

1

2

3

Middle ratings rise regardless of grammatical construction:

4

Testing syntactic variability in a repeated exposure paradigm

5

6 J.M.M. Brown¹,

7 Gisbert Fanselow^{1,2},

8 Rebecca Hall³,

9 Reinhold Kliegl^{1,4}

10

11

12 ¹ SFB 1287 Limits of Variability in Language, University of Potsdam, Potsdam, Germany

13 ² Department of Linguistics, University of Potsdam, Potsdam, Germany

14 ³ Department of Psychology, University of Potsdam, Potsdam, Germany

15 ⁴ Department of Sport Sciences, University of Potsdam, Potsdam, Germany

16

17

18 * Corresponding author:

19 E-Mail: jmmbrown@cantab.net (JMMB)

20

21 **Abstract**

22 People perceive sentences more favourably after hearing or reading them many times. A
23 prominent approach in linguistic theory claims that these types of exposure effects (satiation
24 effects) show direct evidence of a generative approach to linguistic knowledge: only some
25 sentences improve under repeated exposure, and which sentences do improve can be
26 predicted by a model of linguistic competence that yields natural syntactic classes. However,
27 replications of the original findings have been inconsistent, and it remains unclear whether
28 satiation effects can be reliably induced in an experimental setting at all. Here we report four
29 new satiation effects in wh-question constructions across two languages, German and English.
30 First, the effects pertain to zone of well-formedness rather than syntactic class: all
31 intermediate ratings, including calibrated fillers, increase at the beginning of the experimental
32 session regardless of syntactic construction. Second, ratings asymptote below maximum
33 acceptability, that is they do not satiate. Third, these effects are consistent across judgements
34 of superiority effects in English and German. Fourth, wh-question constructions appear to
35 show similar profiles in English and German, despite these languages being traditionally
36 considered to differ strongly in whether they show effects on movement: in both languages
37 violations of the superiority condition can be modulated to a similar degree by manipulating
38 subject-object initiality and animacy congruency of the wh-phrase. We improve on classic
39 satiation methods by distinguishing between the two crucial tests, that is whether exposure
40 selectively targets certain grammatical constructions or whether there is a general repeated
41 exposure effect. We conclude that exposure effects can be reliably induced in rating
42 experiments but exposure does not appear to selectively target certain grammatical
43 constructions. Instead, they appear to be a phenomenon of intermediate gradient judgements.
44

45 **Introduction**

46 Speaker judgments about the well-formedness of sentences form an important source of
47 evidence for evaluating linguistic theories. For more than two decades, linguists have
48 increasingly used experiments to collect acceptability judgments in a systematic way (e.g.,
49 Schütze [1], Cowart [2]). Measuring acceptability in a systematic way has demonstrated that
50 judgements of certain constructions change when this construction is repeatedly presented to
51 the same participant. Such changes are documented as early as Nagata [3] and Carroll [4].
52 The first major treatment in linguistics is Snyder [5], where he argues that constructions that
53 improve under repeated exposure form a natural syntactic class (these constructions "sate" as
54 he calls it, attributing the introduction of the term into syntax to Karin Stromswold). In
55 other words, independently motivated linguistic and psycholinguistic properties determine
56 whether a construction can or cannot sate. Therefore, the availability of satiation for a given
57 construction could be taken as an indication that this construction has a particular linguistic or
58 psycholinguistic property. Which property makes sentences prone to satiation remains
59 unclear. Explanations in the literature range from structural properties within generative
60 syntax (Snyder [5]: 579 classical "subjacency" effects in the sense of Chomsky [6]) to
61 sentence processing difficulty (as in, e.g., Chaves & Dery [7], [8]) to issues of comprehension
62 fluency (e.g., Zervakis & Mazuka [9]).

63 The outcome of the three experiments reported here suggests a different perspective (first
64 taken in a different form by Nagata [3]): satiation does not discriminate between construction
65 types on the basis of some abstract property. Rather, satiation indiscriminately affects all
66 types of sentences within a certain zone of judgments. "Zone of judgements" refers to an
67 interval in the judgment space between "fully acceptable" and "fully unacceptable" such that
68 satiation affects all structures rated within the interval, and none of the structures rated outside
69 the interval. It is plausible to relate this zone to the phenomenon of "intermediate"
70 acceptability, but it remains difficult to anchor this zone with precise acceptability values on,
71 for example, an n-point scale, because specific grades of acceptability are co-determined by

72 various orthogonal design factors in the experiments (such as the nature and quality of the
73 fillers and the target-filler ratio or the range of the n-point scale). Nevertheless, we conjecture
74 that this zone excludes the upper and lower 25% of the judgment space. In any event, our
75 empirical claim does not so much focus on specific values, but rather on the pervasive nature
76 of satiation within such a domain.

77 We draw our conclusions with reference to ratings of a set of standardized filler groups based
78 on those developed by Featherston and colleagues (Featherston [10], Gebrich et al. [11]) for
79 the purpose of calibrating syntactic judgments: each group comprises a set of constructions
80 that are diverse syntactically but share a specific value in the judgment space. All filler groups
81 turn out to satiate when their judgment value falls between the highest and the lowest
82 condition of the core experiment affected by syntactic satiation in our experiments.

83 In the remainder of the introduction, we review some of the literature to highlight the
84 diversity of approaches, both in terms of syntactic manipulations and experimental designs.
85 This review yields a somewhat incoherent and sometimes contradictory profile of results in
86 this field of research. Against this background, we provide the reasons motivating our use of
87 superiority violations and their examination in German and English and the role of fillers in
88 this context.

89 The literature on syntactic satiation effects is already quite large, and has been summarised
90 and discussed in detail in Snyder [12]. Experiments have been quite diverse methodologically,
91 and differ in how they identify syntactic satiation effects. In line with this diversity, the results
92 are quite variable – despite the fact that much of the literature focuses on the same set of
93 “classical” constructions first investigated for syntactic satiation by Snyder ([5]: 576). These
94 constructions are illustrated in (1). in which a gap created by moving some phrase is marked
95 with (t) for ‘trace’. In (1a), the argument *who* moves from within a want-for clause to the
96 front of the sentence. In (1b), *who* is extracted from within a clause introduced by *whether*
97 (wh-island). In (1c), *who* is moved to the front of the sentence across *that* (*that-trace-effect*).
98 In (1d), *what* is extracted from within a complex subject (extraction from a subject island). In

99 (1e), *who* is extracted from within a complex noun phrase, violating the Complex Noun
 100 Phrase Constraint. In (1f), *who* is extracted from within an adjunct (movement out of an
 101 adjunct island), and in (1g), *how many* is extracted from within the noun phrase *how many*
 102 *books* without pied-piping the rest of its constituent (violation of the Left Branch Condition).

103 (1) Classic satiation constructions

- 104 a. *Want-for*: *Who does John want for Mary to meet (t)?
 105 b. *Whether-island*: *Who does John wonder whether Mary likes (t)?
 106 c. *That-trace*: *Who does Mary think that (t) likes John?
 107 d. *Subject island*: *What does John know that a bottle of (t) fell on the floor?
 108 e. *Complex NP island*: *Who does Mary believe the claim that John likes (t)?
 109 f. *Adjunct island*: *Who did John talk with Mary after seeing (t)?
 110 g. *Left branch*: *How many did John buy (t) books?

111 [Snyder [5]: (2a), 576]

112 Relatively few studies investigated other constructions of English (e.g., Sprouse [13], Goodall
 113 [14], Zervakis & Mazuka [9]). In addition, there is some limited amount of work on different
 114 languages. Goodall [14] reports satiation effects for some constructions (e.g., questions
 115 without subject verb inversion) in Spanish but no such effects for, for example, Spanish
 116 counterparts of (1d, e, f). In Danish, satiation effects are reported for sentences with long
 117 movement, and no satiation for sentences with short movement or no movement at all or
 118 sentences with doubly-filled specifiers (Christensen, Kizach & Nyvad [15]).
 119 Table 1 summarises the outcome of a sample of satiation studies with (1) or a subset thereof
 120 as a point of reference.

121 **Table 1. Summary of satiation effects by previous study and construction.**

	Want-for	Whether-island	That-trace	Subject island	Complex NP island	Adjunct island	Left branch condition
Braze (2002) [16]		yes				no	
Chaves and Dery (2014) [7]: exp. 1				yes			
Chaves and Dery (2014) [7]: exp. 2				yes			
Chaves and Dery (2019) [8]: exp. 1				yes (subject gap) no (object gap)			
Chaves and Dery (2019) [8]: exp. 3				yes (parasitic gap) yes (object gap) no (non-parasitic gap)			
Crawford (2011) [17]		yes		no		no	
Francom (2009) [18]: exp. 1	yes	yes	no	yes	no	no	no
Francom (2009) [18]: exp. 2			no	yes	no	no	no
Goodall (2011) [14]			no	no	yes	nos	no
Hiramatsu (2000) [19]	yes	yes	yes	yes	no	no	no
Snyder (2000) [5]	no	yes	no	no	yes	no	no
Sprouse (2007a) [15]		no		no	no	no	
Sprouse (2007b) [16]: exp. 1	no	no	no	no	no	no	no
Sprouse (2007b) [16]: exp. 2	no	no	no	no	no	no	no
Sprouse (2007b) [16]: exp. 3	no	no	no	no	no	no	no
Sprouse (2009) [17]: exp. 1-5		no		no	no	no	
Sprouse (2009) [17]: exp. 6		no		no	no	no	no
Sprouse (2009) [17]: exp. 7		no			no	no	

123 The main observation emerging from Table 1 is how mixed previous results have been. The
124 only consistent results are the negative ones for left branch condition violations (1g). Some
125 studies find satiation effects for certain of the constructions but others do not. Probably, this
126 variability is at least partially due to differences in method.

