman and Arnon (2015).

552 We do see our adult participants “starting big” by learning a rule that manipulates  
553 the larger, unsegmented units, not the smaller ones that require word segmentation.  
554 But our results from the held-out character naming task also add some nuance to this  
555 notion. Participants still notice patterns and pieces within the word-level chunks they  
556 manipulate. In other words, the chunks they learn aren’t fully opaque.

557 This is interesting in connection with anecdotal evidence of many potentially seg-  
558 mentable chunks being learned holistically. For example, people might have a moment  
559 of surprise when they realise, say, that a safety pin is so called because it is safe, or that  
560 dry cleaning is a kind of cleaning which doesn’t involve water. Those situations are evi-  
561 dently different from the suffixing pattern in our artificial language. The suffixes seem to  
562 have been salient enough for people to notice, even if the noticing doesn’t cause learners  
563 to override the chunk they learned.

564 By noticing this pattern at all, though, learners have taken the first step toward such  
565 an analysis. How might we push them to make the leap? Evidence from many learn-  
566 ing domains suggests that learners need more data—more unique examples of a rule in  
567 action, in varying contexts—to move beyond item-by-item learning to systematic rule  
568 induction (see Raviv et al., 2022 for a review). When the rule to be acquired is dispre-  
569 ferred *a priori*, either more data or conflicting data (see next section) may be required to  
570 overcome participants’ strong prior preferences.

571 However, there is a trade-off here. More variable input does lead to better long-term  
572 generalisation, but it also hinders initial learning (Raviv et al., 2022). And in an exper-

573 imental setting, these may be hard to balance. Logistics and finances limit how much  
574 training participants can receive before they need to provide useful data. It would be dif-  
575 ficult within the current experimental design to include the large amount of variability  
576 required, while also ensuring that participants learned the language adequately.

## 577 **4.2 Outlook and future directions**

578 We see two interesting options for follow-up research based on our results and the ob-  
579 servations outlined above. First and most obviously, the production and comprehen-  
580 sion tasks could be interspersed throughout the training phase. The greater attrition  
581 rate that we would expect from this design could be handled either by significant over-  
582 recruitment of participants or in-person administration of the experiment.

583 Second, knowing now that adult learners will tend to adopt the unsegmented anal-  
584 ysis, we might ask: how difficult would it be to pivot from that initial analysis to a seg-  
585 mented one in the face of new, conflicting data? The under-specified sentences shown  
586 during training could be changed partway through to become sentences unambiguously  
587 formed using a case marking rule, incompatible with the unsegmented analysis that par-  
588 ticipants would presumably have learned. Or alternatively, perhaps just one character  
589 alternates thematic roles, or one character is irregular and takes no suffix at all.

590 We could imagine two reasons why a production task might help participants more  
591 swiftly reanalyse their data in the face of this conflicting evidence. First, according to the  
592 noticing mechanism (Swain, 2005), production may help participants more quickly iden-  
593 tify the morphological patterns that are now the only way to fully explain the data. Sec-  
594 ond, the production advantage has also been explained in terms of retrieval from mem-  
595 ory, since retrieval practice is known to strengthen learning (Hopman and MacDonald,  
596 2018; Karpicke, 2012; Karpicke and Roediger, 2008; MacDonald, 2013). Under this mech-  
597 anism, production could help participants recall whatever units they had stored—likely  
598 unsegmented ones (Christiansen and Chater, 2016; Havron and Arnon, 2021; Siegelman  
599 and Arnon, 2015)—to render them more available for reanalysis.

## 600 **5 Conclusion**

601 We began this investigation where two strands of previous research intersect, one show-  
602 ing that that language production helps learners identify and learn rules in their lan-  
603 guage (Hopman and MacDonald, 2018; Izumi, 2002; Swain, 2005), and another showing  
604 that adults struggle to learn morphological rules and prefer word-level ones (Havron  
605 and Arnon, 2021; Jordens et al., 1989; Lupyan and Dale, 2010; Papadopoulou et al., 2011;  
606 Parodi et al., 2004). Bringing these observations together, we wanted to know whether a  
607 production task could help adult learners to identify a more difficult morphological rule  
608 over a more available word-level one.

609 Our results demonstrate that adults prefer to learn a word-level rule for marking  
610 thematic role over a morphological rule, even when they appear to notice morphological  
611 patterns. Contrary to our preregistered hypothesis, practising a new language with a  
612 production task does not steer learners away from this strong preference for word-level  
613 rules. This holds for speakers of both English (Experiment 1) and German (Experiment

614 2), indicating that even adult learners familiar with case marking tend to prefer the word-  
615 level rule over the morphological one.

616 Although we found no evidence that production tasks are better than comprehension  
617 tasks for helping adult learners acquire morphological rules, these two experiments nev-  
618 ertheless clearly illustrate that adults strongly prefer to learn rules operating over larger  
619 units, rather than smaller ones. This supports assumptions made in, e.g., Bentz and Win-  
620 ter (2013) and Lupyán and Dale (2010), for whom adult preference for word-level rules is  
621 central for explaining why languages with more adult learners tend to have less complex  
622 morphology.

623 All in all, although our hypothesis about the role of production for morphological  
624 learning was not borne out, this study has still opened several doors that we believe are  
625 worth passing through to discover more about the interplay of language production and  
626 the kinds of rules that learners learn.

