

Scalar inferences in the acquisition of even

Yadav Gowda¹, Elise Newman¹, Leo Rosenstein¹ and Martin Hackl¹

¹Department of Linguistics and Philosophy, MIT, Cambridge , MA , USA

Correspondence*:

Yadav Gowda
ysg@mit.edu

Elise Newman
esnewman@mit.edu

2 ABSTRACT

Scalar inferences are ubiquitous in human reasoning. Correspondingly, language has many means of expressing and encoding them. One of these means is the focus particle *even*, which utilizes scalar inferences to signal the pragmatic status of asserted content as noteworthy. The vehicles that *even* employs to signal noteworthiness are scalar likelihood inferences. A peculiarity of these inferences is that they are presuppositional in nature (not-at-issue) and yet, they are responsive to the polarity of the sentence expressing the proposition whose likelihood is signaled. This property raises intricate questions about what learners might expect scalar operators of this sort to look like (initial hypothesis space) as well as what type of evidence and learning strategies they have access to as they figure out the specific properties of *even* in adult English. This paper presents a detailed study of this development, combining data from a series of comprehension experiments and corpus studies. We find that children are sensitive to the basic scalar nature of *even* much earlier than previous literature has claimed. We additionally find, however, that children sometimes exhibit non-adult-like responses to *even* sentences, which we argue provide insight into their developing grammar. On this view, the child grammar offers a larger option space for *even* than the adult grammar. Becoming adult-like, in turn, involves eliminating some of these options, namely those that are underutilized in production due to their limited conversational value.

Keywords: Language Acquisition, *Even*, Scalar Inferences, Presupposition, Focus Particles, Additive Particles, Polarity

1 INTRODUCTION

This paper investigates the status of scalar inferences in child grammar through an acquisition study of the scalar focus particle *even*. English *even* provides a window through which to study the grammatical encoding of scalar inferences due to its sensitivity to polarity. It triggers a least-likely presupposition in positive sentences, and a most-likely presupposition in negative sentences as illustrated in (1) and (2), (Horn (1969), Karttunen and Peters (1979), a.o.).

- (1) Jack had even invited AMY to the party.
- Amy was the least-likely for Jack to invite to the party.
- There was at least one other person (out of a salient set) that Jack had invited to the party.
- (2) Kim hadn't even invited SAM to the party.
- Sam was the most-likely for Kim to invite to the party.
- There was nobody else (out of a salient set) that Kim had invited to the party.

32 The meaning conveyed in (1) and (2) consists of three components: the asserted content, which is just the
 33 ordinary meaning of the sentence without *even* (i.e. the meaning of prejacent), and two invited inferences –
 34 scalar and additive in nature – which do not have the status of at-issue content. Rather they exhibit signature
 35 properties of presuppositions (or conventional implicatures), which persist even when the presupposition
 36 trigger occurs in the scope of entailment-cancelling operators. This can be readily seen when transforming
 37 these sentences into questions, (3) and (4), which still give rise to the same invited inferences.¹

- 38 (3) Did Jack even invite AMY to the party?
 39 - Amy was the least-likely for Jack to invite to the party.
 40 - There was at least one other person (out of a salient set) that Jack had invited to the party.
 41 (4) Did Kim not even invite SAM to the party?
 42 - Sam was the most-likely for Kim to invite to the party.
 43 - There was nobody else (out of a salient set) that Kim has invited to the party.

44 Interestingly, even though the invited inferences triggered by *even* project in questions, their precise
 45 character is affected by the polarity of the host sentence. Concentrating on the scalar inferences, which are
 46 the focus of this paper, we observed in (1) and (2) that the presence or absence of negation corresponds to
 47 a “most- or least-likely” inference respectively. The fact that these are different seems to suggest that the
 48 content upon which the scalar inference is based is visible to negation, though not in a way that would
 49 cancel the inference as is the case for asserted content.

50 This interaction between likelihood presupposition and polarity has generated much discussion in the
 51 theoretical literature. It also raises rather interesting questions about how learners acquire the full pattern of
 52 invited inferences triggered by *even*: What is the initial hypothesis space that learners need to navigate as
 53 they acquire *even*?² What is the evidence they have access to and make use of to transform their initial
 54 state of the grammar into the adult state, and what are the learning strategies that allow them to do that? We
 55 aim to address these questions via a series of comprehension and corpus studies (complemented by adult
 56 control experiments) that allow us to establish the developmental path children follow as they learn to track
 57 the different likelihood inferences triggered by *even* across different environments. Our findings portray a
 58 much richer and more intricate acquisition process than previous work has suggested, making it possible
 59 to identify the specific structure of the initial hypothesis space for *even* as well as the grammatical and
 60 pragmatic factors that enable and constrain the learning.

61 To foreshadow, our comprehension studies demonstrate that preschool-aged children show clear evidence
 62 of sensitivity to scalar inferences associated with *even*, contrary to what previous research has suggested.
 63 They also reveal a rather nuanced developmental trajectory traceable already in 4yoS regarding the at times
 64 non-adult-like nature of the scalar inferences. Importantly, our comprehension tasks reveal that not all
 65 non-adult-like behavior should be given the same analysis. While certain responses appear to indicate
 66 true confusion about *even*, other responses appear systematic and informed by their developing grammar,
 67 despite being non-adult-like. This is most clearly evident from the fact that these responses are accompanied
 68 by reasoned justifications involving reference to scalar properties.

¹ There is a lively debate as to the precise content, origin and status of the scalar and additive inferences. Since much of this debate is orthogonal to our purpose here, we label them either as “invited inferences” or as presuppositions. Moreover, we largely abstract away from the additive inference. See among others Horn (1972, 1989); Karttunen and Peters (1979); Rooth (1985); von Stechow (1991); Krifka (1991); Rullmann (1997, 2007); Herburger (2000); Schwarz (2005); Greenberg (2016, 2018); Francis (2018).

² Recent years have seen an impressive growth of cross-linguistic work on additive scalar particles like English *even* (Giannakidou, 2007; Crnič, 2009; Gast and van der Auwera, 2011; Greenberg, 2015, 2018, a.o.). A more complete framing of the initial hypothesis space for *even* would situate the current discussion within that typology. However, since the cross-linguistic inventory of such particles is rather rich and there is no consensus yet on the organizational principles spanning that typology, we will have to leave it to future work.

69 Our adult control studies present a similar finding. They, too, show a systematic error pattern during real
 70 time comprehension (albeit at a lower rate) which appears to be a function of the adult-grammar of *even*.
 71 We argue on the basis of their systematicity and their similarities that both the non-adult-like inferences
 72 generated by children and the corresponding errors in our adult studies reflect options made available by
 73 the basic architecture of the respective grammars of *even*. Hence, they should inform our theories of adult
 74 and child *even*, and put together, they should frame accounts of how *even* is acquired.

75 Our conclusions from the comprehension experiments are further enriched by our corpus studies on child
 76 and child-directed adult use of *even*, presented in Section 5. We find that children who produce *even* do
 77 so essentially error-free, even at 3-4 years of age. We additionally find that the form of our stimuli in the
 78 comprehension studies instantiates a low-frequency use pattern for *even* in both children and adults. The
 79 fact that children nevertheless show partial command of *even* when their comprehension is tested on these
 80 items suggests that children in our age range already have a quite robust grasp of the fundamental fact the
 81 *even* always triggers a scalar inference of some sort.

82 Most striking about the corpus data is, however, that children (and adults) do not exhibit the non-adult-like
 83 behavior that we find in the comprehension experiments. At first sight, this constitutes a rather surprising
 84 production-comprehension asymmetry: the child grammar appears adult-like in production but non-adult-
 85 like in comprehension. We propose to resolve this puzzle by assuming that it is the comprehension data
 86 that faithfully reflect the child grammar of *even*, which is non-adult-like in that it provides a larger space of
 87 options for *even* than the adult grammar. The fact that their production data is fully adult-like should, in
 88 turn, be seen as the result of child speakers choosing not to realize some of the grammatically licit options.
 89 We suggest that they underutilize some of these options because of their limited pragmatic utility. To
 90 explain why production is ahead of comprehension, we suggest that the pragmatic oddity of these options
 91 is more transparent from the perspective of the (child) speaker than from the perspective of the (child)
 92 listener since the speaker knows the intended message while the listener needs to infer it from the form of
 93 the utterance and the presumed conversational goals of the speaker.

2 EVEN'S SCALAR INFERENCES: THEORY AND ACQUISITION

94 2.1 Theoretical background

95 Our study's main focus is the interaction between *even*'s scalar inferences and sentential polarity as it
 96 relates to acquisition. As argued in Karttunen and Peters (1979) the scalar inferences have the status of
 97 a presupposition/conventional impudicature. This can be inferred, for instance, from the fact that, unlike
 98 at issue content, they survive embedding inside entailment cancelling environments such as questions,
 99 (3)(4). Interestingly, and somewhat unexpectedly the precise nature of the scalar inference is sensitive to the
 100 presence of negation (an entailment cancelling operator in its own right). Recalling the examples in (1) and
 101 (2) we see that the scalar inference in (2) is substantially different from the one in (1).

- 102 (1) Jack had even invited AMY to the party.
 103 - Amy was the least-likely for Jack to invite to the party.
- 104 (2) Kim hadn't even invited SAM to the party.
 105 - Sam was the most-likely for Kim to invite to the party.

106 In both (1) and (2), we detect an inference about the likelihood of somebody having been invited to a
 107 party. However, the inference in (2) is essentially the opposite of that in (1). While Amy is understood to
 108 be someone that was unlikely to be invited, Sam is understood to be someone who was very likely to be

109 invited. Both (1) and (2) therefore convey something surprising, namely that someone was invited who
 110 wasn't expected to be invited, or that someone wasn't invited who *was* expected to be invited.

111 There are two families of approaches that attempt to capture this property of *even*. The first type of
 112 approach ([Karttunen and Peters 1979](#), henceforth the *scope theory*) argues that *even* uniformly triggers
 113 a least-likely presupposition, but has a requirement that it outscope (clause-mate) negation³ resulting in
 114 most-likely inference. We can schematize this approach as in (5)-(7). *Even* is assumed to be a clausal
 115 operator that combines with a propositional argument (the prejacent) and a set of alternative propositions
 116 (derived from the prejacent via focus semantics⁴). Its lexical semantics is that of a filter which passes
 117 on the meaning of its propositional argument *p* unchanged but only if *p* is the least likely member in *C*,
 118 (6). For positive sentences this delivers a least-likely inference. For negated sentences, however, with
 119 *even* scoping over negation, as shown in (7) its prejacent will express a negative proposition resulting in a
 120 “least-likely-to-not” inference, which is, of course, equivalent to “most-likely-to” inference.

- 121 (5) [Even [(NOT) S]]
 122 (6) $\llbracket \text{even} \rrbracket^{w,g}(C)(p) \Leftrightarrow p \text{ is the least likely member of } C . p$ ⁵
 123 (7) [Even [Kim had not invited Sam to the party]]
 124 - Least-likely inference: *That Kim hadn't invited Sam to the party* is the least-likely proposition of
 125 the form *that Kim hadn't invited X to the party* \Leftrightarrow Most-likely inference: *That Kim invited Sam to*
 126 *the party* was most-likely.

127 The second approach ([Rooth 1985](#), henceforth the *ambiguity theory*)⁶ proposes that *even* is lexically
 128 ambiguous. One lexical entry is the one assumed by the scope theory, which has an (in principle) unrestricted
 129 distribution. The other lexical item, however, is a Negative Polarity Item (NPI) that comes with a most-likely
 130 presupposition, (9), and is restricted in its distribution to occur below negation, (8)-(10).