127 When we focus on the idea that syntactic satiation is related to a “zone of acceptability”, we
128 again observe mixed results in the literature. Christensen, Kizach, and Nyvad [15] found a
129 positive correlation of order of presentation and judgments in two experiments focusing on
130 Danish wh-movement that are compatible with our hypothesis. In their Experiment 1, all
131 constructions of intermediate acceptability (2.43-3.66 on a 5-point scale) showed the satiation
132 effect but none of the sentences with more extreme values. Experiment 2 sheds a slightly
133 different light on the hypothesis. It was partially identical to Experiment 1, but (among other
134 changes) all structures (e.g., short movement) with an acceptability value above the saturation
135 range of Experiment 1 were removed. Again, a continuous segment of the rating scale (3.86 -
136 4.36) underwent satiation, but the high upper boundary (4.36 on a 5-point scale) does not so
137 much reflect satiation for items with near-perfect acceptability as the absence of more well-
138 formed conditions in the experimental material (compared to Experiment 1, the rating of long
139 movement had gone up by 0.7 points). Thus, the comparison of the two experiments
140 illustrates the difficulty, already mentioned above, of defining the zone of “intermediate
141 acceptability” purely numerically in terms of values on an n-point scale. Moreover, the
142 experiments had their empirical focus on contrasts in the length of movement, so it is not
143 obvious that their interpretation can be generalized beyond this domain.

144 Zervakis and Mazuka [9] tested a more varied set of constructions for satiation effects with a
145 between-subject design. The experimental group rated 5 blocks of 100 experimental sentences
146 on a 7-point scale. The blocks occurred in different orders, such that each of the blocks was
147 rated as the last one by the same number of participants. The control group began by rating
148 400 control sentences and then judged the acceptability of one of the five experimental

149 blocks. A comparison of the experimental with the control group revealed a significant
150 increase in acceptability for all constructions rated higher than 2.67 by the control group up to
151 a garden path construction rated at 4.56. However, a condition called “grammatical binding”
152 falling well into this interval (3.45) failed to satiate, and the simple filler sentences also went
153 up in their ratings although their initial value was at 6.03 (Zervakis & Mazuka [9]: 513).

154 There are also proposals (Snyder [5], Goodall [14]) explicitly arguing against the idea that
155 satiability is a by-product of intermediate acceptability. Snyder [5] identified whether-island
156 (1b) and the Complex Noun Phrase Constraint. (1e) violations as satiating constructions (but
157 not, e.g., *that*-trace violations 1c) by comparing the number of participants giving more “yes”
158 responses to such constructions in the second half of a categorial acceptability task than in the
159 first half, to those showing the opposite dynamics of judgment. In a second experiment,
160 different participants were asked to rate the same material on a 5-point scale. Two
161 constructions did not satiate (*that*-trace (1c) and subject island (1d)) and these violations were
162 rated in between two conditions that did satiate (whether-island (1b) and Complex Noun
163 Phrase Constraint (1e) violations).

164 Snyder’s basic finding was replicated in experiment 1 of Francom [18], where whether-island,
165 Complex Noun Phrase Constraint (1e), *that*-trace (1c) and subject island (1d) violations
166 turned out to have basically the same level of acceptability (between 30 and 40% positive
167 replies in a categorial judgment task, Francom [18]: 35), while only two of them satiated
168 (whether- and subject islands). However, when applying a more balanced design, combined
169 with an increase in the number of items per condition, Francom [18] not only observed
170 changes in the acceptability rankings of the constructions (with the ratings of Complex Noun
171 Phrase Constraint violations falling below those of *that*-trace and subject island violations in
172 comparison to what experiment 1 had revealed), but also found satiation effects restricted to
173 subject island violations (Francom [18]:56-58).

174

175 Goodall [14] provides another example for the sensitivity of acceptability ratings to local
176 context. For English, he reports a considerable increase in the acceptability of Complex Noun
177 Phrase Constraint violations. However, despite comparable acceptability judgements in the
178 first block, this increase was not seen for that-trace violations in the final two of five
179 experimental blocks. Clearly this result is not in agreement with the hypothesis that satiation
180 is a function of the degree of acceptability. The grammar-based account has problems of its
181 own. In Francom ([18], Experiment 2) violations of the Complex Noun Phrase Constraint did
182 not satiate but subject island violations did. The variability between experiments visible calls
183 for an investigation of satiation with a more powerful design.

184 The variability of methods extends also to the intensity of the repeated exposure of the
185 participants to some construction. Zervakis and Mazuka [9] used a block design with a large
186 number of items in a between-subject design; they found quite pervasive satiation effects. We
187 hypothesize that satiation effects are of moderate size and arise only after a considerable
188 amount of exposure.

189 In summary, the review of a select set of studies above highlights the need for employing a
190 standardized experimental paradigm with the ideal of an equal number of items for all levels
191 of gradient acceptability in the material and a maximum of counterbalancing their order of
192 presentation across the experiment. The experimental design must also aim for large statistical
193 power with a large number of subjects and items, within-subject/within-item experimental
194 manipulations, and a sufficient number of blocks of trials. Such a design is expected to yield
195 evidence on the shape of the satiation function: does it arise quickly and asymptote at or
196 considerably below close-to-perfect acceptability?

197 The final topic of the introduction is on the languages and the constructions used in our
198 experiments. As in Goodall [14], we decided to look for satiation effects for a comparable or
199 even identical set of constructions in two languages, rather than only one. If it is the
200 intermediate nature of the judgment rather than specifics of grammatical properties that are

201 responsible satiation effects, the same or a similar satiation functions should emerge for
 202 acceptability ratings in both languages. Conversely, we expect qualitative differences if
 203 grammatical properties matter, as demonstrated by Goodall’s discussion of why Complex
 204 Noun Phrase Constraint violations satiate in English but not in Spanish due to their
 205 grammatical differences, in spite of comparable initial acceptability.

206 There is a tension between the goals of using close-to-identical sentence material in the
 207 experiments (changes ideally being restricted to the use of different lexicalization) and the
 208 goal of reducing the impact of construction type by making level of acceptability the foremost
 209 criterion for deciding on the material. Our experiments aimed for a compromise on these
 210 incompatible demands by basing the main experiment of targets on constructional similarity
 211 and by also including a systematic set of fillers representing six levels of acceptability in both
 212 language by different constructions. English and German are an ideal pair of languages for
 213 our purpose. While they are closely related, and share many constructions in terms of surface
 214 appearance, the grammatical analysis of these constructions for the target sentences may be
 215 quite different as shown, for example, by Haider [22] and many others. This applies in
 216 particular to the grammar of multiple questions,

- 217 (2) a. John knows who saw what.
 218 b. *John knows what who saw ~~what~~.

219 Multiple questions are well-formed in English when the order of the wh-phrases corresponds
 220 to their normal linearization in declarative sentences, as in (2a), in which the wh-subject
 221 precedes the wh-object. The placement of a wh-object in front of a wh-subject usually leads to
 222 a remarkable drop in acceptability (2b), called “superiority effect”, as first observed by Kuno
 223 and Robinson [22]. While there is some disagreement in the literature, the standard hypothesis
 224 in generative syntax is that the unacceptability of (2b) reflects the operation of some
 225 *grammatical* principle. which we will call “superiority condition” in this paper without
 226 committing ourselves to any of the proposals (see, e.g., Häussler et al., [24], and Fanselow,
 227 [25]) for a discussion.

228 The situation is different in German. Like in English, the acceptability of object initial
229 multiple questions was found to be reduced as compared to their subject-initial counterpart in
230 a number of experimental studies, yet syntactic reasons summarised in Haider [26] have led to
231 a widespread conviction that whatever kind of superiority effect one may observe in German,
232 it does not reflect the operation of a syntactic superiority condition but is due to a
233 “conspiracy” of a number of factors, some of which pertain to the realm of language
234 processing (see Häussler et al. [24], for experiments meant to support this perspective). Thus,
235 superiority effects of German and English may arise from different sources. Furthermore,
236 many experiments have revealed the gradient/intermediate nature of judgments related to the
237 superiority effect (see, again, Häussler et al., [24], for an overview). Multiple questions are
238 the only domain for which there is a sufficient number of experimental studies comparing
239 English and German to base the present study on (Featherston [10], Häussler et al. [24]).
240 Finally, the presence of satiation effects for multiple questions has been hypothesized in the
241 previous literature (Hofmeister et al. [28]), but systematic studies concerning satiation in
242 multiple questions have not been undertaken so far.

243 We have good reasons for the choice of the superiority effect, but there are also reasons for
244 deciding against using the “classical” construction set of Snyder [5]. First, there is no
245 counterpart of the *want-for* construction of English in German. Second, German shows
246 regional variation with respect to the status of moving elements out of a complement clause
247 (see, e.g., Fanselow & Weskott [29]). Long movement in the experimental items may thus
248 incur an unwanted impact of sociolinguistic evaluation of the appropriateness of using
249 regional variety examples on judgments of “formal” well-formedness that one would not see
250 in the English counterpart. Rather than trying various options for fixing this problem, we
251 decided to take recourse to short inner-clausal movement and work on acceptability
252 differences for which we are not aware of any dialectal or sociolinguistic variability.

253 For the second dimension of our experiment, that is using fillers representing different levels
254 of (gradient) acceptability, we used constructions identified using a norming study in
255 Featherston, [10], and Gebrich et al. [11]). Based on a set of carefully implemented
256 experiments, they identified five sets of structurally diverse sentences which occupy constant
257 relative points in the judgment space and can thus be used as “calibrators” anchoring these
258 points in the judgment space. Examples for the lower three levels of English are exemplified
259 below, showing that there is no (obvious) property they have in common but the level of
260 acceptability.

261 Level C:

262 Hannah hates but Linda loves eating popcorn in the cinema.

263 Most people like very much a cup of tea in the morning.

264 The striker must have fouled deliberately the goalkeeper.

265 Level D:

266 Who did she whisper that had unfairly condemned the prisoner?

267 The old fisherman took her pipe out of mouth and began story.

268 Which professor did you claim that the student really admires her?

269 Level E

270 Historians wondering what cause is disappear civilisation.

271 Old man he work garden grow many flower and vegetable.

272 Student must read much book for they become clever.

273

274 There is a necessary between-language variability for the fillers because German material is,
275 of course, not based on structural similarity but on level of acceptability. Thus, German level
276 C contains (among other material) unusual binding constellations. The acceptability of the
277 example for level D is problematic because of the sentence does not respect a leftward
278 placement rule for unstressed pronouns. Two examples are given below. Each level contains
279 three constructions as for English.