## References

- 627
- 628 Bentz, C. and Winter, B. (2013). Languages with More Second Language Learners Tend  
629 to Lose Nominal Case. *Language Dynamics and Change*, 3(1):1–27. 2, 3, 22, 26
- 630 Bohman, T. M., Bedore, L. M., Peña, E. D., Mendez-Perez, A., and Gillam, R. B. (2010).  
631 What you hear and what you say: Language performance in Spanish–English bilin-  
632 guals. *International Journal of Bilingual Education and Bilingualism*, 13(3):325–344.  
633 1
- 634 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan.  
635 *Journal of Statistical Software*, 80(1):1–28. 10
- 636 Christiansen, M. H. and Chater, N. (2016). The Now-or-Never bottleneck: A fundamental  
637 constraint on language. *Behavioral and Brain Sciences*, 39:e62. 20, 21
- 638 Croft, W. (2003). *Typology and Universals*. Cambridge University Press, Cambridge, 2  
639 edition. 19
- 640 de Leeuw, J. R., Gilbert, R. A., and Luchterhandt, B. (2023). jsPsych: Enabling an Open-  
641 Source CollaborativeEcosystem of Behavioral Experiments. *Journal of Open Source*  
642 *Software*, 8(85):5351. 5
- 643 Donnelly, S. and Kidd, E. (2021). The longitudinal relationship between conversational  
644 turn-taking and vocabulary growth in early language development. *Child Develop-*  
645 *ment*, 92(2):609–625. 1
- 646 Fedzechkina, M., Jaeger, T. F., and Newport, E. L. (2011). Functional biases in language  
647 learning: Evidence from word order and case-marking interaction. In *Proceedings of*  
648 *the Annual Meeting of the Cognitive Science Society*, volume 33. 3, 7, 26
- 649 Haftka, B. (1996). Deutsch ist eine V/2-Sprache mit Verbendstellung und freier Wortfolge.  
650 In Lang, E. and Zifonun, G., editors, *Deutsch - Typologisch*, pages 121–141. De Gruyter.  
651 14
- 652 Haider, H. (2020). VO-/OV-Base Ordering. In Putnam, M. T. and Page, B. R., editors, *The*  
653 *Cambridge Handbook of Germanic Linguistics*, pages 339–364. Cambridge University  
654 Press, 1 edition. 14
- 655 Havron, N. and Arnon, I. (2021). Starting big: The effect of unit size on language learning  
656 in children and adults. *Journal of Child Language*, 48(2):244–260. 20, 21
- 657 Holmes, V. M. and Dejean De La Bâtie, B. (1999). Assignment of grammatical gen-  
658 der by native speakers and foreign learners of french. *Applied Psycholinguistics*,  
659 20(4):479–506. 2
- 660 Hopman, E. W. M. (2022). *Modality Matters: Generalization in Second Language Learning*  
661 *after Production versus Comprehension Practice*. PhD thesis, University of Wisconsin-  
662 Madison. 2

- 663 Hopman, E. W. M. and MacDonald, M. C. (2018). Production Practice During Language  
664 Learning Improves Comprehension. *Psychological Science*, 29(6):961–971. 2, 20, 21
- 665 Izumi, S. (2002). Output, input enhancement, and the noticing hypothesis: An experi-  
666 mental study on ESL relativization. *Studies in Second Language Acquisition*, 24(4):541–  
667 577. 1, 2, 21
- 668 Jordens, P., de Bot, K., and Trapman, H. (1989). Linguistic aspects of regression in Ger-  
669 man case marking. *Studies in Second Language Acquisition*, 11:179–204. 2, 21
- 670 Kang, S. H. K., McDermott, K. B., and Roediger, H. L. (2007). Test format and corrective  
671 feedback modify the effect of testing on long-term retention. *European Journal of*  
672 *Cognitive Psychology*, 19(4-5):528–558. 2, 19
- 673 Karpicke, J. D. (2012). Retrieval-based learning: Active retrieval promotes meaningful  
674 learning. *Current Directions in Psychological Science*, 21(3):157–163. 21
- 675 Karpicke, J. D. and Roediger, H. L. (2008). The critical importance of retrieval for learning.  
676 *Science (New York, N.Y.)*, 319(5865):966–968. 21
- 677 Keogh, A. and Lupyan, G. (in press). Who benefits from redundancy in learning noun  
678 class systems? In *Proceedings of the 15th International Conference on the Evolution of*  
679 *Language*. 2
- 680 Keppenne, V., Hopman, E. W. M., and Jackson, C. N. (2021). Production-based training  
681 benefits the comprehension and production of grammatical gender in L2 German.  
682 *Applied Psycholinguistics*, 42(4):907–936. 2
- 683 Levshina, N. (2020). Efficient trade-offs as explanations in functional linguistics: some  
684 problems and an alternative proposal. *Revista da ABRALIN*, 19(3):50–78. 3
- 685 Liquin, E. G. and Gopnik, A. (2022). Children are more exploratory and learn more than  
686 adults in an approach-avoid task. *Cognition*, 218:104940. 20
- 687 Lupyan, G. and Dale, R. (2010). Language structure is partly determined by social struc-  
688 ture. *PLoS ONE*, 5(1):e8559. 2, 3, 7, 21, 22
- 689 MacDonald, M. C. (2013). How language production shapes language form and compre-  
690 hension. *Frontiers in Psychology*, 4. 21
- 691 McDaniel, M. A., Anderson, J. L., Derbish, M. H., and Morrisette, N. (2007). Testing the  
692 testing effect in the classroom. *European Journal of Cognitive Psychology*, 19(4-5):494–  
693 513. 2, 19
- 694 McDermott, K. B., Agarwal, P. K., D’Antonio, L., Roediger, H. L., and McDaniel, M. A.  
695 (2014). Both multiple-choice and short-answer quizzes enhance later exam perfor-  
696 mance in middle and high school classes. *Journal of Experimental Psychology: Applied*,  
697 20(1):3–21. 2, 19