- 131 (8) [NOT [even_{NPI} [S]]]
 132 (9) $\llbracket \text{even}_{NPI} \rrbracket^{w,g}(C)(p) \Leftrightarrow p \text{ is the most likely member of } C . p$
 133 (10) [NOT [even_{NPI} [Kim had invited Sam to the party]]]
 134 - Most-likely inference: *That Kim had invited Sam to the party* is the most likely proposition of the
 135 form *that Kim had invited X to the party*.

136 The prejacent of even_{NPI} in (10) is a clausal constituent without negation. The scalar inference will
 137 therefore target a positive proposition, which will yield a most-likely inference since the presuppositional
 138 requirement of even_{NPI} demands of *p* to be the most likely element in *C*. On the ambiguity theory, then,
 139 the inferences attributed to *even* in various environments are determined by a lexical specification of
 140 NPI-hood rather than by scope. Even_{NPI} appears in contexts where NPIs are licensed, and is specified

³ See ([Wilkinson 1996](#), [Lahiri 1998](#), [Guerzoni 2004](#), a.o.). A complete description says that *even* needs to outscope *all* downward-entailing operators, not just negation. Consider example (i) where *even* is inside a conditional. An in situ treatment of *even* makes the wrong (and contradictory) prediction that Mary is less likely to notice one mistake than she is to notice multiple mistakes. In fact, we infer that she is *most-likely* to infer one mistake over many.

(i) If Mary noticed even one_F mistake of yours, it would be a problem. ([Guerzoni 2004](#))
 - Mary noticing one mistake is most-likely compared to noticing multiple mistakes.

On the scope theory, the most-likely inference in (i) is explained if *even* moves outside the scope of the conditional.

⁴ Since the fine mechanics of focus prosody and focus semantics are not central to our paper we abstract away from the details here and refer the reader to e.g. ([Rooth 1996](#)) etc.

⁵ Notation: $\llbracket \alpha \rrbracket = \phi . \psi$ states that the semantic value of α is defined only if ϕ and when defined $\llbracket \alpha \rrbracket = \psi$.

⁶ See ([Rullmann 1997](#), [2007](#), a.o.).

141 to trigger a most-likely inference. In a context where an NPI would not be licensed, regular *even* is used,
 142 giving rise to a least-likely inference.

143 Both of these theories can successfully analyze the examples we have seen so far, and much continued
 144 debate attempts to distinguish them. Of interest here is that they differ in which parts of the grammar are
 145 responsible for a given inference. On the scope theory, examples like (1) carry an unambiguous least-likely
 146 inference because *even* always has a least-likely presupposition. On the ambiguity theory, examples like (1)
 147 are unambiguous due to an additional component of the grammar, namely a constraint on the distribution
 148 of NPIs.

149 The unambiguous most-likely inference in (2) can likewise be accounted for by both proposals. However,
 150 in this case, both proposals depend on an additional component of the grammar, extrinsic to the lexical
 151 specification of *even*, to rule out ambiguity. Relying solely on the basic meaning of *even* as a sentential
 152 operator, which requires only that *even* combines with a propositional node, makes a scope position below
 153 negation (a propositional operator in its own right) in principle suitable on both theories. Thus, both theories
 154 need to appeal to a mechanism extrinsic to the lexical specification of *even* that prevents a least-likely
 155 inference to surface.

156 Anticipating the format of our experimental material, we illustrate this point with a sentence that employs
 157 *even* in pre-subject position to the left of negation, (11). Both theories can explain the attested most-likely-
 158 to/least-likely-to-not inference by assuming that (regular) *even* scopes above negation, (11a). However,
 159 both theories also need to explain why a logical form where (regular) *even* scopes below negation, which
 160 would give rise to a least-likely-to inference, is not attested, (11b).⁷

- 161 (11) Even AMY wasn't invited to the party.
 162 a. [Even [NOT [AMY was invited to the party]]]
 163 b. * [NOT [even [AMY was invited to the party]]]
 164 – Scope Theory: violates scope constraint for *even*
 165 – Ambiguity Theory: Blocked by *even_{NPI}*

166 On the scope theory (11b) is ruled out by a rather specific prohibition against interpreting *even* below
 167 clause-mate negation. On the ambiguity theory, a blocking principle of some sort is required to ensure that
 168 the availability of *even_{NPI}* preempts regular *even* from being inserted in this environment (e.g. because
 169 more specified lexical items, *even_{NPI}* in this case, are prioritized by principles of vocabulary insertion
 170 over less specified lexical items, non-NPI *even*, along the lines of Halle and Marantz (1993)).

171 2.2 Predicting the acquisition profile of even

172 Given two possible inferences (least-likely and most-likely) and two types of environments (positive and
 173 negative, or more generally, upward and downward entailing⁸), we might expect a learner to consider four
 174 possible ways to use *even*. A learner who hypothesizes that *even* is polysemous between a most/least-likely
 175 inference may start with all four use-patterns in the space (c.f. Giannakidou (2007) on Greek) (Table 1A).
 176 A learner who entertains only a least-likely inference associated with *even* (along the lines of the scope
 177 theory), however, should only consider three at the outset (Table 1B).

⁷ The precise mechanism (e.g. some form of reconstruction) by which such a logical form might be generated is not important to our purpose.

⁸ In this paper, we abstract way from the interpretation of *even* in non-monotonic contexts, see e.g. Crnić (2014).

178 The acquisition path that each of these hypothetical learners takes is expected to be different because
 179 they have different starting points. When confronted with a positive *even*-sentence, (12), a learner who
 180 hypothesizes a uniform least-likely inference is expected to appear adult-like, because they should only
 181 detect a least-likely inference in this context. By contrast, a learner who entertains a polysemous *even*
 182 should find sentences like (12) ambiguous between a most-likely/least-likely interpretation. The profile of
 183 this type of learner is therefore expected to be non-uniform – at-times adult-like and at-times non-adult-like
 184 depending on how they choose to resolve the ambiguity.

185 (12) Even AMY was invited to the party.

186 – Uniform least-likely *even*: adult-like

187 – Polysemous *even*: adult-like or non-adult-like depending ambiguity resolution

188 (13) Even AMY wasn't invited to the party.

189 – Uniform least-likely *even*: adult-like or non-adult-like depending on scope

190 – Polysemous *even*: adult-like or non-adult-like depending on ambiguity resolution

191 When confronted with a negative *even* sentence, (13), however, both learners are predicted to detect an
 192 ambiguous most-likely/least-likely inference. For the learner who posits a uniform least-likely inference,
 193 this is because *even* has two in principle available scope positions with respect to negation. For the
 194 polysemous *even* learner, they always posit these two inferences regardless of scope.

195 For a uniform least-likely *even* learner, getting to the adult-like pattern, then, amounts to ruling out
 196 ambiguity in negative environments, presumably by expunging the low scope option for *even* in those
 197 contexts. For a polysemous *even* learner, they must rule out ambiguity in both polarity environments to
 198 arrive at the adult pattern. And they must do so in a specific way: in positive environments the most-likely
 199 variant of *even* must be expunged while in negative environments the least-likely variant is targeted.⁹

200 These different learning profiles should be detectable in a well-controlled acquisition study as soon as a
 201 child has the machinery to reason about likelihood, and identifies *even* as a linguistic device used to convey
 202 likelihood inferences. To our knowledge, however, this has not been successfully investigated. Rather,
 203 existing acquisition studies of *even* aim to address a coarser question: *do children know even?* The findings
 204 of these studies have been taken to indicate that children struggle so much with *even*'s focus and scalar
 205 properties, that it is unreasonable to think one could ever glean the specific shape of the child's hypothesis
 206 space for the likelihood inferences associated with *even*.

207 We will argue, however, that these conclusions are not well supported by the experimental record.
 208 Additionally, we will show that controlling for previously uncontrolled experimental factors reveals a more
 209 intricate structure of the development of the grammar and learning path of *even*. In particular, we uncover
 210 evidence that children as young as 4 years of age *do* entertain a space like those in Table 1, and moreover
 211 that they start with a space that resembles the ambiguity theory of *even*.

212 2.3 Previous work: Kim (2011)

213 Motivation for our study, as well as inspiration for its design, is due in part to results from Kim (2011),
 214 which to our knowledge is the first systematic investigation of *even* in positive and negative environments
 215 in child language. Kim framed her interest in *even* within the context of prior studies on the acquisition

⁹ There is hypothetically a third possibility, which is that a child first hypothesizes a uniform most-likely inference associated with *even*. The learning path of a learner with this profile is unpredictable because the child would effectively have to start over at some point by adding to their hypothesis space.

216 of scalar implicatures. Echoing the consensus view of work in that domain, Kim argues that children as
217 old as 5 are essentially ignorant about *even*. However, we find her results to be ambiguous between two
218 interpretations: 1) children ages 4-5 do not detect *even*'s scalar inferences, or 2) children ages 4-5 *do* detect
219 *even*'s scalar inferences, but nevertheless exhibit non-adult-like behavior that is invited by their developing
220 grammar of *even*.

221 [Kim \(2011\)](#) conducted a comprehension experiment employing a forced-choice task, in which the
222 experimenter tells children stories about 3 characters who are all different sizes. In each story, all of the
223 characters are supposed to do a task, which scales in difficulty with the size of the character. At the end of
224 the stories, either all of the characters succeed or they all fail. The end of the story is accompanied by a
225 prompt of the form, *Even X was/n't able to do Y*. With the help of a puppet, the children are asked to point
226 to X. Figure I shows a sample story setup about bears reaching for cookies on a shelf, with both possible
227 story outcomes.

228 Kim tested 30 children on three positive/negative story pairs of this sort, yielding 90 responses, where each
229 'response' corresponds to a pair of responses to the positive/negative versions of each story (Table 2). There
230 were 3 distinct response profiles corresponding to a given positive/negative pair. Some were completely
231 adult-like ("target characters for both sentence types"), some gave opposite of adult-like responses by
232 choosing the tallest character in positive environments and the shortest character in negative environments
233 ("opposite characters for both sentence types"), and some uniformly chose either the rightmost or leftmost
234 character, regardless of polarity ("always rightmost or leftmost character"). Here "both sentence types"
235 refers to both positive and negative prompts (i.e. *even X was able to do Y/even X wasn't able to do Y*).

236 In addition, Kim ran a control version of the experiment in which *even* was removed from the prompts.
237 No other changes were made to the stories or the task in the control experiment. For example, in the bear
238 story the prompt might be *Larry was (not) able to reach the cookie*. The control study has the interesting
239 property that an adult would presumably struggle to find a felicitous answer to the question, "Who is
240 Larry?", given that "Larry was (not) able to reach the cookie" is equally true of every character. Despite
241 the pragmatic oddity of the task, and the lack of an "adult-like" target answer, Kim shows that children
242 performed very similarly in the control experiment as they did in the main experiment with *even* (Table 2).

243 The fact that children showed low rates of adult-like responses in the *even* experiment, as well as the fact
244 their responses didn't change substantially when *even* was removed in the control experiment, led Kim
245 to conclude that children essentially ignore *even* when participating in these comprehension tasks. She
246 therefore concludes that the children in this age group do not understand *even*.

247 We think, however, that her results are compatible with an alternative explanation and so are not
248 compelled to accept her conclusion. Notice that in her main experiment, none of the children chose the
249 middle character. All response profiles included one or both of the extrema, but never the middle. This is
250 unexpected if their selection were truly random. While it is possible to conclude with Kim that this is an
251 accidental result, or a product of the pragmatics of the task, in an unpublished manuscript, Kenyon Branan
252 argues that this pattern could also reflect a feature of the developing grammar. In other words, it might
253 be the case that children don't choose the middle character because they know that *even* triggers either a
254 least-likely or a most-likely inference, neither of which supports picking the middle character. What they
255 don't know, on this conjecture, are the grammatical conditions that control in the adult grammar which
256 likelihood inference is triggered in which environment.