280 Level C

281 Ich habe dem Kunden sich selbst im Spiegel gezeigt.

282 I have the client himself.REFL himself in.the mirror shown

283 *I showed the client himself in the mirror.*

284 Level D

285 Der Komponist hat dem neuen italienischen Tenor es zugemutet.

286 the composer has the new Italian tenor it expected.of

287 *The composer expected it of the new Italian tenor.*

288 Between-language differences in acceptability ratings for the different levels of filler

289 sentences may pose problems of interpretation for results, but if fillers show a similar

290 behaviour in the two languages, this is most likely due to their position in the judgment space.

291 **Experiment 1: rating of German sentences**

292 Experiment 1 tests the effect of repeated exposure in German indirect multiple wh-questions

293 in subordinate clauses. Targets crossed whether they were subject-initial or object-initial (the

294 latter case tests for possible superiority effects) and whether there was an animacy

295 congruency between subject and object wh-words, yielding a 2 x 2 design. Thus, each target

296 sentence was available in four conditions as shown in (3).

297 (3) German target conditions

298 a. Condition 1: subject-initial; matching animacy (*wer-wen*); well-formed in German

299 Keinesfalls wusste die Haushälterin genau, **wer** bei der Gartenfeier **wen**

300 Certainly.not knew the housekeeper exactly who by the garden.party who

301 ständig angesehen hat.

302 continuously looked.at had

303 *The housekeeper certainly did not know exactly **who** had kept on looking at **who(m)** at*

304 *the garden party.*

305 b. Condition 2: subject-initial; mismatching animacy (*wer-was*); well-formed in German

306 Keinesfalls wusste die Haushälterin genau, **wer** bei der Gartenfeier **was**
 307 Certainly.not knew the housekeeper exactly who by the garden.party what
 308 ständig angesehen hat.
 309 continuously looked.at had
 310 *The housekeeper certainly did not know exactly **who** had kept on looking at **what** at*
 311 *the garden party.*

312 c. Condition 3: object initial; matching animacy (*wen-wer*); not fully well-formed in
 313 German

314 Keinesfalls wusste die Haushälterin genau, **wen** bei der Gartenfeier **wer**
 315 Certainly not knew the housekeeper exactly who by the garden.party who
 316 ständig angesehen hat.
 317 continuously looked.at had
 318 *The housekeeper certainly did not know exactly **who** had kept on looking at **who(m)***
 319 *the garden party.*

320 d. Condition 4: object initial; mismatching animacy (*was-wer*); well-formed in German

321 Keinesfalls wusste die Haushälterin genau, **was** bei der Gartenfeier **wer**
 322 Certainly.not knew the housekeeper exactly what by the garden.party who
 323 ständig angesehen hat.
 324 continuously looked.at had
 325 *The housekeeper certainly did not know exactly **who** had kept on looking at **what** at*
 326 *the garden party.*

327 The first factor manipulates the superiority/crossing movement variable: conditions 1 and 2
 328 contain two wh-phrases and the subject precedes the object. These conditions therefore
 329 illustrate instances where a superiority effect cannot possibly arise, and act as controls against
 330 which to compare potential superiority effect in the object-initial targets. Conditions 3 and 4
 331 contain the crucial manipulation: In these conditions, the object wh-phrase is placed in front
 332 of the subject wh-phrase, meaning that a superiority effect could arise.

333 The second factor manipulates an extra-grammatical feature, namely the animacy of the wh-
334 phrases. Thus, subject and object match in animacy in two sentence types, namely (3a) and
335 (3c), while the other two conditions, namely (3b) and (3d) have subjects and objects that do
336 not match in animacy. We start here from the assumption that differences in subject and
337 object animacy increase the well-formedness of crossing movement, see Fanselow et al. [30]
338 and Häussler et al. [24].

339 As described above, ratings of target sentences are to be analysed not only with respect to the
340 2 x 2 experimental manipulations (and how they change across the blocks of the experiment),
341 but also in the context of a calibrated set of six types of filler sentences [3]. The first five
342 gradations of fillers varied from (A) "almost not well-formed" to (E) "completely well-
343 formed". We included also a new sixth level (F) "uninterpretable and unacceptable" as clearly
344 ungrammatical level. Examples are provided under Material below.

345 Finally, with respect to methodology, we go beyond past designs in this area with a within-
346 subject/within-item counterbalancing scheme built around six blocks of 72 items and a
347 multiple of 24 subjects. As far as we could determine, past research employed only the usual
348 counterbalancing measures for experimental conditions in the theoretical focus. Thus, applied
349 to our design, the four instances of each target sentences are presented equally often in the
350 experiment but such that every subject rates only one instance of each target sentence and
351 rates the same number of items in each of the four conditions. However, in a satiation
352 experiment this is not enough. As items vary in acceptability we must also ensure that they
353 appear equally often in each of the six blocks of the experiment while respecting the usual
354 constraints. In other words, counterbalancing is extended to a 2 x 2 x 6 scheme. The same
355 also holds for the six type of filler sentences. As type of filler varies within subject, but
356 between items, each filler is presented equally often in each block and rated once by each
357 subject; they also rate the six types of fillers equally often in each block. As another
358 innovative design feature we imposed the constraint that first-order transitions between the
359 four target conditions and the six filler levels occurred equally often (i.e., we used a Williams,

360 [31], design). This counterbalancing optimally controls for item differences in acceptability
361 and increases the signal-to-noise ratio for the detection of satiation effects.

362 **Method**

363 **Subjects**

364 A total of 55 students and employees of the University of Potsdam participated in the study.
365 Only the results of German native speakers were included in the evaluation; two participants
366 were excluded due to technical failures. Of the remaining 48 participants, 4 were male and 44
367 female and ranged between 19 and 40 years of age, with an average age of 24 years.
368 Recruitment was carried out via flyers distributed on campus, internet advertisements on
369 various platforms, and an e-mail distribution list from a university experimental laboratory.
370 Participation was voluntary and was remunerated with study credit or eight euros. All
371 participants were naive with regard to the questions and objectives of the study.

372 **Apparatus**

373 The experiment was conducted in the experimental psychology lab at the University of
374 Potsdam. Three soundproof computer booths were used. The experiment was implemented
375 using the Python software PsychoPy version 1.84.2 [32]. The stimulus material was presented
376 on a computer screen along with a scale labelled from 1 to 7. Underneath each of the numbers
377 was a short description of the degree of well-formedness: 1 “überhaupt nicht wohlgeformt”
378 (not at all well-formed), 2 “fast nicht wohlgeformt” (almost not well-formed), 3 “eher nicht
379 wohlgeformt, als wohlgeformt” (more not well-formed than well-formed), 4 “kann man nicht
380 zuordnen” (cannot be classed as well-formed or ill-formed), 5 “eher wohlgeformt, als nicht
381 wohlgeformt” (more well-formed than not well-formed), 6 “nahezu wohlgeformt” (mostly
382 well-formed), 7 “völlig wohlgeformt” (completely well-formed). Sentence ratings from “1” to
383 “7” were entered on a standard computer keyboard with a German keyboard layout.

384 **Material**

385 The material consisted of 120 targets and 252 fillers. that is the ratio of targets to fillers was
 386 roughly 1:2. Subjects rated targets in one of four versions (i.e., a total of 480 different
 387 sentences was constructed from the 120 targets). A complete list of targets and fillers is
 388 provided in the OSF Repository (<https://osf.io/ge2db/>).

389 The targets represent indirect multiple wh-questions in subordinate clauses. The construction
 390 principle yielding the four instances for each of the 120 targets, that is subject/object initial
 391 sentence (2) x animacy match/mismatch of subject and object wh-words (2), is illustrated in
 392 (3) above. In addition, as also shown in the example in (3), half of the target sentences started
 393 and the other half ended with an adverb in the main clause; in the subordinate clause there
 394 was an adverb after the second wh-word. The purpose of adverbs was to make the sentences
 395 sound as natural as possible.

396 The fillers involved six grammatical levels with 42 items each. For the first five gradations
 397 from (A) "interpretable and highly acceptable " to (E) "interpretable but less acceptable than
 398 (D)", we used calibration sets from [3] to create the fillers. That is, we started from single
 399 examples of each construction at each acceptability level (3 constructions x 5 levels = 15
 400 items) [3], and created 195 more items on the same templates. However, we removed the
 401 multiple question construction from level D for obvious reasons. We added a sixth level (F)
 402 "uninterpretable and unacceptable" with 42 items to provide a clearly ungrammatical level.

403 (4) German filler examples (based on the Featherston fillers)

404 1. Level A: Interpretable and highly acceptable

405 In der Mensa essen viele Studenten zu Mittag.

406 In the canteen eat.3PL many students to lunchtime

407 *Many students have lunch in the canteen.*

408 2. Level B: Interpretable but less acceptable than (A)

409 Der Kaiser hat dem Fürsten den Maler empfohlen.

410 The emperor has the prince the artist recommended

411 *The emperor recommended the artist to the prince.*

412 3. Level C: Interpretable but less acceptable than (B)

413 Ich habe dem Kunden sich selbst im Spiegel gezeigt.

414 I have the client himself.REFL himself in.the mirror shown

415 *I showed the client himself in the mirror.*

416 4. Level D: Interpretable but less acceptable than (C)

417 Der Komponist hat dem neuen italienischen Tenor es zugemutet.

418 the composer has the new Italian tenor it expected.of

419 *The composer expected it of the new Italian tenor.*

420 5. Level E: Interpretable but less acceptable than (D)

421 Der Waffenhändler glaubt er, dass den Politiker bestochen hat.

422 The arms.dealer believes he that the politician bribed has

423 *It is the arms dealer that he believes bribed the politician*

424 6. Level F: Uninterpretable and unacceptable

425 Die Tinte wurde für vergossen.

426 *the ink was for spilled*

427 Note that in (4), we follow the convention of providing literal translations for individual items
 428 in glosses, followed by the closest meaningful translation in English on a separate line (here
 429 in italics). Therefore, although the translations for filler levels (A) to (E) are all fully
 430 acceptable and well-formed in English, the German examples themselves from (B) through
 431 (E) are not fully acceptable: they decrease in well-formedness between the highest level (A)
 432 and the lowest level (F). The lowest filler (F) is intended to be uninterpretable, meaning that
 433 there is no meaningful way for all the words to be integrated into the interpretation of the
 434 whole sentence, and therefore no full translation of the whole sentence to English. For filler
 435 (F), we therefore provide only literal translations for individual items in the gloss, and leave
 436 the translation blank.