- 698 Meir, I., Aronoff, M., Börstell, C., Hwang, S.-O., Ilkbasaran, D., Kastner, I., Lopic, R., Ben-  
699 Basat, A. L., Padden, C., and Sandler, W. (2017). The effect of being human and the  
700 basis of grammatical word order: Insights from novel communication systems and  
701 young sign languages. *Cognition*, 158:189–207. 19
- 702 Papadopoulou, D., Varlokosta, S., Spyropoulous, V., Kaili, H., Prokou, S., and Revithiadou,  
703 A. (2011). Case morphology and word order in second language Turkish: Evidence  
704 from Greek learners. *Second Language Research*, 27:173–205. 2, 21
- 705 Parodi, T., Schwartz, B. D., and Clahsen, H. (2004). On the L2 acquisition of the mor-  
706 phosyntax of German nominals. *Linguistics*, 42:669–705. 2, 21
- 707 Pavlov, P. I. (1927). *Conditioned reflexes: an investigation of the physiological activity of*  
708 *the cerebral cortex*. Oxford University Press, London. 20
- 709 R Core Team (2024). *R: A Language and Environment for Statistical Computing*. R Foun-  
710 dation for Statistical Computing, Vienna, Austria. 10
- 711 Raviv, L., Lupyan, G., and Green, S. C. (2022). How variability shapes learning and  
712 generalization. *Trends in Cognitive Science*, 26(6):462–483. 20
- 713 Ribot, K. M., Hoff, E., and Burridge, A. (2018). Language use contributes to expressive  
714 language growth: Evidence from bilingual children. *Child Development*, 89(3):929–940.  
715 1
- 716 Rogers, M. (1987). Learners Difficulties with Grammatical Gender in German as a Foreign  
717 Language\*. *Applied Linguistics*, 8(1):48–74. 2
- 718 Sagarra, N. and Ellis, N. C. (2013). From seeing adverbs to seeing verbal morphology:  
719 Language experience and adult acquisition of L2 tense. *Studies in Second Language*  
720 *Acquisition*, 35(2):261–290. 2
- 721 Siegelman, N. and Arnon, I. (2015). The advantage of starting big: Learning from un-  
722 segmented input facilitates mastery of grammatical gender in an artificial language.  
723 *Journal of Memory and Language*, 85:60–75. 20, 21
- 724 Silverstein, M. (1976). Hierarchy of Features and Ergativity. In *Grammatical Categories*  
725 *in Australian Languages*, pages 163–232. De Gruyter. 19
- 726 Sumner, E. S., Li, A. X., Perfors, A., Hayes, B. K., Navarro, D. J., and Sarnecka, B. W.  
727 (2019). The exploration advantage: Children’s instinct to explore leads them to find  
728 information that adults miss. *PsyArXiv*, page 11. 20
- 729 Swain, M. (2005). The output hypothesis: Theory and research. In *Handbook of Research*  
730 *in Second Language Teaching and Learning*, pages 471–483. Routledge, New York, 1  
731 edition. 2, 21
- 732 Wray, A. (2002). *Formulaic Language and the Lexicon*. Cambridge University Press,  
733 Cambridge. 20
- 734 Wray, A. (2006). Formulaic language. In Brown, K., editor, *Encyclopedia of Language &*  
735 *Linguistics*, pages 590–597. Elsevier, Oxford, 2 edition. 20

## A Exploratory analysis: Removing the ungrammatical exclusion criterion

In the test phase of the experiment, we collected data that would inform two exclusion criteria: we would only keep participants who correctly accepted grammatical sentences and correctly rejected ungrammatical ones. We had defined “ungrammatical” as a word order that diverged from the one in training. Our reasoning was that participants should have learned that the language has SOV order, so they should reject the “ungrammatical” SVO order. Since SVO is also the basic word order of English, participants’ rejection of it would provide the strongest test that they had learned the word order of the artificial language.

However, if a language has case marking, it is likely to also have free word order (Bentz and Winter, 2013; Fedzechkina et al., 2011). It is therefore possible that participants who accepted sentences with a different word order had learned a case marking rule and associated that with a free word order, in which case our exclusion criterion would be removing exactly those participants who learned the segmented analysis we were targeting. This could explain why our results show such a strong preference for the unsegmented analysis.

Here, we lift this exclusion criterion and re-run the analyses described above. This criterion originally excluded 27 comprehension participants and 6 production participants; below we analyse data from 67 participants in the comprehension group and 46 in the production group.

### A.1 Judgement

Figure 12 shows a similar pattern to Figure 3: a general preference for the novel sentences formed using the unsegmented analysis, and greater ambivalence toward ones formed with the segmented analysis.

We fit the same model described in Section 2.4 to this data. The pattern of results (shown in Table 5) remains the same as above. We conclude that the “ungrammatical” word order criterion did not exclude participants who learned a case marking rule and then extrapolated from it that word order was free.

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.41	0.24	-0.07	0.89
Condition	-0.19	0.47	-1.12	0.72
Sentence type	3.48	0.54	2.43	4.58
Condition:Sent. type	0.88	0.56	-0.20	1.98

Table 5: Posterior distributions estimated by a model predicting sentence acceptance by condition, sentence type, and their interaction, now including data from participants originally excluded from Experiment 1 for rejecting sentences with a different word order than seen in training.

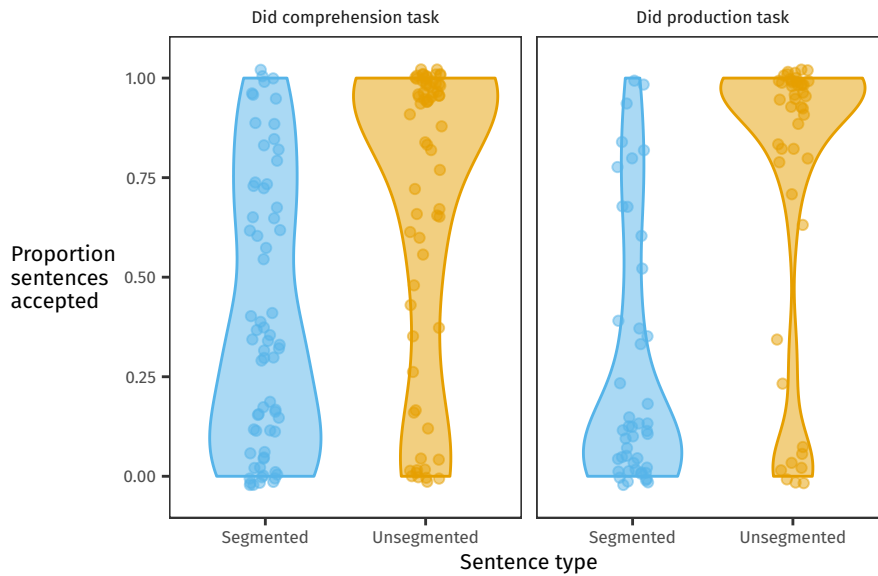


Figure 12: After lifting the ungrammaticality rejection criterion for participants in Experiment 1, the larger pool of participants show the same results: a strong preference for the unsegmented analysis over the segmented analysis, with no clear effect of task.

## 765 A.2 Held-out character naming

766 The results from the held-out character naming analysis also remain extremely similar  
 767 to the ones reported with the original exclusion criteria, as shown in Figure 13 and Table  
 768 6.

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.71	0.21	0.32	1.12
Condition	0.24	0.41	-0.55	1.06

Table 6: Posterior distributions estimated by a model predicting appropriate suffix choice by condition, now including data from participants originally excluded from Experiment 1 for rejecting sentences with a different word order than seen in training.