257 Looking closer at Kim's design and procedures reinforces being more cautious in interpreting her results.
258 Note first that Kim's child participants were asked to point out both extrema characters but never the middle

259 character during the story leading up to the *even* sentence. This introduces a potential bias in favor of the
260 extrema which makes it difficult to assess whether this gap should be seen as an experimental artifact or as
261 a reflection of a non-adult-like grammar.

262 Second, the specific choices Kim made in the design of her materials introduce a potential confound.
263 Concretely, her experiment employed three story types, 2 of which involved reaching tasks, and one
264 of which was a lifting task. For both reaching and lifting stories, the likelihood of success is *directly*
265 proportional to the size of the character. Larger characters are both taller and stronger, and are thus, all else
266 being equal, more likely to succeed at reaching something tall or lifting something heavy. An adult-like
267 response to these story types thus amounts to choosing the smallest character in positive environments
268 (least likely to succeed), and the largest character in negative environments (least likely to not succeed).
269 However, this set up introduces a confound for children who showed a preference for the rightmost or
270 the leftmost character. Such a preference could be interpreted in multiple ways. One interpretation is that
271 children have an irrelevant preference for either the smallest or the largest character (or the rightmost or
272 leftmost character). This interpretation is compatible with Kim's conclusion that children ignore *even*
273 completely when doing the task.

274 However, this behavior is also explained if children are accessing likelihood inferences that are not
275 detectable in adult language, but are present in their hypothesis space for *even*. If for any given example,
276 a child is aware that both a most- and a least-likely inference is in principle available, the child might
277 sometimes choose the most-likely character to succeed, and sometimes the least-likely character to succeed,
278 irrespective of the polarity of the sentence. Choosing either the most- or least-likely character uniformly
279 in both polarity environments would amount to a right-most or left-most preference, in which case this
280 response pattern should not be treated as evidence of naivety. In sum, then, Kim's observations are amenable
281 to two quite different explanations and additional work is required to decide which one is on the right track.

3 COMPREHENSION EXPERIMENT 1

282 Our own comprehension experiments adopt Kim's basic setup, which we think is quite elegant and offers a
283 very natural way of testing children's comprehension of *even*. We are, however, implementing a number of
284 modifications to help overcome the aforementioned limitations in interpreting her results.

285 3.1 Methods

286 Experiment 1 adopts the basic format of Kim's set-up, but includes the modifications and amendments
287 summarized below. Examples of stimuli used in Experiment 1 can be seen in Figure 2

- 288 1. All characters are equally prominent in the story dialogue to avoid biasing children towards the extrema.
289 Unlike Kim's dialogue, which singles out the extrema characters ("Who is the shortest [bear]? Who
290 is the tallest [bear]?"'), our dialogue gives equal weight to all three characters ("This story is about 3
291 squirrel brothers. There's a little one, a medium one, and a big one.")
- 292 2. We include filler stimuli without *even* where the middle character is the correct answer, both to test
293 children's alertness and to give them an opportunity to see that the middle character is a possible
294 answer (Supplementary Material; Figure S1). The filler stimuli involve matching a character to an
295 object based on attributes like color and size.
- 296 3. Children are asked to justify their responses (e.g. "How did you know that was Larry?")

297 4. There are two additional story types in which likelihood to succeed scales *inversely* with the size of the
 298 character (i.e. fitting and filling stories). This distinguishes a preference for the smallest character from
 299 a preference for the least-likely character.¹⁰

300 5. The age range of our study is 3-6yo to look for a possible developmental trajectory for *even* (in contrast
 301 with Kim's range of 4-5yo).

302 We recruited 91 English-speaking children ages 3;1-6;11 (mean = 5;0, age in years;months format) from
 303 Boston-area daycares, preschools, and through the Living Laboratory program at the Museum of Science,
 304 Boston. Three subjects were excluded from the analysis either because they did not complete the task, or
 305 because their responses on the control items (as well as justifications) suggested they were not actively
 306 participating in the task.

307 In total, the experiment consisted of 4 filler stimuli and 8 target stimuli: 4 positive target stimuli (1 each
 308 for reach/lift/fit/fill) and 4 negative target stimuli (1 each for reach/lift/fit/fill). The experimental items were
 309 blocked by polarity of the *even* sentences (NEG vs. POS) and presented in one of two orders (NEG-first
 310 vs. POS-first). Children who were assigned to the NEG-first order heard all of the negative stimuli before
 311 hearing any positive stimuli, and vice versa.

312 As a control, we ran a version of this experiment with 68 adult subjects on Mechanical Turk. We slightly
 313 modified the stimuli from the child study for use on IBEX (Drummond, 2012), creating an introductory
 314 slide which introduced the characters and the situation, and a question slide, which introduced the *even*
 315 sentence and asked the participant to identify the named character. Participants were given 10 seconds to
 316 respond, starting immediately after the question slide was displayed. IBEX recorded question responses,
 317 reaction times, as well as justifications. We followed two exclusion criteria: 1) we excluded participants
 318 who incorrectly answered more than one filler item, 2) we excluded participants who gave the same answer
 319 (e.g. the middle character) throughout the entire study. After these exclusions, 56 subjects remained.

320 Data were analysed using the MCMCglmm package in R (Hadfield, 2010) with a mixed-effects
 321 multinomial logistic regression. We took response type as the dependent variable, and modeled the
 322 fixed effect of polarity and age group {3yo, 4yo, 5yo, 6yo}. Additionally, we modeled a random intercept
 323 by subject and by story type. We opted to use a Bayesian approach to our data analysis for two main reasons.
 324 First, Bayesian approaches make it relatively easy to specify hierarchical models, such as the multinomial
 325 mixed-effects model deployed here, as compared to frequentist data analysis. Second, Bayesian models
 326 with maximal random effects structures converge with less data than frequentist models. 4 chains were
 327 generated per model, and convergence was tested across these chains using the Gelman-Rubin diagnostic
 328 (Gelman et al., 1992; Plummer et al., 2006), as well as by visual inspection of the posterior distributions.
 329 Credible Intervals were calculated at 95% with Highest Posterior Density intervals.¹¹¹²¹³¹⁴

¹⁰ A keen reader might notice that the filling story in Figure 2 is in fact ambiguous between two possible likelihood inferences: 1) likelihood of the basket to be filled up (smallest basket = most-likely) vs. 2) likelihood of the character to collect enough for their basket (largest character = most-likely). Indeed we will see that children and adults alike are sensitive to this ambiguity, which added noise to our initial results, and prompted a change to these stimuli in Experiment 2.

¹¹ For further background on the use of Bayesian methods in the context of linguistics, see Nicenboim and Vasishth (2016).

¹² Credible Intervals, used in Bayesian analysis, are intervals which contain a certain percentage (in this paper, 95%) of the values in the posterior distribution for an unobserved parameter. For instance, a 95% Credible Interval of [10%, 22%] indicates that we can be 95% certain that the value of a parameter falls between 10% and 22%. In this paper, we follow a basic decision rule to take the 95% Credible Interval as an indication of whether a given value can be inferred to be a possible "true" value of a parameter (Kruschke et al., 2012). Thus, the CI here serves a similar role to the *p*-value used in frequentist statistics, in giving us a criterion to reject or accept the null hypothesis.

¹³ Highest Posterior Density is a method for selecting Credible Intervals which, as might be expected, involves selecting the highest density intervals within the posterior distribution. Note that because HPD intervals are chosen by density of the posterior distribution, they will *not* center around the mean.

¹⁴ All code used for the analyses in this paper can be found here: <https://github.com/MITLanguageAcquisitionLab/even>

330 3.2 Results

331 In our original experimental design, polarity (NEG vs. POS) was a within-subject condition and order
332 (NEG-first vs. POS-first) was a between-subject factor. However, our models failed to converge with order
333 as a factor. Because of this, we chose to only take data from the first four target items from each subject,
334 removing order as a condition, and turning polarity into a between-subject factor. This way we can be sure
335 that our reported results on the interaction between response-type and polarity is not affected by possible
336 within experiment learning or other order effects throughout the experiment.¹⁵

337 The results of our experiment differ from Kim's in several ways. The children that we tested offered three
338 types of responses, which we call *adult-like*, *middle*, and *opposite*, corresponding to which character they
339 chose. While Kim's results had no middle responses, in our study, subjects ages 3-5 did sometimes choose
340 the middle character for target items. We also see a steady increase in the number of adult-like responses
341 across age, and a stable number of opposite responses. Figure 3 summarizes the results of Experiment 1.

342 Figure 3B shows the rate of adult-like responses across age, split by polarity. What we find is that this rate
343 increases much more quickly in negative environments than in positive environments. At age 4, there is a
344 clear difference between the positive and negative environments with respect to adult-like responses ([−3%,
345 74%]). The profile of *even* comprehension in negative environments is quite stable across ages 4-6 (4yo:
346 [56%, 98%], 5yo: [56%, 97%], 6yo: [63%, 99%]) while adult-like behavior lags in positive environments
347 at age 4, rising steadily until age 6 (4yo: [12%, 74%], 5yo: [48%, 100%], 6yo: [73%, 100%]).

348 A complementary trend can be seen in the rate of middle responses (Figure 3C), which given that they are
349 never invited by the grammar are an indication of confusion. We see that the number of middle responses
350 remains low, decreasing over time (3yo: [0%, 42%], 4yo: [1%, 27%], 5yo: [0%, 30%], 6yo: [0%, 3%]).

351 Looking at the last response type, opposite responses, reveals a different pattern from both adult-like and
352 middle responses. Opposite responses are stable across the four age groups, hovering at approximately
353 20-25% (Figure 3D) (3yo: [0%, 49%], 4yo: [5%, 38%] 5yo: [0%, 17%], 6yo: [0%, 22%]). Additionally,
354 we do not see a polarity effect for opposite responses.

355 We also asked children to justify their answers and coded their responses as *scalar*, *random* or *none* to
356 reflect whether their justifications contained evidence of scalar reasoning. In general, all justifications that
357 referred to the size of the characters were coded as “scalar”, and all responses were coded jointly by the
358 two first authors. Some sample justifications that we coded as “scalar” include: “Because it’s rare that a
359 tiny thing can lift a big thing”, “teeny one”, “small mouses can usually fit”, “because it’s the biggest”, etc.

360 Children also often provided reasons that did not reference a discernible scale, which we coded as
361 “random”. Some sample “random” justifications include: “look at the pink bunny!”, “because I just knew it”,
362 “he’s two (years old)”, “that one has a little bow”, etc. Some children were too shy to offer a justification, or
363 stated that they didn’t know why they chose a particular character, in which case we coded their responses
364 as “none”.

¹⁵ Though our experimental design included filler items for the purpose of potentially excluding inattentive participants, a problem with one of them prevented us from using these filler items as grounds for excluding subjects in our analysis of Experiment 1. Our original exclusion criteria would have excluded children who answered more than 1/4 filler items incorrectly (allowing for occasional but not systematic lapses in judgment). However, our child participants systematically struggled with one of the filler items, even when it was clear that they were attentive throughout the experiment. It was therefore clear that we should remove that item from our consideration for exclusion. However, doing so increased the likelihood of exclusion from 2/4 to 2/3, which we felt was too significant a difference to implement without probable cause. We therefore chose to be maximally inclusive of all subjects who completed the study.