437 **Design and counterbalancing**

438 The counterbalancing scheme used for the experiment differs from past research and was
439 described at the end of the Introduction. Each subject rated 372 sentences (120 targets + 252
440 fillers) that were distributed across six blocks (i.e., 62 sentences per block; $2 \times 2 \times 5 = 20$
441 targets and $6 \times 7 = 42$ fillers). As already described above, each subject rated the 120 targets
442 in only one of the four conditions and five targets in each of the 2×2 conditions in each
443 block. Similarly, subjects rated the 252 fillers once and seven fillers of each of the six types in
444 each block. The counterbalancing scheme also ensured that all targets were rated equally
445 often in their four conditions in each block and that all fillers and filler levels occurred equally
446 often in each block.

447 Presentation order of targets and fillers also adhered to a Williams [31] design. A Williams
448 design is a type of Latin square design that controls for first-order carry-over effects, ensuring
449 that transitions between the four experimental conditions and the six types of fillers occur
450 equally often. Finally, the item sequence was subject to the following constraints: (a) no
451 immediate repetition of target sentences (i.e., target sentences were bracketed by at least one
452 filler sentence), (b) no immediate repetition of the same experimental condition, and (c) no
453 immediate repetition of the same type of filler sentence.

454 This counterbalancing scheme requires a multiple of 24 subjects. Twenty-four is a typical
455 sample size for psycholinguistic research, but to increase statistical power we recruited twice
456 the minimal number, that is a total of 48 subjects. Statistical power is also high because all
457 design factors (subject/object-initiality of wh-words, animacy-congruency of wh-words,
458 block, level of filler, target vs. filler) vary within subject; and only two of them, target vs.
459 filler and level of filler, are between-item factors.

460 **Procedure**

461 After providing informed consent and collection of demographic information, subjects were
462 instructed in the well-formedness rating procedure. Specifically, they were asked to judge the
463 well-formedness of the sentences according to spoken—not written—language and to use the

464 full spectrum of the seven-point scale. They practiced the procedure with nine examples
465 spanning the scale of well-formedness, including two of the uninterpretable fillers (filler level
466 F). Then, each subject worked through the six blocks of sentences. At the end of each block,
467 the word *pause* appeared on the screen with instructions to press a button when the participant
468 was ready to continue. Participants could pause in front of the computer at this point.
469 Durations of breaks between blocks was under the subjects' control. Most participants took
470 around an hour to complete the experiment.

471 **Statistical analysis**

472 We used the open source software R [33], especially packages *lme4* [34], *tidyverse* [35],
473 *cowplot* [36], *sjPlot* [37], and *broom.mixed* [38] for statistical analyses. Inferential statistics
474 for the analyses of ratings are based on two linear mixed model (LMMs), one only for ratings
475 of targets and one for a joint analysis of targets and fillers.

476 *LMM for rating of targets*

477 The LMM for targets included subject and target as crossed random factors and subject- vs.
478 object-initial word order (2), animacy congruency of wh-words (2), and block (6) as fixed
479 factors. The three factors were varied within-subject and within-items. Across blocks we
480 expected satiation of well-formedness ratings, that is ratings would increase and eventually
481 reach an asymptote. Therefore, a Helmert contrast [39] was specified for levels of block to
482 capture the point at which the rating did no longer change significantly. To this end the first
483 contrast tested block 1 against the average of blocks 2 to 6, the second contrast the second
484 block against the average of blocks 3 to 6, and so on.

485 Random effects associated with within-subject or within-item factors potentially give rise to
486 variance components (VCs) and correlation parameters (CPs) of the mean rating and of
487 experimental effects in the random-effect structure of the LMM. They need to be included to
488 guard against false positives, but many VCs and CPs are not supported by the data and, if
489 included, reduce statistical power. We selected an LMM following the strategy outlined in
490 [40]; see also [41]). The significance of fixed effects which are the focus here did not depend

491 on the specifics of the random-effect structure. Details of model selection are documented in
492 the analysis scripts in the OSF repository.

493 *LMM for targets in the context of fillers*

494 The second LMM included ratings of both target and filler sentences. We specified a block
495 (first block vs. average of blocks 2 to 6) x type of sentence (10) design; type of sentence
496 comprised four types of target and six types of filler sentences. Sentence types were assigned
497 to levels of a factor according to their overall mean rating of well-formedness. Then
498 sequential-difference contrasts across these levels were specified for the factor and
499 interactions between block and nine contrasts; five of these nine block x type-of-sentences
500 contrasts tested the null hypotheses that increase in well-formedness ratings from block 1 to
501 block 2 did not differ (C1 to C5 [39]. As the six types of filler sentences varied between
502 items, only two of the nine contrasts for sentence types and their interactions with block
503 yielded estimates for VCs and CPs; all nine sentence types varied within subjects. Model
504 selection, that is determination and inclusion of VCs and CPs followed the strategy outlined
505 in [40]. Due to the complexity of this LMM, we used the JuliaStats/MixedModels.jl package
506 [42] for model fitting and model selection. Details of model selection are documented in the
507 analysis scripts in the OSF repository.

508 Given the large number of subjects, items, and observations, the usual t-distribution
509 approximates the normal distribution. Therefore, we report test-statistics (estimate / standard
510 error) as z-values and interpret absolute values larger than 2 as significant.

511 **Results and Discussion**

512 **Targets**

513 The first analysis focuses on the targets involving multiple questions. The four conditions of
514 targets using constructions of (a) “wer-wen”, (b) “wer-was”, (c) “wen-wer”, and (d) “was-
515 wer” for the same sentence frames map onto a 2 x 2 design with the main effects subject-vs-
516 object initial targets (so; a+b-c-d) and animacy congruency of wh-words (an; a+c-b-d) as well
517 as the interaction between these two effects (a-b-c+d). In addition, this 2 x 2 design was

518 repeated across six blocks of trials. Fig 1a displays the change in ratings of well-formedness
519 for the four types of targets shown in red; Fig 1a also shows performance for the levels of
520 fillers, but we initially focus only on the targets. The LMM fixed-effect estimates and
521 statistics are given in a supplement (S1 Table).

522 **Fig 1. Well-formedness ratings (a) by block and (b) for first block and mean of later blocks for German**
523 **targets and fillers.**

524 As expected, the subject-initial targets were rated as more acceptable than object-initial
525 targets, yielding a significant main effect of order ($b=0.52$, $z=8.97$). Overall, there was no
526 significant effect of animacy congruency ($b=-0.05$, $z=-1.62$), but there was a significant
527 interaction with order ($b=0.13$, $z=5.00$): there was a very clear preference for the inanimate
528 “was-wer” construction over the animate “wen-wer” construction for object-initial targets,
529 whereas – at least in the first block-- the reverse preference held for subject-initial targets.
530 Indeed, the pattern of change across blocks followed a very simple structure: overall, there was
531 a significant increase in ratings from the first to the average of blocks 2 to 6 ($b=0.31$, $z=3.24$),
532 and this increase was different for the interaction of order and animacy ($b=-0.15$, $z=-2.60$): As
533 is clearly visible in Fig 1, ratings increased less strongly for “wer-wen” and more strongly for
534 “wer-was” constructions than the other three conditions. None of the other contrasts defined
535 for the change in ratings across blocks, nor – with one exception – none of the other interaction
536 terms involving these contrasts were significant (all z -values < 2.00). The exception is a marked
537 violation of parallelism between block 4 and 5: only the “was-wer” construction exhibited in
538 increase in rating ($b=0.10$, $z=2.34$). We consider this interaction as spurious.

539 There are three main results. First, there is clear evidence for large differences in the
540 judgement of well-formedness of multiple questions with a clear preference for subject-initial
541 than object-initial targets, in line with previous experiments [24, 27, 30] and at the same time
542 there is no reliable evidence that these preferences changed much after one block of trials.

543 Second, for object-initial targets, mismatched animacy is preferred to matched animacy. The
544 condition with mismatched animacy is almost on par with subject-initial conditions, in line
545 with the findings in Häussler et al. [24] and Fanselow et al. [30]. It is only in the condition
546 with matched animacy that a superiority effect can be seen, and even here the condition is
547 more acceptable than three of the filler levels (including filler levels D and E that are
548 interpretable), suggesting that the sentence is not categorically ill-formed. This finding is
549 consistent with the claim that German has no superiority condition in the grammar and shows
550 only selective superiority effects in cases where processing difficulty is increased (Haider
551 [26]).

552 Finally, there is a general increase in well-formedness from the first to the second block of
553 trials and this increase was stronger for “wer-was” constructions than the other three. Thus, it
554 appears that the difference in well-formedness with an animacy mismatch in subject-initial
555 targets can be overcome with modest exposure with the mismatch condition.

556 **Targets in the context of fillers**

557 Is the increase in well-formedness from block 1 to block 2 related to sentence type – and
558 therefore to natural syntactic classes - or is it rather a reflection of the general level of well-
559 formedness? We use the ratings of fillers to address this question. As with targets, there was
560 no significant change in well-formedness ratings from block 2 to block 6 for the levels of
561 fillers relevant for the comparison with targets. Therefore, we averaged the performance of
562 the final five blocks for the ten types of sentences. Fig 1b shows the corresponding pattern of
563 changes between block 1 and the mean of blocks 2 to 6.