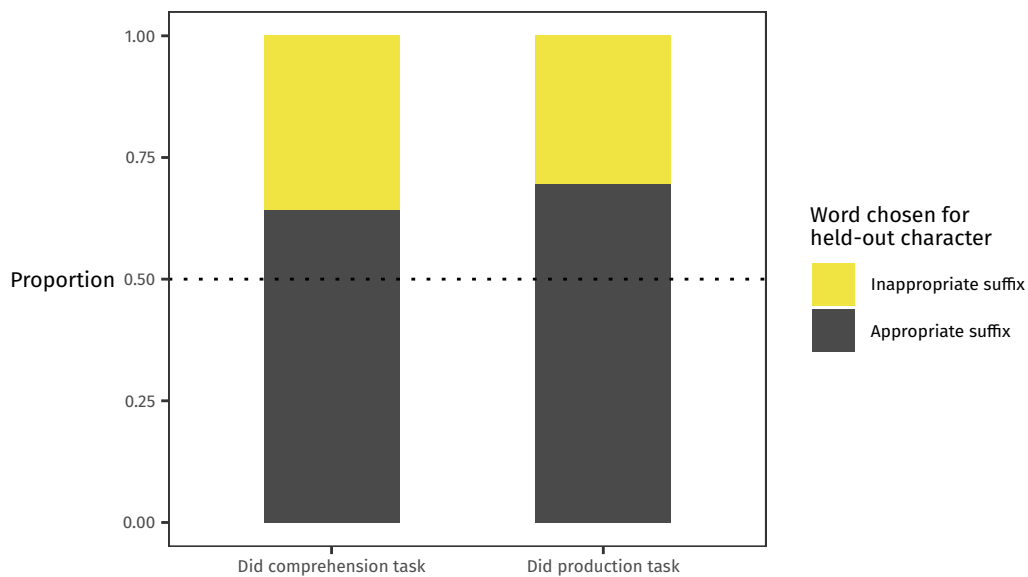


Figure 13: Proportion of Experiment 1 participants in each group, now including participants previously excluded from the analysis based on the ungrammaticality rejection criterion, who labelled the held-out character with a word containing the appropriate suffix. The same pattern holds as in the original analysis.

## 769 B Overlaps in exclusion criteria

770 The following table shows how many of the 183 participants recruited for Experiment  
 771 1 were caught by each combination of exclusion criteria. (Gram. = incorrectly rejected  
 772 grammatical sentences; Ungram. = incorrectly accepted ungrammatical sentences; Prac-  
 773 tice = low accuracy on practice phase; Notes = self-reported taking notes.)

Gram.	Ungram.	Practice	Notes	Comprehension	Production
				40	40
			×	0	1
		×		5	12
	×			27	6
	×		×	1	0
	×	×		8	8
×				5	1
×		×		4	6
×	×			5	4
×	×	×		5	5

775 The following table shows how many of the 135 participants recruited for Exper-  
 776 iment 2 were caught by each combination of exclusion criteria. (The ungrammatical  
 777 sentences criterion was not used on its own to exclude participants in Experiment 2.)

	Gram.	Ungram.	Practice	Notes	Comprehension	Production
					35	36
			×		4	10
		×			11	10
778		×	×		2	1
	×				4	2
	×		×		1	4
	×	×			4	2
	×	×	×		7	2

## 779 C Analysis of all participants

780 In this appendix, we report the same analyses as in Sections 2.4 and 3.4 run on the data  
 781 from all originally-recruited participants, imposing none of the preregistered criteria for  
 782 exclusion.

### 783 C.1 Experiment 1

784 We recruited 183 participants in total for Experiment 1: 100 in the COMPREHENSION con-  
 785 dition and 83 in the PRODUCTION condition.

#### 786 C.1.1 Judgement

787 Figure 14 visualises the proportion of times each participant accepted each type of sen-  
 788 tence at test. The same model described above was fit to this data; its posterior estimates  
 789 are summarised in Table 7, and the conditional posterior distributions over the proba-  
 790 bility of accepting a sentence are shown in Figure 15.

791 Overall, we see a similar pattern to the original analysis: participants in both the  
 792 COMPREHENSION and the PRODUCTION condition accept the unsegmented sentences more  
 793 than the segmented sentences.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	0.27	0.16	-0.04	0.58
Condition	0.09	0.31	-0.54	0.69
Sentence type	2.07	0.32	1.45	2.71
Condition:Sent. type	0.23	0.32	-0.40	0.86

Table 7: The posterior probability distributions estimated by the model for the sentence acceptance data from all 183 participants recruited for Experiment 1. Values are on the log-odds scale.

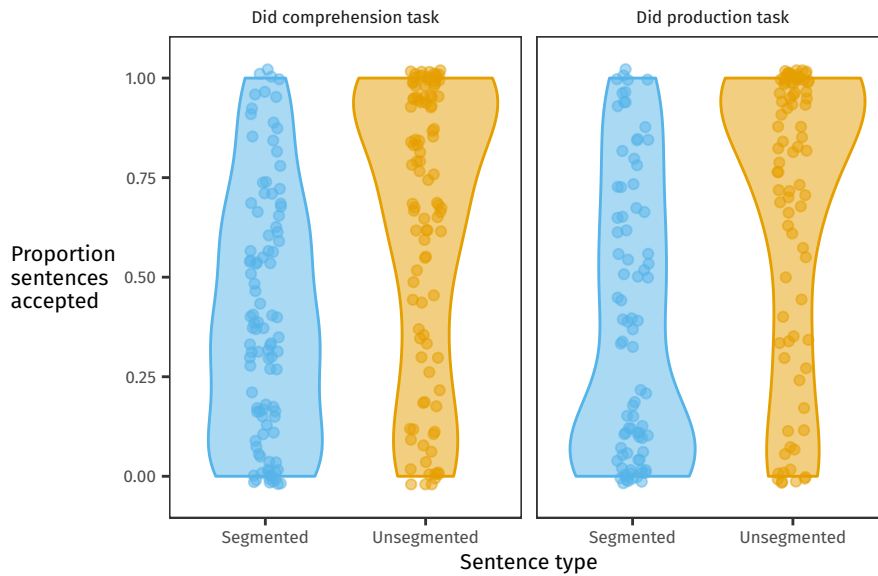


Figure 14: All 183 participants recruited for Experiment 1 accepted novel sentences that followed the unsegmented analysis more frequently than sentences that followed the segmented analysis, regardless of task. Each dot represents one participant’s proportion of accepted sentences of each type.

	Estimate	Est’d error	Lower 95% CrI	Upper 95% CrI
Intercept	0.46	0.15	0.17	0.76
Condition	0.11	0.31	-0.48	0.72

Table 8: The posterior probability distributions estimated by the model for all 183 participants’ held-out character naming data in Experiment 1. Values are on the log-odds scale.