365 Table 3 summarizes the number of justifications of each type that were provided for each response pattern.
 366 Notably, scalar justifications were predominantly offered for adult-like and opposite responses, while
 367 random justifications were the most frequent justification type for middle responses.

368 The overall results of the adult control study appear in Table 4. The adults, as expected, predominantly
 369 pick the least-likely character in positive environments (88%) and the most-likely character in negative
 370 environments (82%). This shows a slight asymmetry in favor of positive stimuli – however, it is not
 371 statistically significant.¹⁶

372 3.3 Analysis and discussion

373 We want to draw attention to several features of these results, some of which resolve open questions in
 374 Kim's experiment, and some of which raise new ones.

375 First, we argue that these results refute Kim's conclusion that children ages 4-5 show no evidence of
 376 comprehending *even*. Our results show a clear upward trajectory in the rate of adult-like responses between
 377 the ages 3-6, as seen in Figure 3A.

378 Furthermore, our results provide a clear indication as to when children begin to reliably exhibit sensitivity
 379 to the scalar properties of *even*. Given that there are three possible response types, chance behavior should
 380 be 33%. While 3yo do not give adult-like responses at an above chance rate ([17%, 83%]), starting at 4
 381 years old, children provide adult-like responses at well above chance, nearing adult-like levels by age 6
 382 (4yo: [45%, 89%], 5yo: [60%, 100%], 6yo: [76%, 100%]).

383 An analysis of the mean difference between the rate of adult-like, middle and opposite responses across
 384 ages shows that by age 6, the rate of adult-like responses has increased significantly since ages 3-4, while
 385 the rate of middle responses has decreased significantly (Table 5A). These results reinforce our claim that
 386 children progressively acquire an adult-like understanding of *even*.

387 Looking just at the progression of adult-like responses in negative contexts, we see that the only significant
 388 difference is between the 3yo on the one hand, and the 4-, 5-, and 6yo on the other (Table 5B). There is
 389 no significant difference between the rate at which 4yo provide adult-like responses to negative stimuli
 390 and the rate at which 5- and 6yo do, suggesting that children acquire negative *even* sentences earlier than
 391 positive ones. This result contradicts the expectation that negation adds a computational cost. Not only are
 392 children unphased by the addition of negation, but they appear to have an easier time deducing the right
 393 inference in its presence.

394 Looking at adults in the control study, we do not see a polarity effect. There is no significant difference in
 395 rates of different response types between positive and negative stimuli.

396 Adults and children both show two types of error behavior, with a higher rate of opposite responses. In
 397 adults, the rate of middle responses is extremely low, but the rate of opposite responses is nearly as high as
 398 the rate of opposite responses in 6yo.¹⁷

399 The contrasting behavior of middle vs. opposite responses in both adults and children seems to indicate
 400 that the learning trajectory for *even* contains a stage at which elements of the grammar of *even* are already
 401 in place, but are not yet fully adult-like. The fact that middle responses disappear over time but opposite

¹⁶ A similar asymmetry in the same direction can be seen in terms of response times, with negative stimuli taking longer than positive stimuli, but this difference is again not statistically significant.

¹⁷ Just as in the child study, we asked adults to justify their responses. Unsurprisingly since almost all errors were opposite responses, all of the justifications were scalar except for two participants that chose the middle character and wrote *Just a random guess* and *I don't know* respectively.

402 responses are stable through age 6 suggests that opposite responses are principled at this stage, while
403 middle responses are not. In sum, we propose to analyze middle responses as indications of genuine
404 confusion, while opposite responses should be analyzed as licensed by the developing grammar of *even*.
405 Moreover, given that they are detectable at the earliest stages where children exhibit a stable appreciation
406 of *even*, we propose to analyze them as an indication of what the initial hypothesis space for *even* looks
407 like, Table 1A.

408 The decrease in middle responses over time correlates with the increase in adult-like behavior, which
409 supports our view of middle responses as a measure of confusion. As children become more adult-like, they
410 become less confused. We suggest that by contrast, opposite responses are not an indicator of confusion
411 and thus do not decrease noticeably as adult-like behavior increases. Notice that they persist to some extent
412 even in adults.

413 This characterization of the middle and opposite responses is supported by the justifications that children
414 provided for each response type. Modeling the rate of justification types by response type shows that
415 children provide significantly more scalar justifications to adult-like and opposite responses as compared to
416 middle responses, as seen in Table 5C. Furthermore, there is no significant difference in the rate of scalar
417 justifications between adult-like and opposite responses. The stability of opposite responses throughout our
418 above-chance-performing age ranges (4-6yos), combined with how well-reasoned their justifications are,
419 suggests that children know that *even* is associated with a space like Table 1A.

420 Recall from Section 2 that the in principle space of inferences associated with *even* on the *adult* grammar
421 included opposite inferences, which motivated the notion of the hypothesis space explored here. This
422 prediction is readily borne out in our adult study; adults essentially only make errors in the form of opposite
423 responses, which is expected given their grammar. The exciting finding from our child study is that children
424 also favor this error behavior over other potential error patterns, suggesting that they too access a space
425 like Table 1A. That said, our results contain some unexplained noise that merits further discussion. We
426 therefore conducted a second comprehension study, whose main purpose was to investigate the potential
427 sources of noise and their impact on our results.

428 One major source of noise in our data was the amount of variation in adult-like behavior across the
429 different story types. In particular, the filling stories were accompanied by more non-adult-like responses
430 than the other stories (Supplementary Material; Figures S2, S3).

431 A clue for this pattern comes from the justifications provided by children and adults. Specifically, our
432 subjects frequently indicated that there was a salient alternative interpretation of the filling stories. On this
433 alternative interpretation, apparent “opposite” responses are actually adult-like, suggesting that our initial
434 reported rates of error responses were actually somewhat inflated.

435 In the filling stories, the smallest character should have the easiest time filling their basket because
436 their basket is the smallest. On an alternative construal, however, the smallest character might have the
437 most difficulty collecting the requisite number of acorns (because of their age, or general inexperience)
438 necessary to fill their basket. Several adult and child subjects interpreted the filling stories on the latter
439 characterization, providing well-reasoned justifications like those in (14).¹⁸

440 (14) *Adult-like justifications for opposite responses in filling stories*

441 a. “The first squirrel is Sammy because the youngest one would probably have the hardest time
442 filling his basket and I assume the smallest squirrel to be the youngest.” (adult)

¹⁸ Note that one of the children used *even* in an adult-like manner to justify the opposite response (14p).

443 b. “because even though his basket is little he can even fit a lot of acorns in there” (child)

444 Experiment 2 therefore includes new stimuli to replace the ambiguous filling stories, in order to eliminate
445 this source of noise.

446 An additional concern about our experimental design pertains to the role of the ability modal in the
447 prompts. Recall that our test subjects were asked to respond to prompts of the form *Even X_F was/n’t able*
448 *to do Y*. This element of the design was inherited from Kim’s experiment for the purpose of attempting to
449 replicate her results. However, now that we have failed to replicate her results, we return to this feature of
450 the design with some scrutiny. Of interest is the fact that ability is itself a gradable notion. To interpret a
451 sentence with an ability modal, one must therefore have access to a kind of scalar reasoning.^[19] A question
452 we should ask, then, is whether the overt expression of ability accounts for any of the children’s behavior
453 on its own. To investigate the relevance of this possible confound, we removed any overt scalar notions
454 from our stimuli in Experiment 2.^[20]

4 COMPREHENSION EXPERIMENT 2

455 4.1 Methods

456 Experiment 2 contains two large scale changes to the stimuli from Experiment 1: the removal of the
457 ability modal in the prompts, (15), and the replacement of the filling stories with *spilling* stories (Figure 4).

- 458 (15) Sample prompts from Experiment 2 (ability modal removed)
- 459 a. Even Benny_F has gotten an apple!
- 460 b. Even Franky_F hasn’t gotten the socks on!

461 The spilling stories are about likelihood to spill/drop something heavy. Like the fitting stories, height is
462 inversely proportional to likelihood to spill. The smallest character is the most likely to be the weakest, and
463 therefore the most likely to drop or spill the heavy bucket/basket.

464 These two changes combined are expected to avoid the problems of the previous filling stories because
465 they remove the ambiguity about whether the likelihood scale should refer to the abilities of the characters
466 vs. properties of their baskets/cups.

467 In order to make our data comparable to Experiment 1, we made polarity a between factor. Each subject
468 evaluated eight target stimuli and four filler items, where all target stimuli were of the same polarity
469 environment. This doubles the number of data points for each polarity condition.

470 We presented these stimuli to subjects in either of two orders, one in which the story types proceeded
471 as *fit*, *reach*, *lift*, *spill*, and the other in reverse order. We collected data from 80 children, ages 3;1-5;11
472 (mean = 4;7). Unlike in Experiment 1 (see fn. 15), where children systematically struggled with a particular
473 filler item, no such difficulty was apparent in any of the modified filler items in Experiment 2. Thus, we
474 were able to use performance on these filler items as an exclusion criterion in Experiment 2. Subjects who

^[19] See Greenberg (2015), anticipated by Rullmann (2007), for an account of *even* that replaces the likelihood scale with a scale introduced by a (contextually salient) degree predicate.

^[20] In fact, we already have some reason to believe that the scalar properties of *be able* do not play a role in opposite responses. First, opposite responses are present in our adult data. While it is possible that some adults ignore *even* entirely in our experiment, this seems less likely. Additionally, we attempted to run a version of Experiment 1 with children where we did not pronounce *even* (e.g. *Jessiepillar was able to reach the shelf*). We ultimately stopped the experiment because most of the children became very confused and didn’t want to finish the task. This suggests that children who provide opposite responses are indeed sensitive to the presence of *even*.

475 failed to correctly answer at least 3 out of 4 filler questions were excluded from our results (14 exclusions),
476 as were those whose justifications and behavior during the study suggested that they were not paying
477 attention to the task (4 exclusions). After exclusions, 62 subjects remained. The remaining data is evenly
478 distributed by polarity, block order, and age group, with five or six subjects per age group/polarity/order
479 combination. We focused on these earlier years because two interesting effects apparent in Experiment
480 1 either disappeared or stabilized by age 6: 1) the stable preference for opposite responses over middle
481 responses, and 2) the polarity asymmetry in the rate of adult-like responses.

482 As with Experiment 1, we performed an adult control study with 85 participants on Mechanical Turk
483 with amended stimuli on IBEX (see section 3.1 for information on the changes made). After exclusion
484 criteria were applied, data from 60 participants remained.

485 Data for the child and adult experiments were again analysed using a mixed-effects multinomial logistic
486 regression. We modeled the fixed effect of polarity and age group as well as order (forward and reverse).
487 Additionally, we modeled a random intercept by subject and a random intercept and slope by story type.
488 See Section 3.1 for additional information on the statistical analysis.

489 4.2 Results

490 Figure 5A summarizes the rate of each response type by age for Experiment 2. The rate of adult-like
491 responses again increases steadily with age (3yo: [8%, 62%], 4yo: [53%, 91%], 5yo: [79%, 100%]), while
492 middle responses decrease steadily with age (3yo: [11%, 81%], 4yo: [0%, 22%], 5yo: [0%, 3%]). We also
493 see a fairly stable population of opposite responses (3yo: [3%, 37%], 4yo: [5%, 33%], 5yo: [0%, 18%]).