564 The four types of targets were presented together with six different levels of fillers A to E
565 which were expected to cover the spectrum from clearly grammatical (A and B) to
566 increasingly ungrammatical (C to F). Experience with “unusual” grammatical targets that are
567 licensed by the grammar (i.e., object-initial targets) should lead to a larger gain in ratings of
568 well-formedness than typical grammatical fillers like A and B (i.e., the latter should already
569 be rated at a very high level – there is no or only little room for improvement) and also to a

570 larger gain than for ungrammatical fillers like C to F (i.e., they should be rated at a low level
571 and stay there because they are not licensed by the grammar).

572 Critical tests refer to interactions between block and neighboring sentence types. The contrast
573 for the two subject-initial targets illustrates interactions in support of grammatical activation:
574 “wer_wen” is rated higher than “wer_was” ($b=0.30$, $z=4.26$) and shows less of a gain between
575 blocks ($b=-0.44$, $z=-3.65$). However, in the middle range of well-formedness ratings there is
576 no evidence for differential gains: “wer_was” is rated higher than filler level C ($b=0.32$,
577 $z=2.42$), but there is no significant interaction with block ($b=-0.15$, $z=-1.09$). Filler level C is
578 rated higher than the object-initial targets “was_wer” ($b=0.42$, $z=2.91$), but exhibits a small,
579 but significantly larger gain ($b=0.27$, $z=2.12$). For the two object-initial targets we observe
580 higher rating ($b=0.40$, $z=4.24$) for “was_wer” than “wen_wer”, but no significant difference
581 in gain ($b=-0.15$, $z=-1.16$). There is one ambiguity: “wen_wer” is rated significantly higher
582 than filler level D ($b=0.46$, $z=2.70$) and significantly benefitted from exposure whereas filler
583 level D stayed at the low level of acceptability ($b=0.34$, $z=2.54$). In general, the results do not
584 allow us to reject the hypothesis that initial gains are due to the general level of well-
585 formedness, but there is some ambiguity for low levels of well-formedness.

586 The experiment fulfilled our expectations only partially. There is a zone in the judgment
587 space (3.5 to 4.5) in which we see pervasive satiation. Furthermore, level E satiates, extending
588 the satiation zone below a mean rating of 2, but also note that the rise of E occurs with a
589 delay. Level D did not show evidence for a rise, despite its higher initial acceptability than
590 level E. Hence, the behavior of levels D and E was unexpected.

591 In summary, the German data shows that (i) regardless of exposure, the strength of superiority
592 effects can be modulated by manipulating animacy (here we replicate previous findings in the
593 literature, e.g., Häussler et al. [24]); (ii) exposure effects can be reliably induced
594 experimentally; and (iii) exposure effects may be a property of intermediate judgements rather
595 than of certain types of syntactic constructions.

596 **Experiment 2: rating of English sentences**

597 For German, many authors assume that object-initial multiple questions are indeed
598 grammatical because the superiority condition (or the more general constraints implying it)
599 can be circumvented or fail to apply because of peculiarities of German sentence structure.
600 When crossing movement is less acceptable, such differences in intuitions are taken to result
601 from grammar-external factors such as increased processing complexity (e.g., Haider [26]).
602 Do patterns of exposure effects change in a language that has a grammatical superiority
603 condition? In this section, we report a (partially) parallel study to the German one in English,
604 where it is still controversial whether superiority violations are ungrammatical or not (e.g.,
605 Hofmeister et al. [28], Häussler et al. [24] for recent discussion). We test whether we find the
606 same or a different pattern of exposure effects to German.

607 The animacy variation in the German items had the purpose of having effects of crossing
608 movement with different strength. Animacy variations seem to have no such effect in English
609 (Häussler et al [24]), but there is another differentiation with a similar consequence in
610 English, viz. so-called discourse linking (d-linking). A discourse linked wh-phrase such as
611 “which book” asks for a specific item among a contextually introduced set of objects.
612 Discourse-linking changes the well-formedness of wh-island constructions, and also the well-
613 formedness of multiple wh-questions by eliminating penalties for crossing movement, so that
614 crossing movement may no longer be ungrammatical even in English (Pesetsky [44] and [45],
615 cf. also Featherston [27] for experimental evidence supporting this assumption in the
616 generative literature.

617 With these considerations, the experiment again corresponds to a 2 x 2 design with subject- or
618 object-initial wh-words and with or without discours-linked wh-phrases. Examples of the four
619 English target conditions are given in (5).

620 (5) English target conditions

621 a. Condition 1: subject-initial; non-discourse-linked object: *who* – *what*;

622 well-formed in English

623 *The housekeeper forgot **who** had dropped **what** during the party.*

624 b. Condition 2: subject-initial; discourse-linked wh-phrases: *which* N_{subj} – *which* N_{obj};

625 well-formed in English

626 *The housekeeper forgot **which guest** had dropped **which glass** during the party.*

627 c. Condition 3: object-initial; non-discourse-linked object: *what* – *who*;

628 ill-formed in English

629 *The housekeeper forgot **what who** had dropped during the party.*

630 d. Condition 4: object-initial; discourse-linked wh-phrases: *which* N_{obj} – *which* N_{subj};

631 not grammatically ill-formed in English

632 *The housekeeper forgot **which glass which guest** had dropped during the party.*

633 The first factor ‘superiority’ can be seen by comparing (5a-b) with (5c-d). In the examples in

634 (5a) and (5b), the subject appears before the object and thus fulfills the superiority condition.

635 Examples (5c) and (5d) show crossing movement, and at least (5c) violates the superiority-

636 condition in most if not on all accounts. In terms of linear order, (5d) should fall into the same

637 category, but at least some models take the superiority condition to be inapplicable here. The

638 second factor of discourse-linking can be seen by comparing (5a) and (5c) with (5b) and (5d).

639 In (5a) and (5c), the subject and object are the non-discourse-linked wh-words *who* and *what*,

640 whereas in (5b) and (5d) the subject and object are the discourse-linked DPs *which guest* and

641 *which glass*.

642 **Method**

643 **Subjects**

644 A total of 48 English native speakers participated in the study. Fourteen self-identified as

645 male, 32 self-identified as female and 2 subjects selected “other” under gender. Ages ranged

646 between 18 and 48, with an average age of 31 years. Recruitment was carried out through the

647 web-based recruitment platform, Prolific. Participation was voluntary and was remunerated

648 with £6.80. All participants were naive with regard to the questions and objectives of the
649 study.

650 **Apparatus**

651 The experiment was conducted using Ibox software on the web-based Ibox Farm server
652 (<https://spellout.net/ibexfarm>, developed by Alex Drummond). For the lab-based study, we
653 had generated distinct questionnaires for each participant. To retain distinct counterbalancing
654 in this web-based version, we created 48 distinct questionnaires with unique links.

655 Participants first clicked on a welcome page and then received a unique link, in randomized
656 order. Clicking on the link led them to the questionnaire. The stimulus material was presented
657 on a computer screen along with the numbers 1 to 7 arranged horizontally in small boxes. A
658 short description of the degree of well-formedness was given under the numbers 1 “not at all
659 well-formed” and 7 “completely well-formed”. No description was given for numbers 2 to 6.
660 Sentence ratings of “1” to “7” could either be entered on a keyboard or by pointing and
661 clicking on the number.

662 **Material**

663 As in the German study, material consisted of 120 targets and 252 fillers. The targets were
664 translations of the German material with some adjustments. Instead of animacy we
665 implemented discourse linking as a second factor. We also left out adverbs because including
666 adverbs in English did not make the sentences sound more natural, and changed the content
667 nouns in some sentences. A complete list of targets and fillers is provided in the OSF
668 Repository (<https://osf.io/ge2db/>).

669 As for German sentences, targets represent indirect multiple wh-questions in subordinate
670 clauses as shown above in (5). The fillers represented six levels of well-formedness, as with
671 German. Each level of well-formedness was made up of 42 items. The top five levels were
672 each made up of three constructions that were identified in Gebrich et al. [11] as consistently
673 rating at that gradient level. To create the fillers, we used the three example items from

674 Gebrich et al. [11] for each level and created additional items until we reached the desired
675 number of 42. We then added an additional sixth gradation (F) "not at all well-formed" to
676 better reflect the lower end of the spectrum of acceptability.

677 (6) English filler examples (based on Featherston fillers)

678 a. Level A:

679 The patient fooled the dentist by pretending to be in pain.

680 There is a statue in the middle of the square.

681 The winter is very harsh in the North.

682 b. Level B:

683 Before every lesson the teacher must prepare their materials.

684 Jane does not boast about her being elected president.

685 Jane cleaned her motorbike with which cleaning cloth?

686 c. Level C:

687 Hannah hates but Linda loves eating popcorn in the cinema.

688 Most people like very much a cup of tea in the morning.

689 The striker must have fouled deliberately the goalkeeper.

690 d. Level D:

691 Who did she whisper that had unfairly condemned the prisoner?

692 The old fisherman took her pipe out of mouth and began story.

693 Which professor did you claim that the student really admires her?

694 e. Level E:

695 Historians wondering what cause is disappear civilisation.

696 Old man he work garden grow many flower and vegetable.

697 Student must read much book for they become clever.

698 f. Level F:

699 The ink was for spilled.

700 For the construction we used native speaker intuitions of naturalness from one of the
701 investigators as the guiding principle. We again based the fillers on calibrated sets of
702 distractor items developed by an independent research group, outlined for English in Gebrich
703 et al. [11]. As in German, we started from the one example per construction cited in the paper
704 for each of the top 5 levels of well-formedness, created further items on that same template,
705 and then created a sixth uninterpretable filler level to reach a total of 252 fillers.

706 **Design and counterbalancing**

707 Counterbalancing was identical to the German study. In order to implement the individualized
708 counterbalancing scheme over the internet, we created unique questionnaires for each
709 participant ID, and assigned participants to IDs using a random link generator in Ibex Farm.

710 **Procedure**

711 Subjects were instructed to judge the well-formedness of the sentences according to their
712 judgements of spoken language, rather than written language that they might find in a
713 textbook, and to use the full spectrum of the seven-point scale. They practiced the procedure
714 with ten examples. Then, each subject worked through the six blocks of sentences; durations
715 of breaks between blocks was under the subjects' control.