## 794 C.1.2 Held-out character naming

795 Figure 16 illustrates the proportion of participants in each condition who named the  
 796 held-out character using the appropriate suffix—the one that doesn’t appear elsewhere  
 797 in the sentence. As in the original analysis, more than half of the participants in both  
 798 groups chose the word containing the appropriate, and the model estimates that both  
 799 groups have very similar probabilities of selecting the appropriate suffix (see the poste-  
 800 rior summaries in Table 8 and the conditional posterior distributions in Figure 17).

## 801 C.2 Experiment 2

802 We recruited 135 participants in total for Experiment 2: 68 in the COMPREHENSION con-  
 803 dition and 67 in the PRODUCTION condition.

### 804 C.2.1 Judgement

805 In Figure 18, we show the proportion of times each participant accepted each type of  
 806 sentence at test. We see the same pattern as in the original Experiment 2 data and in

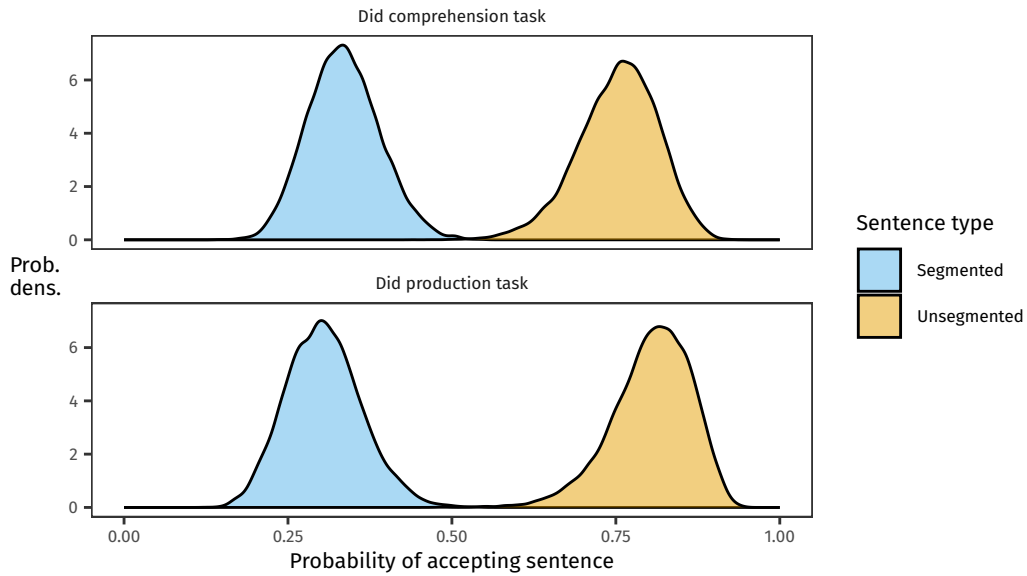


Figure 15: Conditional posterior probability distributions of the probability that all 183 participants recruited for Experiment 1 would accept a sentence. UNSEGMENTED sentences are more likely to be accepted than SEGMENTED sentences, regardless of whether participants did a comprehension or production task.

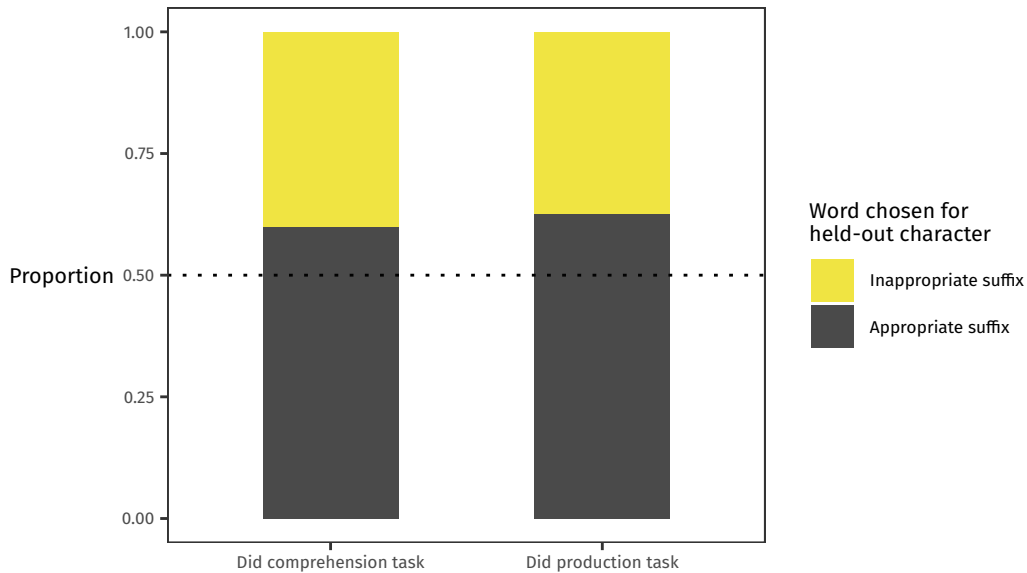


Figure 16: In the held-out character naming task of Experiment 1, more than half of all 183 participants selected the word with the appropriate suffix.