494 Looking at each response type individually by polarity reveals that the polarity asymmetry from
495 Experiment 1 is far less pronounced. The rate of adult-like responses is roughly equal between the
496 two polarity conditions in 4yos, which was where the asymmetry was most pronounced in Experiment 1.
497 5yos do perform noticeably better in negative environments than positive environments, but this result is
498 not statistically significant (Figure 5B). The same pattern holds for middle responses (Figure 5C). There is
499 no polarity asymmetry at ages 3 or 4, but a slight asymmetry becomes visible at age 5. Opposite responses
500 show no sensitivity to polarity, and are even more stable than in Experiment 1 (Figure 5D).

501 Finally, the justifications for each response type, coded as *scalar/random/none*, are summarized in Table
502 6. As in Experiment 1, scalar justifications were offered far more for adult-like and opposite responses than
503 for middle responses. Justifications for middle responses were most often random.

504 A comparison of adult-like behavior across each of the story types shows that there is less variation across
505 items compared to Experiment 1 for both adults and children (Supplementary Material; Figures S4, S5).
506 Most notably, the profile of responses is no longer substantially different for spilling stories than for the
507 other story types.

508 The results of the adult control study are summarized in Table 7. As in Experiment 1, adults were slightly
509 less error-prone in positive environments.²¹ However, this asymmetry is, again, not statistically significant.

510 4.3 Analysis and discussion

511 Experiment 2 confirms the main finding from Experiment 1, namely that children begin to comprehend
512 *even* earlier than Kim (2011) suggests. Once again, 3yos do not choose the adult-like response at an above
513 chance rate [8%, 62%], while 4yos [53%, 91%] and 5yos [79%, 100%] perform at well above chance.

²¹ As in the adult control version of Experiment 1, there is also a non-statistically significant polarity asymmetry in terms of reaction time.

514 Looking at mean differences between age groups, 4- and 5-year-olds give adult-like responses at a
515 significantly higher rate than 3-year-olds (Table 8A). However, opposite responses remain steady throughout
516 age groups.

517 The opposite responses show essentially the same profile in Experiment 2 as they did in Experiment 1.
518 Many of the children gave opposite responses, including the older ones, and these were often accompanied
519 by scalar justifications. The middle responses, on the other hand, showed some sensitivity to age and
520 polarity, and were primarily given random justifications, Table 8B.

521 Experiment 2 was successful in that it replicated these results from Experiment 1 with substantially
522 less noise. The justifications indicate that both our adult and child participants interpreted *even* primarily
523 based on the size of the characters and not based on any less obvious criteria such as perceived age or
524 general competence. Figures S4 and S5 in the Supplementary Material show that story type also no longer
525 correlates with any particular error pattern.

526 Experiment 2 also revealed that the polarity asymmetry in child comprehension observed in Experiment 1
527 was detectable in Experiment 2 only as a trend that did not reach statistical significance. 5yos, in particular,
528 exhibited an advantage for negative environments on adult-like comprehension similar to that observed in
529 Experiment 1 for 4yos. Determining what might be responsible for the similarities and differences of this
530 environmental effect across the two experiments is unclear to us and deserves further investigation.

531 Interestingly, opposite responses in our adult control exhibit a noticeable sensitivity to polarity.
532 Specifically, opposite responses are observed primarily in negative environments while positive
533 environments generated vanishingly few.²² After filtering out answers associated with justifications that
534 indicated guessing or another interpretation of the story, there were 10 opposite responses observed in
535 negative environments, but only 2 opposite responses in positive environments.²³²⁴

536 To summarize, both Experiments 1 and 2 have enabled us to identify two distinct error profiles in adults
537 and children. Middle responses indicate simple confusion or inattention. Opposite responses appear to be
538 licensed by the grammar.

539 An important difference between adult and child error behavior is the polarity sensitivity of opposite
540 responses. Children offered opposite responses in *both* positive and negative environments, suggesting that
541 they access a space of inferences like that in Table 1A. By contrast, adults basically only offered opposite
542 responses in *negative* contexts, suggesting that they access a space like that in Table 1B.

5 CORPUS STUDIES

543 In order to better understand children's experience with and usage of *even*, we conducted two corpus
544 studies, in which we compare the features of our stimuli to tokens of *even* found in child and adult corpora.
545 Our investigation resulted in three relevant findings:

546 1. Children produce *even* essentially adult-like (error-free) as early as 3 years of age, hence much earlier
547 than they comprehend *even* in an adult-like manner.

²² Attempts at modelling this contrast with a multinomial mixed-effects model unfortunately failed to converge.

²³ Our scoring of justifications was charitable towards adult-like interpretations or true guesses so as to minimize the chance of inflating the rate of opposite responses.

²⁴ Recall that the rate of opposite responses in Experiment 1 was inflated due to the ambiguity of the filling stories. Reanalysis of the results from Experiment 1 that accounts for this inflation in fact demonstrates a similar polarity asymmetry: opposite responses are primarily concentrated in negative environments (3 in POS vs. 17 in NEG).

- 548 2. Adults show a use-bias for *even* in negative environments (here, ‘use-bias’ refers to a bias towards using
 549 a word in a particular environment). Children, by contrast, initially favor *even* in positive environments
 550 and acquire the adult use-bias for negative *even* in the ages of 4-5.
 551 3. The form of stimuli in our comprehension study instantiate a *low*-frequency use pattern for *even* in
 552 child and child-directed speech.

553 To study the distribution of *even* in both child-produced speech and child-directed speech²⁵, we examined
 554 token instances of *even* in the American English CHILDES corpus (MacWhinney, 2000).²⁶ Data were
 555 analyzed using childe-coder (Gowda, 2020; Sanchez et al., 2018), which presents instances of target
 556 words along with their broader contexts, and provides an interface for users to save metadata about these
 557 instances in a database (Supplementary Material; Figure S6). Instances of *even* produced by speakers
 558 (marked Target_Child for child-produced speech, and marked Mother, Father, Adult, Uncle, Grandmother,
 559 Aunt, Grandfather, Family_Friend, or Teacher for child-directed speech) were coded by mutual agreement
 560 among the authors for several criteria, including:

- 561 1. Presence of negation: {Yes, No, Unclear}
 562 2. Order of negation and *even*: {N/A, *even*-NEG, NEG-*even*, Unclear}
 563 3. Whether negation is sentential: {Yes, No, Unclear}
 564 4. Order of *even* and the subject: {*even*-Subj, Subj-*even*, Unclear}
 565 5. The focus associate of *even*: {Subject, Object, Verb, Adjunct, Unclear}
 566 6. The likelihood inference: {Most-Likely, Least-Likely, Unclear}

567 In addition to examining the surrounding context of each instance of *even* to determine the correct coding,
 568 we also made use of the metadata and audio recordings available in CHILDES.

569 5.1 Results

570 Our comprehension studies focused on the interaction between polarity and the likelihood inference in
 571 *even* sentences. Because we coded tokens of *even* in production by both polarity and the contextually salient
 572 likelihood inference, we can similarly investigate this interaction in production. Results for child-produced
 573 speech are summarized in the first four rows of Table 9. Results for child-directed adult speech are in the
 574 last row.

575 In Table 9, for both adults and children, the rates of “opposite” uses of *even* – most-likely inferences
 576 for *even* in positive environments and least-likely inferences in negative environments – are so low
 577 that they are essentially absent from our findings. This is quite striking given that we see robust rates of
 578 opposite responses in comprehension and suggests that opposite responses are, in fact, a comprehension
 579 phenomenon.²⁷

580 Another salient pattern we observe in both child and child-directed corpus data is the overall prevalence
 581 of negative-*even* sentences (Table 10). Children start off using *even* in positive environments more often
 582 than in negative environments at age 3. According to a chi-squared test for homogeneity, the distribution of

²⁵ Here, we use ‘child-directed speech’ to refer to all utterances produced by non-child speakers in the CHILDES corpora. Thus, this data includes all speech in CHILDES that a child was exposed to, not necessarily just speech that was specifically directed at a child.

²⁶ Due to metadata consistency issues, data from the MacWhinney corpus was excluded from the analysis of both child-produced and child-directed speech.

²⁷ In actuality, because our coding schema kept track of the polarity of the environment rather than upward/downward entailment, there were higher numbers of most-likely inferences in apparent positive environments, as indicated in Table 9 by the gray figures. However, these are not true opposite responses because most-likely inferences are predicted by the grammar in downward entailing environments more generally, whether or not there is sentential negation. Indeed, a study of most-likely responses in positive sentences for adults shows that they all involve downward entailing environments.

583 positive and negative uses of *even* in 3yos is significantly different from adults ($p < .001$), while there is
 584 no significant difference ($p > .1$) between the distribution of positive and negative *even* sentences in 4, 5,
 585 and 6yos and adults. That is, by age 4, children appear to have an adult-like use-bias for *even* in negative
 586 environments.

587 Lastly, we investigated several features of our comprehension study stimuli in the corpus and found that
 588 our stimuli instantiate a low-frequency usage pattern for *even*. Our stimuli can be described according to
 589 the following feature specification:

- 590 1. Focus associate = subject
 591 2. *Even* precedes the subject
 592 3. *Even* precedes negation

593 We parametrically compared our stimuli to other *even* constructions by coding for whether a given
 594 utterance with *even* focused the subject, whether it contained sentential negation, and the linear order of
 595 *even*, the focus associate, and negation (if applicable). Table S1 in the Supplementary Material shows that
 596 subject focus was relatively infrequent compared to VP or object focus. When the subject was focused,
 597 however, pre-subject *even* was preferred to post-subject *even*.

598 In addition to the placement of the subject with respect to *even*, we can compare the placement of *even*
 599 with respect to negation (Supplementary Material; Table S2). *Even* follows negation in the majority of
 600 cases. When this requirement interacts with the previous tendency for subject-associating *even* to appear
 601 pre-subject, we see a preference for maintaining both constraints, resulting in *not even*. Examples with
 602 both subject focus and sentential negation were therefore most often of the form in (16b). Our stimuli, by
 603 contrast, took the form in (16a)²⁸.

- 604 (16) a. Even Linda_F didn't write to me. (dispreferred)
 605 b. Not even Linda_F wrote to me. (preferred)

606 5.2 Analysis and Discussion

607 In our corpus studies, children as young as 3yo produced *even* as if they were adults.²⁹ Taking their
 608 behavior at face value, and assuming Snyder (2007)'s principle of Grammatical Conservatism, this should
 609 not be possible unless they have identified a grammatical basis for *even*'s scalar inferences.

610 (17) **Grammatical Conservatism:** Children do not begin making productive use of a new grammatical
 611 construction in their spontaneous speech until they have both determined that the construction is
 612 permitted in the adult grammar, and identified the adult's grammatical basis for it. (Snyder 2007)

613 The fact that they exhibit non-adult-like behavior in comprehension should therefore not indicate that
 614 they lack a grammar entirely, as is argued by (Ito, 2012; Kim, 2011, e.g.). Indeed we argue that children
 615 must have a grammar of sorts for *even*, which happens to invite at times non-adult-like behavior in
 616 comprehension. Moreover, children must have some appreciation for which cells in the space of inferences
 617 allowed by their grammar are not available to adults, or else they would not be so adult-like in production.
 618 We will explore this in Section 6.

²⁸ Table S1 in the Supplementary Material contains all tokens of *even* in which any type of negation was present; sentential or constituent. *Even-neg* order was mostly available for clauses with constituent rather than sentential negation.

²⁹ Although in this paper we only present CHILDES data from ages 3-6, adult-like production of *even* is apparent as early as 2 years old, with 32 instances of POS *even* sentences and 16 instances of NEG *even* sentences.