716 **Statistical analysis**

717 Statistical analysis for the English study followed the same procedure as in Experiment 1,
718 except that the order of levels for type of sentence was different and consequently also the
719 contrasts resulting from the application of the sequential-difference contrast to this factor.

720 **Results and Discussion**

721 **Targets**

722 The four types of targets using constructions of (a) "who-what", (b) "which N_{subj}-which N_{obj}",
723 (c) "what-who", and (d) "which N_{obj}-which N_{subj}" for the same sentence frames map onto a 2
724 x 2 design with the main effects subject-vs-object initial targets (so; a+b-c-d) and discourse
725 linking of subject and object (dlink; a+c-b-d) as well as the interaction between these two

726 effects (a-b-c+d). In addition, this 2 x 2 design was repeated across six blocks of trials. Fig 2a
727 displays the change in ratings of well-formedness for the four types of targets shown in red;
728 Fig 2a also shows performance for the levels of fillers. The LMM fixed-effect estimates are
729 given in a supplement (S2 Table).

730 **Fig 2. Well-formedness ratings (a) by block and (b) for first block and mean of later blocks for English**
731 **targets and fillers (web experiment).**

732 Here again, the two subject-initial targets were rated as more acceptable than the two object-
733 initial ones, yielding a significant main effect of order ($b=0.67$, $z=12.25$). Overall, there was
734 also a significant difference between targets with discourse-linked wh-phrases to targets
735 without discourse-linked wh-phrases ($b=-0.35$, $z=-6.76$) qualified by a significant interaction
736 with order ($b=0.43$, $z=9.49$): Subject-initial word order was clearly preferred with or without
737 discourse-linking to the other two conditions; for object-initial word order there was a clear
738 preference for the discourse-linked construction (“which N_{obj} – which N_{subj}”).

739 The change in ratings across blocks showed an overall profile that was less clear than for
740 German targets: There was a nominal, but not significant increase from block 1 to the average
741 of the following blocks ($b=0.16$, $z=1.77$) and, counter to expectations, a significant overall
742 negative difference between block 3 and the final three blocks ($b=-0.11$, $z=-2.47$). However,
743 there were two significant interactions between the first and second contrast for block with
744 order. There was a significant increase in well-formedness from the first block to the average
745 of the others for object-initial but not for subject-initial targets ($b=-0.29$, $z=-4.95$) and a weaker
746 effect for the change from the second to the average of the rest ($b=-.14$, $z=-.3.17$). The second
747 interaction was “helped” to some degree by a small reduction of well-formedness of subject-
748 initial targets in the final blocks. Finally, there was also a significant interaction for the second
749 contrast of block with discourse linking due to an increase in well-formedness for the non-
750 discoursed linked targets (“what-who”; $b=0.10$, $z=2.28$). These are the targets with the lowest
751 well-formedness of the four conditions.

752 **Targets in the context of fillers**

753 As with German sentences, it is instructive to examine the well-formedness of English targets
754 in the context of fillers (see Fig 2). Unlike for German sentences we do not see the prototype
755 pattern for the two subject-initial targets; indeed, there is no evidence for differences in rating
756 ($b=0.14$, $z=1.23$) and they do not show evidence for a differential change in rating ($b=0.10$,
757 $z=0.81$); also, filler level B is statistically not distinguishable from SO_wh (both $z < 1$).

758 In contrast, the two types of object-initial targets clearly increased from block 1 to the average
759 of the other blocks, but this increase was similar to fillers of type C and D. In the post-hoc
760 LMM, the lines for these four sentence types were statistically parallel, meaning that there
761 was no interaction with block for the three pairwise contrasts (all $|z\text{-values}| < 1.34$). Thus, the
762 English sentences in the middle range of well-formedness between 2.8 and 4.5., show the
763 expected behavior since there is a pervasive satiation effect in this zone of judgment. Unlike
764 for German sentences, there is no ambiguity in this respect even for fillers of level D. Level E
765 can be added to the satiation zone as well. So, the predictions of our hypothesis are fulfilled
766 by Experiment 2. Just like in German, level F fails to follow this pattern, in starting out at a
767 level similar to E, but not showing any indication of satiation.

768 In summary, in English targets, just like in the German experiment, there is a large difference
769 in the judgement of well-formedness of multiple wh-questions with a clear preference for
770 subject-initial over object-initial targets, in line with previous experiments (Sprouse [20],
771 Hofmeister et al. [28]). For object-initial targets there was a very clear preference for
772 discourse-linked than not-discourse-linked targets. There was positive initial change for the
773 well-formedness of object-initial targets, but this change might simply reflect a general
774 middle-raise in well-formedness also observed for fillers with a similar initial rating of well-
775 formedness.

776 **Experiment 3: Rating of English sentences (replication in** 777 **lab)**

778 In Experiment 3, we report a replication of the English experiment. Time-wise, this
779 experiment was carried out after Experiment 1 and before the web-based Experiment 2. We
780 chose to run Experiment 2 because the counterbalancing scheme was not rendered as intended
781 in Experiment 3. Although the targets were not counterbalanced as intended across blocks in
782 Experiment 3, some aspects of the counterbalancing were preserved: each participant still saw
783 an individual questionnaire, as well as more than one target item per block. Fillers were
784 counterbalanced across blocks as intended, and ordering restrictions such as making sure that
785 targets were not adjacent were maintained. We submit that past studies counterbalanced target
786 conditions, but, as far as we could tell from procedural descriptions, earlier studies did not
787 counterbalance conditions across blocks either.

788 Most importantly, one important result of Experiments 1 and 2 was that the change in well-
789 formedness between initial blocks of trials for targets might simply reflect their level of well-
790 formedness. This argument was based on the absence of evidence for differences in change
791 when compared to ungrammatical sentences with comparable ratings of well-formedness. The
792 argument rests on arguing a null hypothesis of parallel changes. Such results are in need of
793 replication.

794 **Method**

795 **Subjects**

796 A total of 48 subjects (28 female, 20 male) participated in the study. Ages ranged between 16
797 and 51, with an average age of 23 years. Recruitment was carried out through the participant
798 pool (SONA) at University College London (Psychology and Language Sciences) and the
799 University of Cambridge (Language Sciences). Participation was voluntary and was
800 remunerated with £6. All participants were naive with regard to the questions and objectives
801 of the study.

802 **Apparatus**

803 The experiment was conducted in the Speech and Language Sciences Lab at University
804 College London and at Trinity Hall, Cambridge. The experiment was implemented using the
805 Python software PsychoPy [32] version 1.84.2. The stimulus material was presented on a
806 computer screen along with a scale labelled from 1 to 7. Underneath each of the numbers was
807 a short description of the degree of well-formedness: 1 “not at all grammatically natural”, 2
808 “almost not grammatically natural”, 3 “more grammatically unnatural than natural”, 4 “cannot
809 be rated”, 5 “more grammatically natural, than unnatural”, 6 “almost grammatically natural”,
810 7 “completely grammatically natural”. Sentence ratings from “1” to “7” were entered on a
811 standard computer keyboard with a German keyboard layout.

812 **Material**

813 Material was identical to Experiment 2.

814 **Design and counterbalancing**

815 We did not counterbalance the targets across blocks in this study, due to an error in generating
816 the questionnaires. Some of the features of the counterbalancing scheme used in Experiment 1
817 and Experiment 2 were nonetheless retained, such as different questionnaires for each
818 participant, and counterbalancing of fillers across blocks.

819 **Procedure and statistical analysis**

820 Subjects were instructed to judge the well-formedness of the sentences according to their
821 judgements of spoken language, rather than written language that they might find in a
822 textbook, and to use the full spectrum of the seven-point scale. They practiced the procedure
823 with ten examples. Then, each subject worked through the six blocks of sentences; durations
824 of breaks between blocks was under the subjects’ control. Statistical analysis for Experiment
825 3 followed the procedure for Experiment 2.

826 **Results and Discussion**

827 **Targets**

828 Fig 3a displays the change in ratings of well-formedness for the four types of English targets
 829 shown in red and the six types of fillers in blue. The LMM fixed-effect estimates and statistics
 830 are given in a supplement (S3 Table). Main effects of order ($b=.74$, $z=12.93$) and discourse-
 831 linking ($b=-0.45$, $z=-8.87$) as well as the interaction of these two factors ($b=0.56$, $z=10.97$)
 832 were replicated in the lab experiment: again, subject-initial word order was clearly preferred
 833 with or without discourse-linking to the other two conditions; for object-initial word order
 834 there was a clear preference for the discourse-linked construction (“which N_{obj} – which
 835 N_{subj} ”).

836 **Fig 3. Well-formedness ratings (a) by block and (b) for first block and mean of later blocks for English**
 837 **targets and fillers (lab experiment).**

838 The change in ratings across blocks was much clearer than in Experiment 2: The first two
 839 block contrasts were significant (b1: $b=0.39$, $z=4.54$; b2: $b=0.14$, $z=3.53$) indicating an
 840 overall increase of well-formedness from the first to the average of the rest and a second,
 841 smaller increase from block 2 to the rest. Again, there were two significant interactions
 842 between the first and second contrast for block with order. There was a significant increase in
 843 well-formedness from the first block to the average of the others for object-initial but not for
 844 subject-initial targets ($b=-0.13$, $z=-2.24$) and for the change from the second to the average of
 845 the rest ($b=-.12$, $z=-.3.09$). As in Experiment 2, the increase in well-formedness across the
 846 initial blocks was larger for object-initial than subject-initial target targets (see Fig 3b). Unlike
 847 in Experiment 2 there was also an increase for subject-initial target targets.

848 **Targets in the context of fillers**

849 The final question is whether the different increases for subject- and object-initial targets are
 850 different than those of fillers of similar well-formedness. The results replicate Experiment 2.

851 In fact, the pattern of results of Experiment 3 comes as close to our expectations as
852 Experiment 2 because there is a continuous zone of acceptability in which all constructions
853 show satiation. The satiation zone reached out higher in Experiment 3 than Experiment 2, so
854 that only filler level A was not affected. Just as in the preceding two experiments, filler level
855 F behaved in an unpredicted way by not showing any signs of rising judgments. We also note
856 that the unexpected negative effect and the interaction for the second contrast of block with
857 discourse linking reported for Experiment 2 did not replicate.