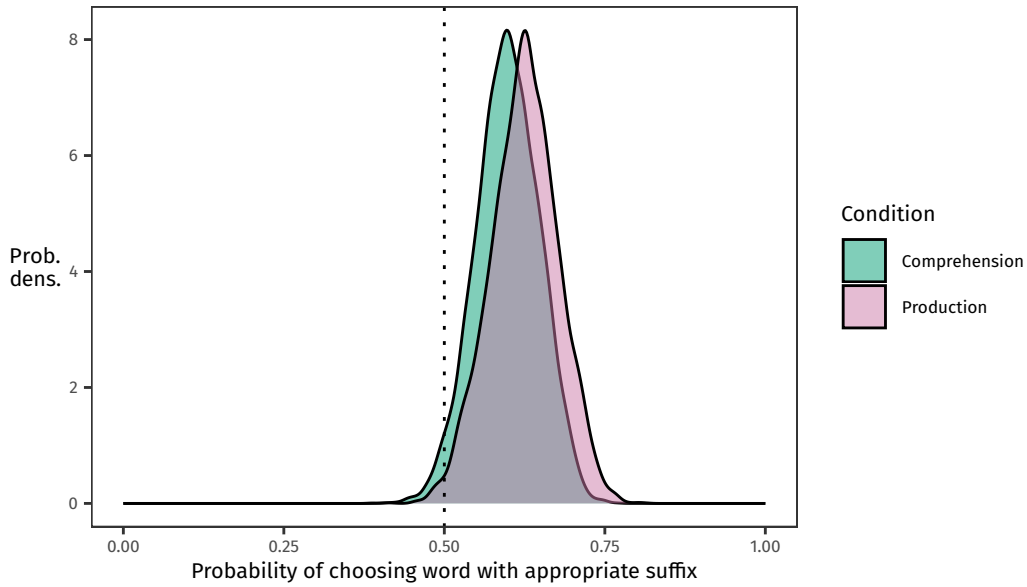


Figure 17: Conditional posterior probability distributions over the probability of selecting a word that contains the appropriate suffix in Experiment 1, shown for all 183 originally recruited participants.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	0.15	0.14	-0.13	0.43
Condition	0.32	0.27	-0.21	0.87
Sentence type	2.51	0.41	1.72	3.33
Condition:Sent. type	0.19	0.41	-0.61	0.98

Table 9: The posterior probability distributions estimated by the model for the sentence acceptance data from all 135 participants recruited for Experiment 2. Values are on the log-odds scale.

807 the data of all 183 participants from Experiment 1: participants prefer the unsegmented  
 808 sentences over the segmented ones, regardless of task. Table 9 summarises the posterior  
 809 distributions estimated by the same model as above, and Figure 19 shows the conditional  
 810 posterior distributions.

### 811 C.2.2 Held-out character naming

812 Figure 20 shows that, like the original analysis, around three-quarters of participants  
 813 in each condition named the held-out character using the appropriate suffix. Table 10  
 814 summarises the posteriors estimated by the same model described above, and Figure 21  
 815 shows the conditional posterior distributions.



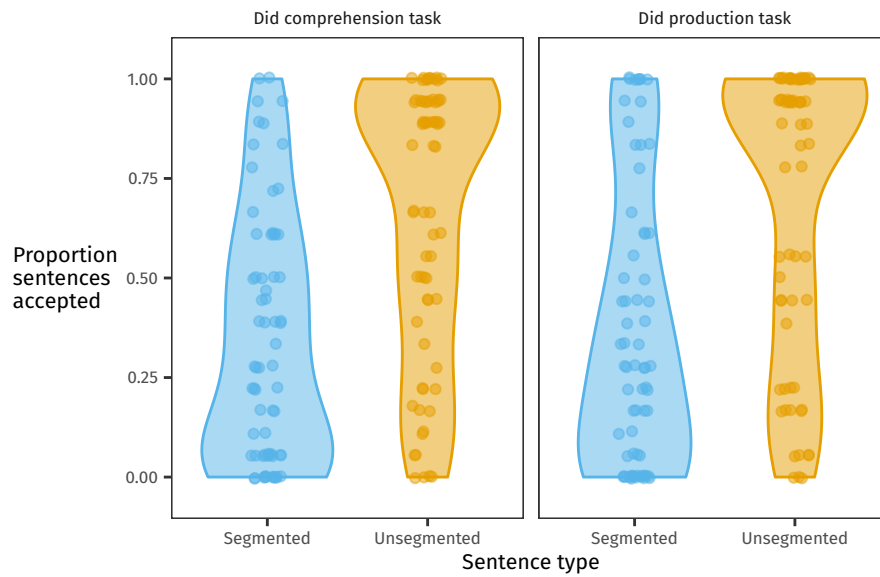


Figure 18: All 135 participants recruited for Experiment 2 accepted novel sentences that followed the unsegmented analysis more frequently than sentences that followed the segmented analysis, regardless of task. Each dot represents one participant's proportion of accepted sentences of each type.

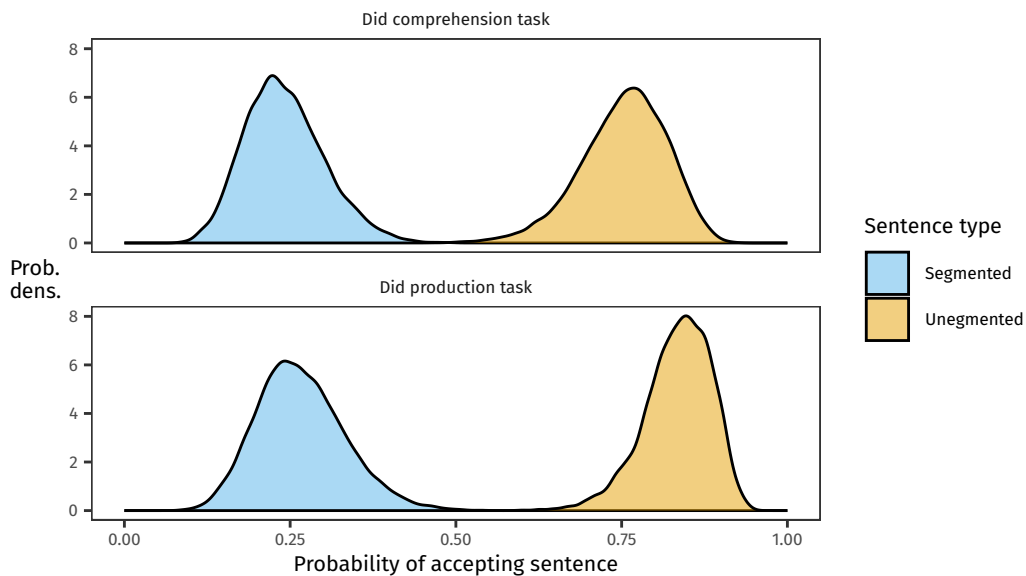


Figure 19: Conditional posterior probability distributions of the probability that all 135 participants recruited for Experiment 2 would accept a sentence. UNSEGMENTED sentences are more likely to be accepted than SEGMENTED sentences, regardless of whether participants did a comprehension or production task.