619 This conclusion is supported by another finding from our corpus studies, namely the fact that our
620 comprehension task stimuli instantiate a low-frequency use pattern for *even* in both child and adult speech.
621 Children's lack of experience with our *even* sentences, however, apparently did not deter them. Children
622 ages 4-6 still offered adult-like responses as the dominant response pattern in our comprehension studies.
623 We must therefore conclude that their command of *even* is quite sophisticated. They are able to abstract
624 away from the particular *even* sentences that they hear and use most frequently, and generalize to other
625 less-frequent uses.

626 Lastly, our corpus studies offered a perspective on the polarity asymmetry observed in our comprehension
627 studies. Adults apparently produce *even* approximately 1.5 times more often in negative environments than
628 positive environments, which translates into a comparative abundance of negative *even* in children's input.
629 We propose that children show (slightly) better rates of comprehension in negative environments because
630 they have more experience trying to interpret negative *even*. This effect is, however, less pronounced than
631 our other findings because of the natural tension between inherent knowledge and experience. The proposed
632 child grammar of *even* affords children the same abstract knowledge of *even* in negative as well as positive
633 environments, which accounts for their overall adult-like competence. Their performance with *even* in real
634 time, however, can be marginally affected by their confidence with each polarity environment.³⁰

6 LEARNING SCALAR INFERENCES

635 The previous sections presented novel findings from a series of comprehension and corpus studies which
636 significantly enrich our knowledge of the empirical landscape of how *even* is acquired. The picture that
637 emerges reveals a surprisingly intricate developmental path.

638 Our evidence from the comprehension studies suggests that children as young as 4 years of age
639 systematically draw scalar inferences that are comparable in nature to those generated by the adult
640 grammar of *even*. Interestingly, their knowledge of the environmental conditions controlling when to draw
641 which type of inference (least-likely-to in positive environments and most-likely-to/least-likely-to-not in
642 negative environments) is not completely adult-like at this age. This gives rise to a comprehension behavior
643 that runs at times directly opposite to that of adult speakers.

644 With regard to the basic effect of polarity on the nature of the scalar inference, we saw a rather striking
645 lack of opposite scalar inferences in the corpus data. Even occurrences of *even* produced by 3yo conform to
646 the adult pattern. This is quite surprising since we know from the comprehension experiments that opposite
647 inferences are allowed by their developing grammar as late as age 6. Furthermore, children ages 4-6 were
648 adult-like in production to the extent that they showed an adult-like use-bias for negative *even*, indicating
649 that they are also sensitive to the conversational settings in which *even* sentences are predominantly used.

650 To explain why the acquisition of *even* unfolds along such an intricate path is a non-trivial task. It
651 involves specifying the initial hypothesis space that the learners start out with as well as the final state that
652 characterizes the adult grammar of *even*. Furthermore, learners must identify the relevant evidence as well
653 as develop or adopt strategies that together enable them to transform the former cognitive structure into the
654 latter. Though we are not yet in a position to offer a full account of the acquisition of *even* that lives up to
655 all of these demands, we do think that our findings allow us to make significant progress towards that goal.

³⁰ Recall that this effect is fragile in our experiments (i.e. statistically present in Experiment 1 but only present as a trend in Experiment 2). This fragility might be due to a rather fine-grained effect of experience. It is conceivable that the absence of the ability modal and the presence of the present perfect make our stimuli in Experiment 2 even less well-represented in the input than the stimuli in Experiment 1. If this is true, it would thereby delay the advantage of negative environments to a later age when enough of such cases have been encountered.

656 6.1 Initial hypothesis space for even

657 The findings from our comprehension studies provide persuasive information about the nature of the
658 initial hypothesis space: it needs to allow for all four combinations of likelihood inferences and polarity
659 of the environment to be expressible by *even*, Table 1A. A simple way of implementing such a grammar
660 would be to postulate a polysemous *even* which can freely occur in positive and negative environments.

661 Our argument in support of this conclusion is straightforward. We found both adult-like and non-adult-like
662 scalar inferences in the earliest stages of comprehending *even* in both positive and negative environments.
663 In other words, we saw that all four cells in Table 1A are utilized as soon as learners start to appreciate the
664 scalar nature of *even*. Importantly, we saw that all four inference patterns occurred stably throughout an
665 extended period of learning, and we found them to be regularly accompanied by reasoned justifications that
666 referenced the relevant scale properties. This shows that the non-adult-like responses (just like the adult-like
667 responses) were sanctioned by the developing grammar. Thus, they should be analyzed as exemplars that
668 are predicted by the initial hypothesis space rather than as errors whose source is unrelated to the grammar
669 of the learner.

670 6.2 Adult grammar of even

671 From the perspective of the two competing views on the adult grammar described in Section 2, the
672 initial hypothesis space in Table 1A can be described as an as yet unconstrained form of the grammar
673 predicted by the ambiguity theory. Recall that the main tenet of the ambiguity theory is that *even* can in
674 principle carry a least-likely as well as a most-likely inference. The distribution of these variants needs
675 to be constrained via the addition of a grammatical feature (in this case an NPI-feature on most-likely
676 *even*). Without that addition, the distribution of *even* would remain unconstrained allowing for all logically
677 possible combinations to be realized, Table 1A.

678 An as yet unconstrained version of the grammar predicted by the scope theory, by contrast, is not a
679 viable option for the initial hypothesis space since it makes only three of the four required cells available,
680 Table 1B. The reason is, again, transparent. The main tenet of the scope theory is that *even* can only
681 carry a least-likely inference. In combination with negation we can generate a most-likely inference if
682 *even* out-scopes negation since the resulting least-likely-to-not inference is equivalent to a most-likely-to
683 inference. However, without the presence of negation, the scope theory can only generate a least-likely
684 inference, which leaves the most-likely inferences in positive environments unaccounted for.

685 Though we have identified (on empirical grounds) a greater similarity between the initial hypothesis
686 space for *even* and the ambiguity theory, it would be unjustified to conclude at this point that the ambiguity
687 theory has to be correct for the adult grammar of *even*. Both theories of adult *even* are, in fact, compatible
688 with the initial state postulated for children. They simply require different transformations on the initial
689 state. Deciding which one offers the better account for the adult grammar, therefore, depends on what
690 the actual steps are that allow learners to transform the initial hypothesis space (Table 1A) into either the
691 constrained version of Table 1A or Table 1B.

692 To arrive at an ambiguity grammar of adult *even*, the learner must replace the polysemous *even* with two
693 separate *evens*, each specified for a particular likelihood inference and level of polarity sensitivity (NPI or
694 unmarked). Under the scope theory, learning amounts to eliminating the possibility of *even* triggering a
695 most-likely inference altogether. This is arguably a simpler transformation on the initial hypothesis space.
696 However, it yields the target grammar only at the expense of adopting an unprecedented constraint on the
697 syntactic scope of *even*.

698 Importantly, both theories also have to explain why other logically possible combinations of likelihood
 699 inferences and polarity-sensitivity/scope constraint are never selected by learners of English *even*. Below
 700 we argue that a plausible source to rule out unattested combinations is the limited conversational utility
 701 of those combinations. Interestingly, these considerations will also provide us with a possible account of
 702 why acquiring the adult grammar takes relatively long and what might be responsible for the puzzling
 703 production-comprehension asymmetry we have observed.

704 6.3 Pragmatics of likelihood-inferences

705 Throughout the paper we have described the scalar inferences triggered by *even* in terms of likelihoods –
 706 “least-likely-to” in positive environments and “most-likely-to/least-likely-to-not” in negative environments.
 707 In doing so, we adopted the terminology of Karttunen and Peters (1979) which is intuitive and sufficiently
 708 transparent to characterize the differences between the various scalar inferences we have encountered.
 709 Whether likelihood is (always) the correct way to characterize the dimension of the scalar inferences of
 710 *even* is, however, debated in the literature. Alternatives include various formulations of expectedness,
 711 noteworthiness, informativity as well as scales introduced by gradable predicates.³¹ We cannot provide
 712 a full assessment of this debate here. Instead, our aim is to clarify the connection between the inferred
 713 relative likelihood of a proposition and its noteworthiness in a given conversational situation. This will be
 714 sufficient to diagnose the conversational status of the adult-like as well as the opposite scalar inferences we
 715 have observed.

716 Taking a closer look at the opposite inferences preschoolers draw in our comprehension studies, we
 717 observe that they are not simply absent from the adult grammar of *even*, but are in fact conversationally
 718 odd if we try to render their content anyway. Compare, for example, (18), which features the content of the
 719 adult-like scalar inferences via an appositive relative clause, to its rather odd sounding counterpart in (19),
 720 which features the “opposite” content of *even*’s inferences.

721 (18) Adult-like inferences

- 722 a. Everybody has reached the book, including Jessiepillar, who was the least likely to have done
 723 so.
- 724 b. Nobody has reached the book, including Jessiepillar, who was the most likely to have done so.

725 (19) Opposite inferences

- 726 a. # Everybody has reached the book, including Jessiepillar, who was the *most likely* to have
 727 done so.
- 728 b. # Nobody has reached the book, including Jessiepillar, who was the *least likely* to have done
 729 so.

730 We propose that a pragmatic explanation of this contrast can provide insight into the factors that help
 731 learners constrain their initial hypothesis space for *even*. To see how, let us examine the conversational
 732 context in which our *even* sentences were uttered.

733 Recall that the three characters in our stories either all succeeded at the relevant task or they all failed.
 734 This fact was highlighted explicitly with a universal statement immediately preceding the *even* sentence.
 735 The truth-conditional content of the *even* sentence was therefore redundant, which puts the burden to

³¹ See Fillmore (1965); Fauconnier (1975); Kay (1990); Rullmann (2007); Giannakidou (2007); Greenberg (2016), among others. Refinements of, or alternatives to likelihood-based characterization can be envisioned that are consistent with our findings. Choosing among them is, however, not topical for us here.

736 provide conversational utility for the utterance squarely on its not-at-issue content. A question we might
 737 ask is whether a pragmatic requirement on conversational utility constrains the space of possible likelihood
 738 inferences at all.

739 Given the oddity of (19), we argue that it does. Moreover we propose that this oddity is derived if a
 740 connection between a character's likelihood of success with some notion of propositional noteworthiness is
 741 important to adult-like competence with *even*. To see why, it is important to consider the context in which
 742 the *even* sentence occurs very carefully.

743 After the universal statement but before the *even* sentence is uttered, the context contains a proposition of
 744 the form, *Every x in C is such that x has reached the book*. If the modal horizon against which the likelihood
 745 inference were evaluated contained all and only the verifying situations characterized by the universal
 746 statement, the likelihood inference would be moot. Comparing the relative likelihood of *Jessiepillar has*
 747 *reached the book* to *Some other x in C has reached the book* is almost nonsensical because both sentences
 748 are already true in the context.

749 What allows the likelihood inference to be meaningful and useful is to consider a context in which
 750 it was *not* given that any character would reach a book. In other words, for the *even* sentence to have
 751 a sensible inference, it must take as its context variable a set of propositions that does *not* contain the
 752 universal statement that preceded the *even* sentence in our experiment. But what drives this move? Why
 753 does a listener bother to accommodate a different common ground in which to make sense of the likelihood
 754 inference, rather than just ignore it? We propose that this move is connected to a notion of propositional
 755 noteworthiness.

756 We propose that when an adult listener hears an *even* sentence in our context, they detect its redundancy
 757 and ask, what makes this proposition worth repeating? I.e. what makes this result surprising or deserving
 758 of comment? It is this question that allows the listener to specify a useful common ground in which to
 759 consider the likelihood inference. The logic goes as follows: in order for “Jessiepillar has reached the book”
 760 to be noteworthy, it must be unexpected. Jessiepillar must therefore be the character whose success was
 761 least likely.