858 **General Discussion**

859 **Summary of combined results**

860 We found parallel patterns between German and English for both exposure effects and
861 modulations of their size by crossing movement and by factors related to semantics and
862 discourse. Specifically, we found:

- 863 a. A rise in initial block(s) only;
- 864 b. A rise in a continuous zone of acceptability irrespective of sentence type, that is shifts
865 were found (i) in target conditions that respected superiority; (ii) target conditions that
866 violated superiority; and (iii) filler levels with no wh-elements at all;
- 867 c. Comparable shifts in German and English, both in initial blocks and in continuous
868 zones of well-formedness;
- 869 d. Effects of D-linking in English that were on par with animacy effects in German,
870 specifically:
 - 871 i. A decrease for German with superiority violations with matching animacy that was
872 comparable to regular (non-d-linked) superiority violations in English
 - 873 ii. An increase for English with superiority violations where both wh-phrases were D-
874 linked that was comparable to acceptable (mismatched animacy) superiority
875 violations in German.

876 **Satiation effects track zone of well-formedness, not grammatical**
877 **construction**

878 Our experiments sought to establish whether there are “satiation effects” in the sense of
879 studies building on Snyder [5] caused by repeated exposure, and whether these effects are
880 zone- or construction-sensitive. All three experiments showed a clear rise in acceptability for
881 some of the conditions of the target experiment and for some of the filler groups between the
882 first and the remainder of the blocks. It is difficult to conceive of an alternative account for
883 these effects observed across three experiments in two languages than the assumption that
884 ratings increased with the number of items perceived in a certain target experimental
885 condition or in a certain filler group.

886 Is satiability driven by syntactic properties or does it indiscriminately affect all constructions
887 of a certain zone of acceptability? The results of our experiments by and large support the
888 latter hypothesis. Both in the German and the two English experiments, there was a
889 continuous zone (with two exceptions discussed below) in the range of acceptability such that
890 all constructions in that zone satiated, and no items outside that zone showed an increase in
891 acceptability. This was the case even though the constructions affected have no obvious
892 grammatical factor in common. In German, for instance, ratings rise for the two object-initial
893 multiple questions (target conditions 3 and 4), for the well-formed subject-initial questions
894 with mismatching animacy (target condition 2), and for fillers C which involve unusual
895 binding options and ill-formed orders of dative and accusative objects. One obvious factor
896 these constructions have in common is their level of acceptability. In English we see the same
897 rise in acceptability across dissimilar constructions as in German, despite the effects
898 pertaining to “harder” grammatical constraints. For instance, crossing movement improves,
899 but so do constructions that do not respect verb-object adjacency constraints, constructions
900 where the subject pronoun *who* is ungrammatically extracted across the complementiser *that*,
901 and where pronouns do not match the gender of their antecedent, for example, *fisherman* -

902 *her*. What unifies all these constructions is that participants give them ratings in the same
903 range.

904 Our results come with two problematic aspects, however. First, filler group D did not satiate
905 in German. This could merely indicate that the satiation zone of Experiment 1 begins at quite
906 a high point in the judgment space (the result would then be less of a problem) or it could
907 point to the existence of set of constructions that are well in the satiation zone but
908 nevertheless fail to satiate, if the dynamics of judgments for filler group E is a subcase of
909 satiation. For the latter option, we have little to offer as to why it fails to conform to the
910 expected picture. The English experiments do not come with such problematic constructions.

911 The other problem is related to group E, for which ratings have gone up in all three
912 experiments. As such, this is not a problem – it would merely show that the satiation zone
913 begins quite low in our experiments, but the non-congruent behaviour of group F is
914 incompatible with the idea of a zone of judgment linked to satiation because E and F start out
915 at roughly the same acceptability value.

916 This discrepancy may reflect that our decision to add an even lower level F of acceptability to
917 the five filler groups designed by Featherston and colleagues was infelicitous. First, our
918 expectation that the newly added group F might receive the lowest ratings only slightly above
919 1 was not fulfilled, it thus failed to serve the function it was meant to have in our experiments.
920 Second, and probably more importantly, these fillers were created *ad hoc* from a mixture of
921 filler material from previous unrelated experiments, and were thus neither standardized nor
922 calibrated. We are tempted to make this unfortunate property responsible for the irregular
923 behaviour of filler group F. Thus, we have to leave the uncertainty in determining the lower
924 bound of the zone of satiation unresolved.

925 As for the upper bound, Experiments 1 and 2 show satiation below a value of 4.5, while
926 satiation extended to a level of 5.5 in Experiment 3, i.e., well in the area approaching full

927 acceptability. The interpretation that the satiation zone always originates in the intermediate
928 range of acceptability but may extend further into higher or lower areas comes with the least
929 number of commitments, and we lack evidence that would allow us to go for a stronger
930 interpretation.

931 Apart from the two difficulties just discussed, our conclusion concerning the postulation of a
932 “zone of satiation” can be criticized in at least two ways. With respect to the current
933 experiments, it might be objected that the satiating constructions do have a common
934 grammatical or psycholinguistic property overlooked by the present authors, or that the
935 factors making satiation possible are manifold, such that all our relevant items simply happen
936 to fall under one or the other of these factors. Both possibilities are real, but difficult to refute
937 in the absence of a concrete proposal as to what these factors have been.

938 **Dynamics of satiation**

939 Our results go against the expectation that satiation effects rise continuously to an asymptotic
940 level of maximal well-formedness. There is an initial early increase followed by stability
941 clearly below the maximum. This result may well be the most remarkable feature of the plots
942 of judgments across blocks and across the three experiments. The size of initial rise of ratings
943 is not very dramatic either; they increase only by roughly .5 on a 7-point scale, meaning that
944 repeated exposure does not change the quality of the judgment. The relatively small size of
945 the effect is in agreement with most of the experimental results reported in the literature. The
946 stability we found could be (1) due to predominantly stable judgements, (2) due to
947 participants developing a rating pattern in the first two blocks that is used for the remainder of
948 the experiment (entrenchment), or (3) due to a “mere exposure effect” (Zajonc, [46]).

949 If the first option applies, participants might have stopped deliberately rating the item they are
950 presented with and instead produced memorised values to assign to types of sentences. There
951 is currently no accepted way to control for such memorising and entrenchment. Some limited

952 discussion in previous works like Sprouse [13] presents a potential argument that memorising
953 and entrenchment undermine a balanced design and the use of sets of items of the same type,
954 as is the case with control items and fillers. Alternatively, according to the second option,
955 participants might have taken their time to find a stable mapping from their perceived
956 intuitions to a 7-point scale. Once that mapping is established ratings remain constant. Both
957 explanations, however, do not explain why the stabilization of the judgments is always
958 upwards oriented.

959 The “mere exposure effect” refers to observations that experimentally manipulated frequency
960 of exposure to nonsense words and syllables increases their ratings of positive affective
961 connotations (e.g., familiarity, liking, etc.) without any reinforcement (Zajonc, [46]). Bader
962 and Häussler [47] obtained such a correlation for printed frequency of sentences and their
963 rating of well-formedness. We envision that this account could open a promising line of
964 research, especially if consequences for ratings of acceptability due to differences in the
965 quality of prosodic representations between sentences varying in degree of well-formedness
966 are taken into account as well.

967 **The integrated analysis of targets and fillers**

968 We end with a methodological note. Is it appropriate to compare within-item targets to
969 between-item fillers? A norm in the field of experimental syntax is to test only within-item
970 differences (Sprouse et al., [48]). Within-item factorial designs are indeed a powerful way to
971 control for item variability. In a principled theory-driven workflow, within-item designs allow
972 us to test theory-relevant contrasts specified *a priori* with tight control of known sources of
973 the variability. Fillers, however, should be as unrelated as possible to the targets. Deriving
974 them as within-item variants of targets However, the fillers used in these experiments are not
975 just any old fillers that we averaged into groups after testing. Rather, we had *a priori*
976 expectations about the gradient acceptability of these fillers based on the norming studies in
977 Featherston [10] and Gebrich et al. [11]. We also carefully built the additional fillers that we

978 used to the template (*Eichsatz*) used in the norming study. Moreover, results did indeed meet
979 expectations that we had before testing; they differed significantly in the expected order.
980 Therefore, although part of the analysis presented in this paper is exploratory – in particular,
981 we decided which target conditions to compare to which other target or filler conditions after
982 seeing the results – we chose the levels of acceptability before testing based on previous
983 research results. It may turn out that future research will uncover interesting commonalities
984 between fillers that we have missed here. Such a development is part and parcel of
985 programmatic research however, and does not undermine the validity of comparing targets to
986 normed levels of gradient fillers, provided that norming is carried out before testing and the
987 differences between acceptability levels come out as expected.

988 **Conclusion**

989 In conclusion, this paper revisited old questions about satiation, superiority effects and
990 processing complexity of exposure effects across a single experimental session. The
991 combination of (1) targeting superiority violations, (2) integrating six blocks of trials as an
992 experimental factor into the counterbalancing scheme, (3) going beyond previous research
993 with respect to number of subjects and number of items, (4) integrating carefully selected
994 levels of filler sentences in a secondary joint analysis of targets and fillers, (5) replicating the
995 overall profile of means in two languages, that is German (Experiment 1) and English
996 (Experiments 2 and 3), and (6) replicating, for English, the overall profile of means between
997 web (Experiment 2) and lab (Experiment 3) should provide a useful reference platform for
998 follow-up research. Moreover, all our data and scripts are available for additional exploratory
999 re-analyses from different theoretical perspectives. We neither claim that we resolved all open
1000 questions – indeed there are a few results that are inconsistent with our perspective – nor that
1001 the results generalize far beyond the experimental setting (e.g., results may change with
1002 comprehension questions, see Zervakis & Mazuka, [9]).