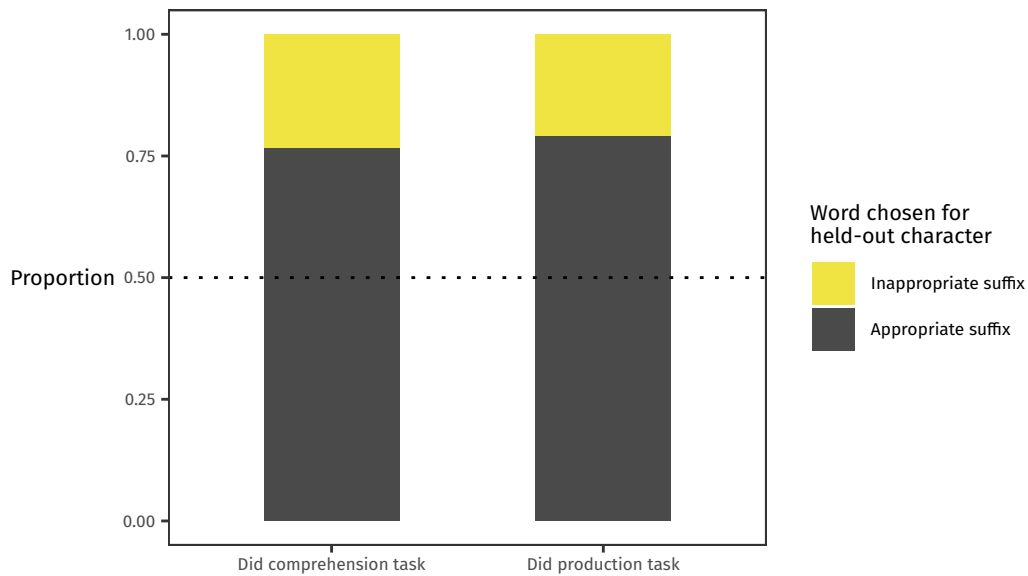


Figure 20: In the held-out character naming task of Experiment 2, at least three-quarters of all 135 participants selected the word with the appropriate suffix.

	Estimate	Est'd error	Lower 95% CrI	Upper 95% CrI
Intercept	1.26	0.21	0.86	1.67
Condition	0.14	0.41	-0.65	0.95

Table 10: The posterior probability distributions estimated by the model for all 135 participants' held-out character naming data in Experiment 2. Values are on the log-odds scale.

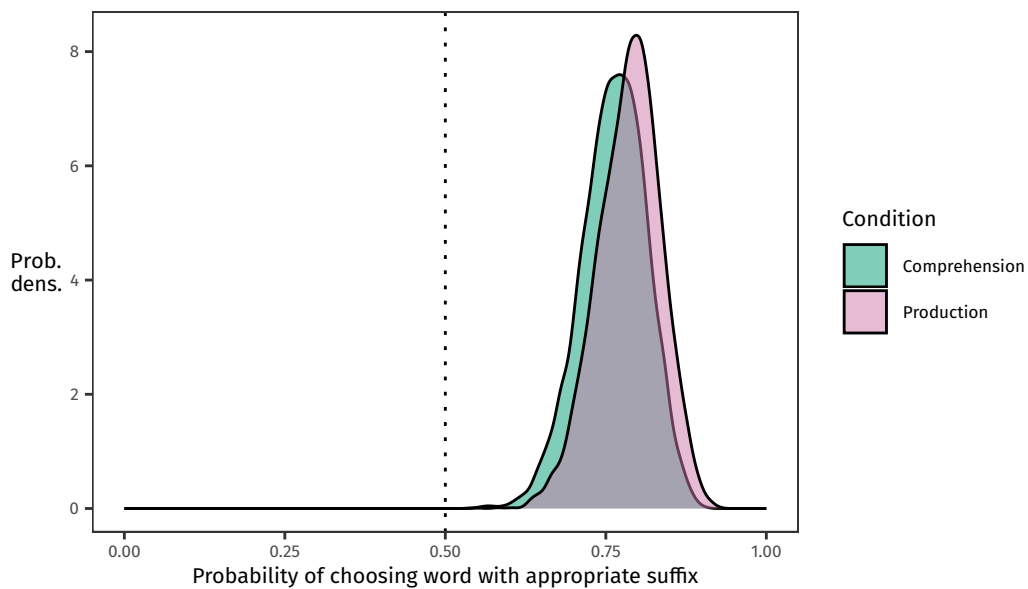


Figure 21: Conditional posterior probability distributions over the probability of selecting a word that contains the appropriate suffix in Experiment 2, shown for all 135 originally recruited participants.

## 816 **D Bayesian model specifications**

### 817 **D.1 Judgement**

```
818 brm(  
819   sentence_accepted ~ sent + cond + sentcond + (sent | ppt_id),  
820   family = bernoulli(),  
821   prior = c(  
822     prior(normal(0, 1.5), class = Intercept),  
823     prior(normal(0, 2), class = b),  
824     prior(normal(0, 5), class = sd, coef = Intercept, group = ppt_id),  
825     prior(normal(0, 5), class = sd, coef = sent, group = ppt_id),  
826     prior(lkj(2), class = cor, group = ppt_id)  
827   )  
828 )
```

### 829 **D.2 Held-out character naming**

```
830 brm(  
831   match ~ cond,  
832   family = bernoulli(),  
833   prior = c(  
834     prior(normal(0, 1.5), class = Intercept),  
835     prior(normal(0, 2), class = b)  
836   )  
837 )
```

838 **E Exploratory analysis: Participants who know case**  
 839 **marking languages**

840 In the post-experiment debrief questionnaire, we asked participants if they knew or  
 841 understood any other languages beyond English. If they self-reported knowing a case  
 842 marking language, we placed them into a separate group from the participants who did  
 843 not. Fifteen participants out of 80 reported that they know or understand the following  
 844 case marking languages: Arabic, Czech, German, Latin, Polish, Romanian, Slav, Somali,  
 845 Tunisian, Turkish, and Urdu.

846 **E.1 Judgement**

847 Figure 22 visualises the proportion of sentence acceptance judgements for each partici-  
 848 pant, split by condition and further by whether each participant knows a case marking  
 849 language.