762 Notice, however, that there is no way to connect to a most-likely inference on this logic, hence capturing
 763 the oddness of (19). Considering characters who were *likely* to succeed in no way explains why emphasizing
 764 those characters' success is interesting.

765 If this is the type of pragmatic reasoning that adults employ generally in a conversation, opposite responses
 766 are predicted to be infelicitous, thus providing insight into why the grammar of *even* is constrained to just
 767 two adult-like inferences. Therefore, the content that our child comprehenders end up with when they
 768 select the opposite character is odd from the perspective of the adult grammar.

769 To clarify, for Jessiepillar to be singled out even though her height doesn't justify it, as happens when
 770 children provide opposite responses, is of course logically possible.³² However, for adult speakers this
 771 requires a different “backstory”, e.g. Jessiepillar might be the most/least motivated of the three to do what
 772 it takes to reach the book making her the most/least likely to succeed/fail despite her height. Our stories
 773 did not provide any useful information about the characters other than their height, however. Thus, adult
 774 comprehenders are stuck with Jessiepillar's height as the only available basis for anchoring the likelihood

³² Our cases are not characterized by entailment relations between alternatives. In such cases, opposite responses would only be possible at pains of accepting a contradiction. See e.g. Lahiri (1998).

775 inference triggered by *even*. They therefore pick the shortest character to be Jessiepillar when all candidates
 776 succeed and the tallest when all of them fail.

777 For our child comprehenders the situation is different. While they may be practiced enough
 778 conversationalists to consider a common ground in which a likelihood inference is meaningful, their
 779 grammar overgenerates. Because their grammar allows for a least-likely as well as a most-likely
 780 specification of the scalar inference regardless of the polarity of the sentence, they consider both the
 781 tallest and the shortest character as grammatically viable candidates for Jessiepillar. This is only possible,
 782 however, if we assume that children are less attuned to the relevant conversational pragmatics than adults
 783 are, thus allowing both options to remain viable in our conversational setting. Thus, they can in principle
 784 choose the character at either end of the scale.

785 With regard to why children are less adept at this kind of pragmatic reasoning than adults, a number
 786 of options seem plausible. For instance, it may be that “explicit” Theory of Mind level reasoning about
 787 the motivations of speakers, which is required to detect this sort of pragmatic oddness, is not yet fully
 788 developed or sufficiently practiced.³³ Alternatively, or additionally, it may be that children’s processing
 789 resources for this kind of reasoning are still not up to full capacity. Whatever the true underlying causes,
 790 it seems reasonable to characterize their pragmatic reasoning as more tolerant than those of adults
 791 towards not knowing exactly what the speaker had in mind when they issued their *even* sentence.³⁴
 792 A willingness to proceed in a state of partial ignorance leaves both inferences in play making both extrema
 793 live possibilities. Of course, this follows only if the developing grammar generates both types of inferences
 794 in both environments to begin with. On our proposal this is so because their grammar allows for all four
 795 cells of the initial hypothesis space to be expressible by *even*.³⁵

796 The fact that opposite choices occur at all is therefore predicted on the basis of children’s more limited
 797 conversational experience compared to adults. However, the fact that opposite choices occur less frequently
 798 than adult-like choices may be seen as a reflection of those choices being less optimal even from the
 799 perspective of the learner. After all, these opposite choices require a willingness to proceed without having
 800 figured out exactly what the speaker meant with their *even* sentence.³⁶

801 Turning to the question of how children actually acquire the adult grammar of *even*, the following
 802 picture emerges: learning, under the present view, is a function of becoming more adept at recognizing
 803 conversational goals and more intolerant when those goals are not identified during comprehension. In other
 804 words, the more pragmatically skilled a learner is, the better they will be at recognizing and recording the
 805 specific conversational setting in which *even* sentences are used. This growing conversational confidence
 806 favors the adult grammar, which does not support opposite inferences, and correspondingly discriminates
 807 against a grammar that does support opposite inferences. Eventually, the absence of evidence in favor of

³³ Though it is now widely accepted that some aspects of Theory of Mind reasoning, often called “implicit” Theory of Mind inferences, are in place earlier than our age range (cf. Onishi and Baillargeon (2005)) it is plausible that the relevant skills in our task include assessing a speaker’s conversational assumptions and goals, which has been argued to come online much later in development and hence are likely to scale with the mount of practice young reasoners have, see e.g. Perner and Roessler (2012).

³⁴ See Katsos and Bishop (2011) for a similar notion that children are pragmatically more tolerant than adults.

³⁵ We have evidence from the adult control studies that the adult grammar, by contrast, does not provide access to all four cells. Rather, the adult error pattern we have observed is compatible with only three of the four cells in the space (Table I B). This potentially indicates that the adult grammar is most like the scope theory, since only one type of error - least-likely-to inferences in negative environments - occurred. Most-likely-to inferences in positive environments never occurred, which is straightforwardly predicted by the scope theory. The ambiguity theory, by contrast, rules them out by a mechanism – licensing of NPI – that is known to be error prone during processing (Drenhaus et al., 2005).

³⁶ If we assume that whenever they proceed in a state of ignorance they are guessing which character they should pick our observed rate of opposite responses of 25% translates into a rate 50% of not being able to figure out what exactly the speaker had in mind. To assess whether this is a reasonable estimate would, however, require a more fleshed out theory of what makes pragmatic reasoning of this sort difficult for our learner.

808 opposite inferences is taken by the learner to suggest an adjustment to the grammar to ensure that sentences
 809 that would generate opposite inferences are no longer generated to begin with.³⁷

810 Turning to the question from Section 6.2, i.e. why learners do not consider other logically possible adult
 811 grammars with different combinations of likelihood inferences and polarity, two factors emerge: 1) no data
 812 from the input supports such grammars, and 2) the learner’s own pragmatic knowledge discourages them.

813 Last but not least, the present perspective also allows us to sketch a plausible account of the production-
 814 comprehension asymmetry. Recall that we never see children use *even* in a non-adult way, even at the
 815 earliest stages. Specifically, they underutilize the opposite inferences that their grammar apparently licenses.

816 A key difference between production and comprehension is that the speaker knows what that intended
 817 message is, while the comprehender has to figure out what the message is that the speaker intended.³⁸ Our
 818 account allows us to exploit this difference directly: we proposed that opposite inferences surface during
 819 comprehension when the listener is unable to figure out what the speaker’s conversational goals might be
 820 (e.g. why they singled out Jessiepillar), but is nevertheless willing to go along with the task at hand (in our
 821 case selecting one of the three characters).³⁹⁴⁰ Importantly, we did not require that children not appreciate
 822 the connection between likelihood and expectedness. Indeed, the fact they they mostly interpret our target
 823 sentences in an adult way suggests that they do, just not as reliably as adults. This means that we can
 824 reasonably assume that when they issue an *even* sentence they do so with the intent to convey information
 825 that rides on the intuitive connection between likelihood and expectedness.

7 CONCLUSION

826 This paper has advanced a view of the scalar inferences associated with *even* in child grammar that stands
 827 in contrast with much prior literature. Previous acquisition studies of *even* and scalar inferences treated
 828 children’s non-adult-like behavior as evidence of simple confusion about *even* and likelihood inferences in
 829 discourse. Our studies present evidence to the contrary. Children are not simply confused. They are, in fact,
 830 rather keenly sensitive to the scalar nature of *even* and generate robustly both least-likely and most-likely
 831 inferences very early on. They are non-adult-like only in that they exhibit more tolerance to uses of *even*
 832 during comprehension where the speaker’s conversational goals carried by *even* are left unresolved.

833 On our view, there is nothing difficult per se about detecting scalar inferences. As soon as children learn
 834 to associate them with the particle *even*, they immediately access a relevant hypothesis space of scalar
 835 inferences associated with *even* along the lines of Rooth (1985)’s ambiguity theory (Table 1A). This is
 836 evident from their “error” patterns in our comprehension studies, as well as the absence of such errors in
 837 production.

838 We presented two comprehension studies that shared Kim (2011)’s “Guess who?” format. In these
 839 experiments, children as young as 4yo performed well above chance. Most notably, children in our age

³⁷ Our proposal is part of the growing literature on the role of pragmatics in language acquisition which has uncovered a great deal of pragmatic sophistication that young learners bring to the task of language acquisition in variety of different situations, ranging from referential word learning (e.g. Horowitz and Frank 2015; Sullivan and Barner 2015; Sullivan et al. 2019) to speech act pragmatics and its role in the acquisition of propositional attitudes (Hacquard and Lidz 2018). What the precise relation of our proposal is to the type of pragmatic reasoning in these cases is not immediately obvious since expunging a grammatical option that is underutilized because speakers tend to not highlight likely outcomes bears little resemblance to determining when and how a novel word is used by a speaker to refer to a novel object or the cases Hacquard and Lidz (2018)’s “pragmatic syntactic bootstrapping hypothesis” is meant to account for.

³⁸ See especially Hendriks (2014) for discussion of production-comprehension asymmetries in language acquisition.

³⁹ See Aravind (2018) for evidence that children in our age range who are otherwise quite astute at picking up the status of presupposed information as contextually entailed nevertheless prefer a fully redundant reading of an utterance over one that is informative only at pains of accommodation.

⁴⁰ The error pattern we observed in our adult comprehenders – opposite responses in negative environments – can be understood in a similar vein if we assume, as seems plausible, that our web-based task environment invites shallow processing as a form of satisficing, (Ferreira, 2003). Shallow processing will generate scalar inferences for *even* but may not require a full reconstruction of the (imagined) speaker’s conversational purpose in singling out Jessiepillar.

840 range predominantly chose either the most- or the least-likely character to succeed in any given story, but
841 rarely chose the middle character. They additionally justified their choices with comments that demonstrated
842 a sensitivity to scalar properties about the characters. This behavior is predicted by the space of inferences
843 in Table IA, but is not expected if children are merely guessing.

844 Further motivation for this treatment of child non-adult-like behavior comes from our adult control
845 studies. Adults likewise were occasionally susceptible to “opposite” likelihood inferences, and justified
846 these choices with normal reference to the scalar properties of the characters. We argued that this wouldn’t
847 be possible unless these inferences were made in principle available by the grammar.

848 In addition to the comprehension experiments, we conducted two corpus studies that examined tokens
849 of *even* in child and child-directed speech. Strikingly, neither adults nor children exhibited “opposite of
850 adult-like” uses of *even*. We argue on the basis of this production-comprehension asymmetry that children
851 not only hypothesize a space of inferences like that in Table IA, but they also command some of the
852 knowledge necessary to constrain this space (or else they wouldn’t be so adult-like in production).

853 Learning to transform the space in Table IA to the adult grammar amounts to making use of that
854 knowledge in comprehension as well as production. This is a gradual process of becoming increasingly
855 intolerant to certain inferences, as they become increasingly confident in their ability to reason about, and
856 identify, speakers’ conversational goals.