1003 The main results are that after an initial rise, there was a remarkable quantitative asymptotic
1004 stability of acceptability ratings within experiments and there was also remarkable qualitative
1005 agreement between experiments with German and English sentences for strong differences in
1006 syntactic violations. The results corroborated claims of satiation effects and strongly suggest
1007 that the nature of these effects is different than claimed in previous studies: rather than a
1008 continuous rise to maximal acceptability, judgments rise only initially and asymptote
1009 significantly below the maximum. Thus, the effects appear to be primarily linked to an
1010 intermediate zone of well-formedness, not to natural syntactic classes.

1011 **Acknowledgements**

1012 We would like to thank Kim-Laura Speck for help developing the counterbalancing scheme
1013 used in the experiments, and Johannes Rothert for extracting the values in Table 1 from the
1014 relevant publications.

1015 **References**

- 1016 1. Schütze C. The empirical base of linguistics: grammaticality judgments and linguistic
1017 methodology. Chicago: University of Chicago Press; 1996.
- 1018 2. Cowart W. Experimental Syntax: applying objective methods to sentence judgments.
1019 Thousand Oaks: SAGE publications; 1997.
- 1020 3. Nagata H. The relativity of linguistic intuition: The effect of repetition on
1021 grammaticality judgments. *Journal of Psycholinguistic Research*. 1988;17(1):1–17.
- 1022 4. Carroll JM. Complex compounds: phrasal embedding in lexical structures.
1023 *Linguistics*. 1979;17:863–877.
- 1024 5. Snyder W. An experimental investigation of syntactic satiation effects. *Linguistic
1025 Inquiry*. 2000;31:575–82.
- 1026 6. Chomsky N. *Lectures on Government and Binding*. Dordrecht: Foris; 1981.
- 1027 7. Chaves RP, Dery JE. Which subject islands will the acceptability of improve with
1028 repeated exposure? In: Santana-LaBarge RE, editor. *Proceedings of the 31th West
1029 Coast Conference on Formal Linguistics*. Cascadilla Proceedings Project; 2014. p. 96–
1030 106.
- 1031 8. Chaves RP, Dery JE. Frequency effects in subject islands. *Journal of Linguistics*.
1032 2019; 55(3):475–521.
- 1033 9. Zervakis J, Mazuka R. Effect of repeated evaluation and repeated exposure on
1034 acceptability ratings of sentences. *Journal of Psycholinguistic Research*. 2013;42(6):
1035 505–525.
- 1036 10. Featherston, S. Why linguistics needs boiling and freezing points. In: Featherston S,
1037 Winkler S, editors. *The Fruits of Empirical Linguistics*. Berlin:
1038 De Gruyter; 2009
- 1039 11. Gebrich H, Schreier V, Featherston S. Standard items for English judgement studies:
1040 syntax and semantics. In: Featherston S, Hörnig R, von Wietersheim S, Winkler S,

- 1041 editors. *Experiments in Focus: Information Structure and Semantic Processing*. Berlin:
1042 de Gruyter; 2019. p. 305–328.
- 1043 12. Snyder W. Satiation. In: Goodall G, editor. *The Cambridge Handbook of*
1044 *Experimental Syntax*. Cambridge: CUP; 2021.
- 1045 13. Sprouse J. A program for experimental syntax: finding the relationship between
1046 acceptability and grammatical knowledge. University of Maryland; 2007.
- 1047 14. Goodall G. Syntactic satiation and the inversion effect in English and Spanish wh-
1048 questions. *Syntax*. 2011;14(1):29–47.
- 1049 15. Christensen KR, Kizach J, Nyvad AM. Escape from the island: Grammaticality and
1050 (reduced) acceptability of wh-island violations in Danish. *Journal of Psycholinguistic*
1051 *Research*. 2013;42(1):51–70.
- 1052 16. Braze FD. Grammaticality, acceptability and sentence processing: A psycholinguistic
1053 study. University of Connecticut. Storrs, CT; 2002.
- 1054 17. Crawford J. Using syntactic satiation to investigate subject islands. In: Choi J, Hogue
1055 EA, Punske J, Tat D, Schertz J, Trueman A, editors. *Proceedings of the 29th West*
1056 *Coast Conference on Formal Linguistics*. Cascadilla Proceedings Project; 2012. p. 38–
1057 45.
- 1058 18. Francom JC. *Experimental syntax: Exploring the effect of repeated exposure to*
1059 *anomalous syntactic structure: Evidence from rating and reading tasks*. University of
1060 Arizona. Tucson, AZ; 2009.
- 1061 19. Hiramatsu K. *Assessing linguistic competence: evidence from children's and adults'*
1062 *acceptability judgments*. University of Connecticut; 2000.
- 1063 20. Sprouse J. Continuous acceptability, categorical grammaticality, and experimental
1064 syntax. *Biolinguistics*. 2007;1:123–134.
- 1065 21. Sprouse J. Revisiting Satiation: evidence for an equalization response strategy.
1066 *Linguistic Inquiry*. 2009;40:329–41.
- 1067 22. Haider H. *Symmetry breaking in syntax*. Cambridge: CUP; 2012.

- 1068 23. Kuno S, Robinson J. Multiple wh-questions. *Linguistic Inquiry*. 1972;3:463– 487.
- 1069 24. Häussler J, Grant M, Fanselow G, Frazier L. Superiority in English and German:
1070 Cross-Language Grammatical Differences? *Syntax*. 2015;18:235– 265.
- 1071 25. Fanselow G. Acceptability, Grammar, and Processing. In: Goodall G, editor. *The*
1072 *Cambridge Handbook of Experimental Syntax*. Cambridge: CUP; 2021.
- 1073 26. Haider H. The superiority conspiracy: four constraints and a processing effect. In:
1074 Stepanov A, Fanselow G, Vogel R, editors. *Minimality Effects in Syntax*. Mouton de
1075 Gruyter; 2004. p. 147–176.
- 1076 27. Featherston S. Universals and grammaticality: wh-constraints in German and English.
1077 *Linguistics*. 2005;43:667–711.
- 1078 28. Hofmeister P, Jaeger TF, Arnon I, Sag IA, Snider N. The source ambiguity problem:
1079 distinguishing the effects of grammar and processing on accept- ability judgments.
1080 *Language and Cognitive Processes*. 2013;28:48–87.
- 1081 29. Fanselow G, Weskott T. A Short Note on Long Movement in German. *Linguistische*
1082 *Berichte*. 2010;222:129-140.
- 1083 30. Fanselow G, Schlesewsky M, Vogel R, Weskott T. Animacy effects on crossing wh-
1084 movement in German. *Linguistics*. 2011;49:657–683.
- 1085 31. Williams EJ. Experimental designs balanced for the estimation of residual effects of
1086 treatment. *Australian Journal of Scientific Research*. 1949;2:149-168.
- 1087 32. Peirce JW. PsychoPy: Psychophysics software in Python. *Journal of Neuroscience*
1088 *Methods*. 2007;162(1–2):8–13.
- 1089 33. R Core Team. *R: A language and environment for statistical computing*. Vienna,
1090 Austria: R Foundation for Statistical Computing; 2020. URL [https://www.R-pro-](https://www.R-project.org/)
1091 [ject.org/](https://www.R-project.org/).
- 1092 34. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using
1093 lme4. *Journal of Statistical Software*. 2015.

- 1094 35. Wickham H, Averick M, Bryan J, Chang W, McGowan L, François R, et al.
1095 Welcome to the tidyverse. *Journal of Open Source Software*. 2019;4(1686).
1096 Doi:10.21105/joss.01686.
- 1097 36. Wilke CO. Cowplot: streamlined plot theme and plot annotations for ‘ggplot2’.
1098 <https://cran.r-project.org/web/packages/cowplot/index.html>; 2019.
- 1099 37. Lüdtke D. sjPlot: data visualization for statistics in social science. [https://cran.r-](https://cran.r-project.org/web/packages/sjPlot/)
1100 [project.org/web/packages/sjPlot/](https://cran.r-project.org/web/packages/sjPlot/); 2020.
- 1101 38. Bolker B, Robinson D. broom.mixed: tidying methods for mixed models.
1102 <https://cran.rproject.org/web/packages/broom.mixed/index.html>; 2020.
- 1103 39. Schad DJ, Vasishth S, Hohenstein S, Kliegl R. How to capitalize on a priori contrasts
1104 in linear (mixed) models: a tutorial. *Journal of Memory and Language*. 2020;110.
- 1105 40. Bates D, Kliegl R, Vasishth S, Baayen RH. Parsimonious mixed models. arXiv
1106 preprint arXiv:1506.04967.; 2017.
- 1107 41. Matuschek H, Kliegl R, Vasishth S, Baayen H, Bates D. Balancing type I error and
1108 power in linear mixed models. *Journal of Memory and Language*. 2017;94:305–315.
- 1109 42. Bates D, Alday P, Kleinschmidt D, Calderon JBS, Noack A, Kelman T, et al.
1110 *JuliaStats/MixedModels.jl*: v2.3.0. <https://doi.org/10.5281/zenodo.3727845>; 2020.
- 1111 43. Pesetsky D. Wh-in-situ: movement and unselective binding. In: ter Meulen A,
1112 Reuland E, editors. *The Representation of (In)definiteness*. Cambridge, MA: MIT
1113 Press; 1987. p. 98–129.
- 1114 44. Pesetsky D. *Phrasal movement and its kin*. Cambridge, MA: MIT Press; 2000.
- 1115 45. Zajonc RB. Attitudinal Effects of Mere Exposure. *Journal of Personality and Social*
1116 *Psychology*. 1968;9:1-27.
- 1117 46. Bader M, Häussler J. Toward a model of grammaticality judgments. *Journal of*
1118 *Linguistics*. 2010;46:273–330

- 1119 47. Sprouse J, Caponigro I, Greco C, Cecchetto C. Experimental syntax and the variation
1120 of island effects in English and Italian. *Natural Language and Linguistic Theory*.
1121 2016;34:307-344.

1122 **Supporting Information**

1123 **S1 Table. LMM parameter estimates for ratings of German target sentences**

1124 **(Experiment 1)**

1125 **S2 Table. LMM parameter estimates for ratings of English target sentences**

1126 **(Experiment 2)**

1127 **S3 Table. LMM parameter estimates for ratings of English target sentences**

1128 **(Experiment 3)**