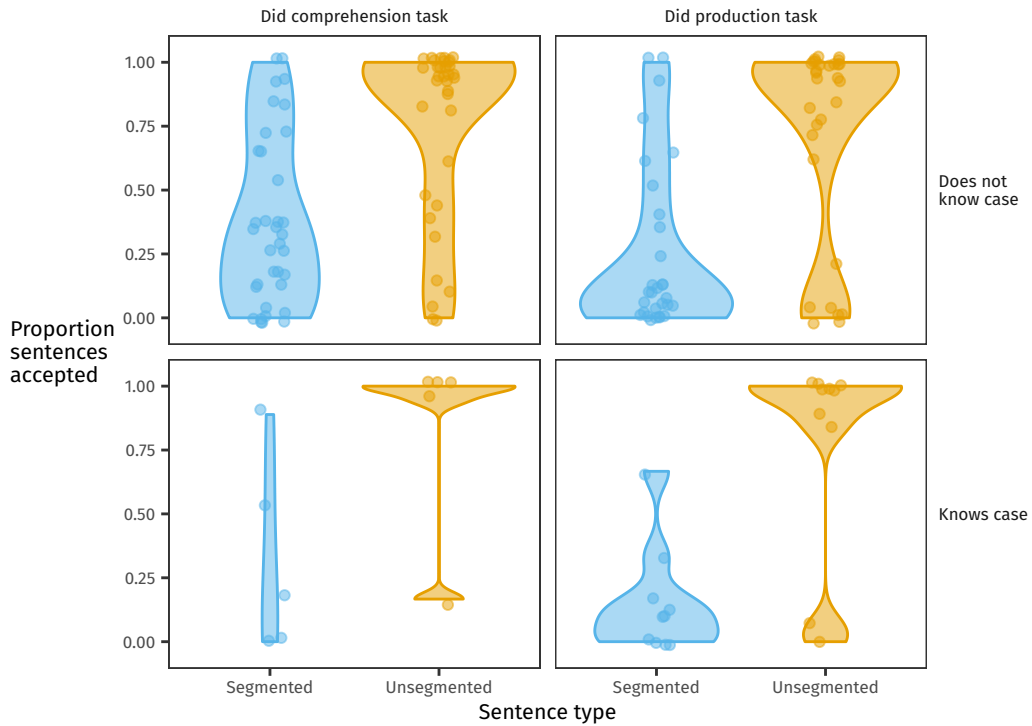


Figure 22: Participants in Experiment 1 who self-reported knowing a case marking language show a similar pattern of sentence acceptance to participants who do not know a language with case marking.

850 We fit the same Bayesian model as described in Section 2.4 to this data, adding in  
 851 an additional sum-coded predictor for knowledge of case ( $-0.5$  when the participant  
 852 does not know a case marking language,  $+0.5$  when they do), and all two- and three-  
 853 way interactions with the predictors sentence type and condition (scaled to  $\pm 0.5$ ). Table  
 854 11 summarises the posterior distributions of the population-level effects estimated by  
 855 the model. In short, the previously-estimated effects remain qualitatively the same, and

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.53	0.35	-0.15	1.25
Condition	-0.76	0.66	-2.05	0.52
Sentence type	4.36	0.86	2.68	6.03
Case	-0.03	0.69	-1.36	1.37
Condition:Sent. type	0.52	0.84	-1.11	2.18
Condition:Case	0.09	0.67	-1.21	1.42
Sent. type:Case	0.44	0.83	-1.21	2.04
Cond.:Sent. type:Case	0.28	0.85	-1.38	1.93

Table 11: Posterior distributions estimated by a model predicting sentence acceptance by condition, sentence type, and knowledge of a case marking language, and all interactions between them.

	Estimate	Est.Error	Q2.5	Q97.5
Intercept	0.56	0.32	-0.06	1.19
Condition	0.10	0.62	-1.11	1.27
Case	-0.59	0.62	-1.82	0.65
Condition:Case	-0.56	0.62	-1.78	0.64

Table 12: Posterior distributions estimated by a model predicting appropriate suffix choice by condition, knowledge of a case marking language, and their interaction.

856 the model indicates great uncertainty about any association between prior knowledge  
857 of case marking languages and acceptance of sentences formed using the segmented  
858 analysis.

## 859 E.2 Held-out character naming

860 Figure 23 illustrates that the 15 participants who know a case marking language select  
861 the word with the appropriate suffix less often than the larger group of 65 participants  
862 who do not know case. However, we fit a model estimating appropriate suffix choice  
863 as a function of condition, knowledge of case, and their interaction (scaled to  $\pm 0.5$ ), and  
864 the posterior distribution estimates in Table 12 indicate that we cannot be certain about  
865 any differences between participant groups.

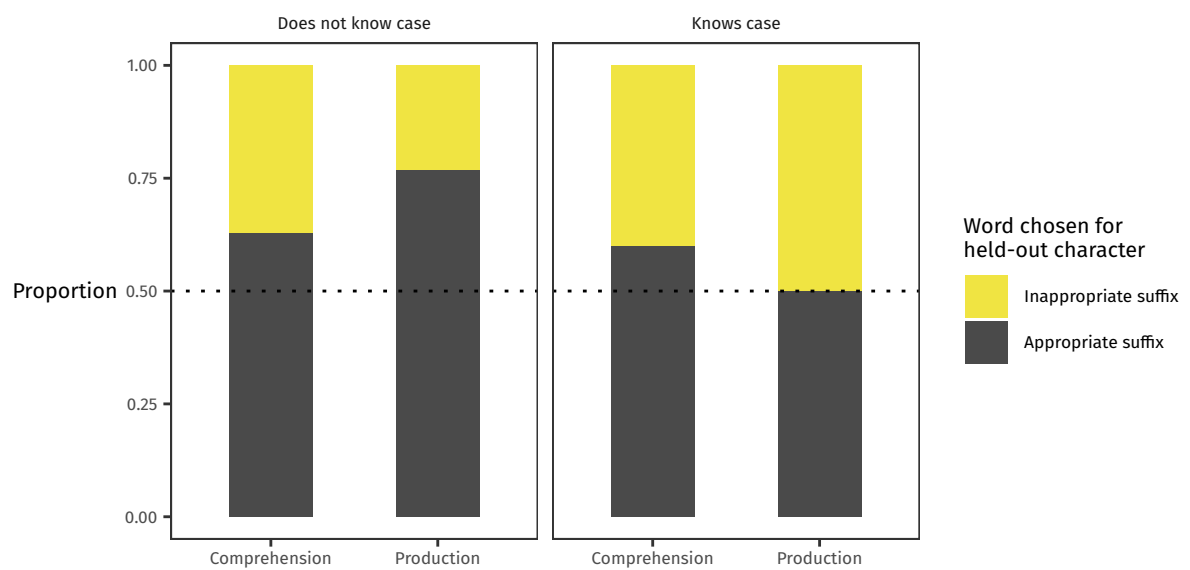


Figure 23: The 15 participants in Experiment 1 who know a case marking language give overall less appropriate responses to the held-out character naming task, with production participants selecting the appropriate suffix less than participants in the comprehension group.