REFERENCES

- 857 Aravind, A. (2018). *Presuppositions in context*. Ph.D. thesis, MIT
- 858 Crnič, L. (2014). Non-monotonicity in NPI licensing. *Natural Language Semantics* 22, 169–217.
859 doi:10.1007/s11050-014-9104-6
- 860 Crnič, L. (2009). *Getting even*. Ph.D. thesis, MIT
- 861 Drenhaus, H., Frisch, S., and Saddy, D. (2005). Processing negative polarity items: When negation
862 comes through the backdoor. In *Linguistic Evidence* (Mouton de Gruyter). 145–164. doi:10.1515/
863 9783110197549.145
- 864 Drummond, A. (2012). Ibex: A web interface for psycholinguistic experiments. Available at: <https://github.com/addrummond/ibex> [Accessed April 9, 2018].
- 865
- 866 Fauconnier, G. (1975). Pragmatic scales and logical structure. *Linguistic Inquiry*, 353–375
- 867 Ferreira, F. (2003). The misinterpretation of noncanonical sentences. *Cognitive Psychology* 47, 164 – 203.
868 doi:10.1016/S0010-0285(03)00005-7
- 869 Fillmore, C. J. (1965). Entailment rules in semantic theory. *POLA report* 10
- 870 Francis, N. (2018). Presupposition-denying uses of *even*. In *Proceedings of SALT 28*, eds. S. Maspong,
871 B. Stefánsdóttir, K. Blake, and F. Davis (Ithaca, NY: CLC Publications), 161–176. doi:10.3765/salt.
872 v28i0.4409
- 873 Gast, V. and van der Auwera, J. (2011). Scalar additive operators in the languages of Europe. *Language*
874 87, 2–54. doi:10.1353/lan.2011.0008
- 875 Gelman, A., Rubin, D. B., et al. (1992). Inference from iterative simulation using multiple sequences.
876 *Statistical science* 7, 457–472. doi:10.1214/ss/1177011136
- 877 Giannakidou, A. (2007). The landscape of *even*. *Natural Language and Linguistic Theory* 25, 39–81.
878 doi:10.1007/s11049-006-9006-5
- 879 Gowda, Y. (2020). childe-coder: Utility for coding data from CHILDES corpora. Available at: <https://github.com/tlonic/childe-coder> [Accessed March 13, 2020].
- 880

- 881 Greenberg, Y. (2015). *Even*, comparative likelihood and gradability. In *Amsterdam Colloquium*, eds.
882 T. Brochhagen, F. Roelofsen, and N. Theiler. vol. 20, 147–156
- 883 Greenberg, Y. (2016). A novel problem for the likelihood-based semantics of *even*. *Semantics and*
884 *Pragmatics* 9, 1–28. doi:10.3765/sp.9.2
- 885 Greenberg, Y. (2018). A revised gradability semantics for *even*. *Natural Language Semantics* 26, 51–83.
886 doi:10.1007/s11050-017-9140-0
- 887 Guerzoni, E. (2004). *Even*-NPIs in yes/no questions. *Natural Language Semantics* 12, 319–343. doi:10.
888 1007/s11050-004-8739-0
- 889 Hacquard, V. and Lidz, J. (2018). Children’s attitude problems: Bootstrapping verb meaning from syntax
890 and pragmatics. *Mind & Language* 34, 73–96. doi:10.1111/mila.12192
- 891 Hadfield, J. D. (2010). MCMC methods for multi-response generalized linear mixed models: The
892 MCMCglmm R package. *Journal of Statistical Software* 33, 1–22. doi:10.18637/jss.v033.i02
- 893 Halle, M. and Marantz, A. (1993). Distributed Morphology and the pieces of inflection. In *The View from*
894 *building 20: essays in linguistics in honor of Sylvain Bromberger*, eds. K. L. Hale, S. J. Keyser, and
895 S. Bromberger (Cambridge, MA: MIT Press), no. 24 in Current studies in linguistics. 111–176
- 896 Hendriks, P. (2014). *Asymmetries between language production and comprehension*, vol. 42 of *Studies in*
897 *Theoretical Psycholinguistics* (Cambridge, Massachusetts: Springer). doi:10.1007/978-94-007-6901-4
- 898 Herburger, E. (2000). *What counts: focus and quantification* (Cambridge, Massachusetts: MIT Press).
899 doi:10.7551/mitpress/7201.001.0001
- 900 Horn, L. (1969). A presuppositional analysis of *only* and *even*. In *Chicago: Chicago Linguistic Society*.
901 vol. 5, 98–107
- 902 Horn, L. (1972). *On the semantic properties of logical operators in English*. Ph.D. thesis, University of
903 California: Los Angeles
- 904 Horn, L. (1989). *A natural history of negation* (University of Chicago Press)
- 905 Horowitz, A. C. and Frank, M. C. (2015). Young children’s developing sensitivity to discourse continuity
906 as a cue for inferring reference. *Journal of Experimental Child Psychology* 129, 84–97. doi:10.1016/j.jecp.
907 2014.08.003
- 908 Ito, M. (2012). Japanese-speaking children’s interpretation of sentences containing the focus particle
909 *datte* ‘even’: Conventional implicatures, QUD, and processing limitations. *Linguistics* 50. doi:10.1515/
910 ling-2012-0004
- 911 Karttunen, L. and Peters, S. (1979). Conventional implicature. In *Syntax and semantics*, eds. C.-K. Oh and
912 D. A. Dineen (Academic Press), vol. 11: Presupposition. 1–56
- 913 Katsos, N. and Bishop, D. V. (2011). Pragmatic tolerance: Implications for the acquisition of
914 informativeness and implicature. *Cognition* 120, 67–81. doi:10.1016/j.cognition.2011.02.015
- 915 Kay, P. (1990). *Even*. *Linguistics and Philosophy* 13, 59–111. doi:10.1007/BF00630517
- 916 Kim, S. (2011). *Focus particles at syntactic, semantic and pragmatic interfaces: The acquisition of only*
917 *and even in English*. Ph.D. thesis, University of Hawaii
- 918 Krifka, M. (1991). A compositional semantics for multiple focus constructions. In *Proceedings of SALT 1*,
919 eds. S. K. Moore and A. Z. Wyner (Ithaca, NY: CLC Publications), 127–158. doi:10.3765/salt.v1i0.2492
- 920 Kruschke, J. K., Aguinis, H., and Joo, H. (2012). The time has come: Bayesian methods for data
921 analysis in the organizational sciences. *Organizational Research Methods* 15, 722–752. doi:10.1177/
922 1094428112457829
- 923 Lahiri, U. (1998). Focus and negative polarity in Hindi. *Natural language semantics* 6, 57–123. doi:10.
924 1023/A:1008211808250

- 925 MacWhinney, B. (2000). *The CHILDES Project: Tools for analyzing talk. transcription format and*
926 *programs*, vol. 1 (Psychology Press)
- 927 Nicenboim, B. and Vasishth, S. (2016). Statistical methods for linguistic research: Foundational ideas—part
928 II. *Language and Linguistics Compass* 10, 591–613. doi:10.1111/lnc3.12207
- 929 Onishi, K. H. and Baillargeon, R. (2005). Do 15-month-old infants understand false beliefs? *Science* 308,
930 255–258. doi:10.1126/science.1107621
- 931 Perner, J. and Roessler, J. (2012). From infants to children's appreciation of belief. *Trends in Cognitive*
932 *Sciences* 16, 519 – 525. doi:10.1016/j.tics.2012.08.004
- 933 Plummer, M., Best, N., Cowles, K., and Vines, K. (2006). CODA: Convergence diagnosis and output
934 analysis for MCMC. *R News* 6, 7–11
- 935 Rooth, M. (1985). *Association with focus*. Ph.D. thesis, University of Massachusetts Amherst
- 936 Rooth, M. (1996). Focus. In *The Handbook of Contemporary Semantic Theory*, ed. S. Lappin (Oxford:
937 Blackwell Publishers). 271–297
- 938 Rullmann, H. (1997). Even, polarity, and scope. In *Papers in Experimental and Theoretical Linguistics*,
939 eds. M. Gibson, G. Wiebe, and G. Libben (University of Alberta), vol. 4. 40–64
- 940 Rullmann, H. (2007). What does even even mean? *Ms. University of British Columbia*
- 941 [Dataset] Sanchez, A., Meylan, S., Braginsky, M., MacDonald, K. E., Yurovsky, D., and Frank, M. C.
942 (2018). childe-db: a flexible and reproducible interface to the Child Language Data Exchange System.
943 doi:10.31234/osf.io/93mwx
- 944 Schwarz, B. (2005). Scalar additive particles in negative contexts. *Natural Language Semantics* 13,
945 125–168. doi:10.1007/s11050-004-2441-0
- 946 Snyder, W. (2007). *Child language: the parametric approach*. Oxford linguistics (Oxford University
947 Press)
- 948 Sullivan, J. and Barner, D. (2015). Discourse bootstrapping: preschoolers use linguistic discourse to learn
949 new words. *Developmental Science* 19, 63–75. doi:10.1111/desc.12289
- 950 Sullivan, J., Boucher, J., Kiefer, R. J., Williams, K., and Barner, D. (2019). Discourse coherence as a cue
951 to reference in word learning: Evidence for discourse bootstrapping. *Cognitive Science* 43, e12702.
952 doi:10.1111/cogs.12702
- 953 von Stechow, A. (1991). Current issues in the theory of focus. In *Semantics: An International Handbook*
954 *of Contemporary Research*, eds. A. von Stechow and D. Wunderlich (Walter deGruyter). 804–825
- 955 Wilkinson, K. (1996). The scope of even. *Natural Language Semantics* 4, 193–215. doi:10.1007/
956 BF00372819

FIGURE CAPTIONS

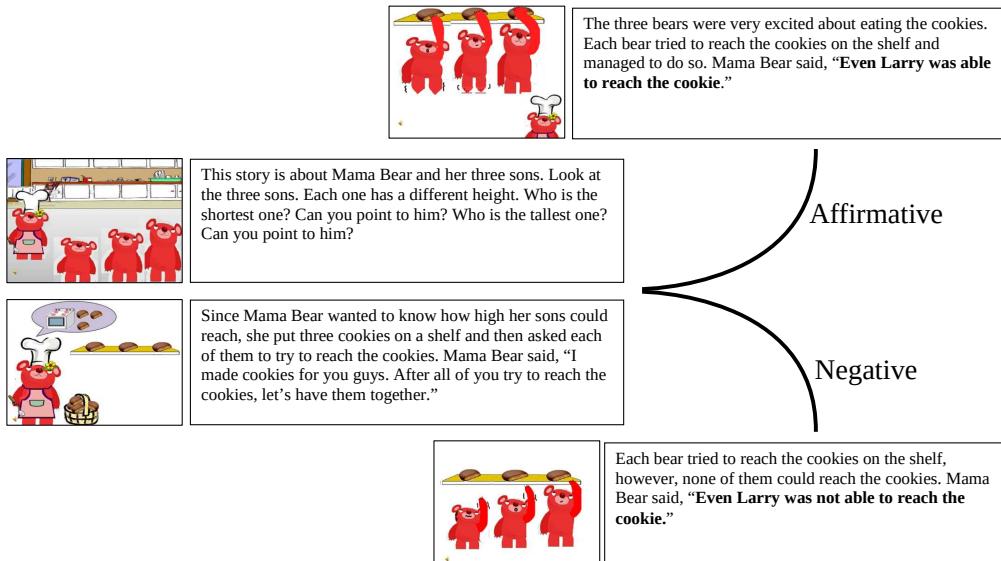


Figure 1. Sample story setup. (Kim, 2011)

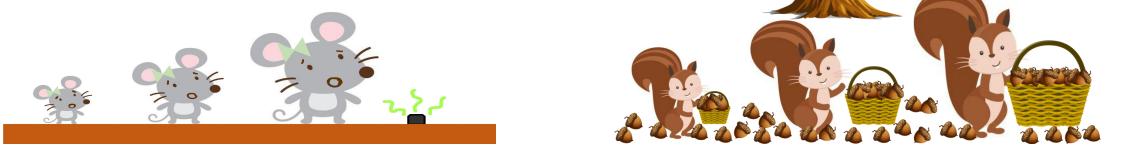
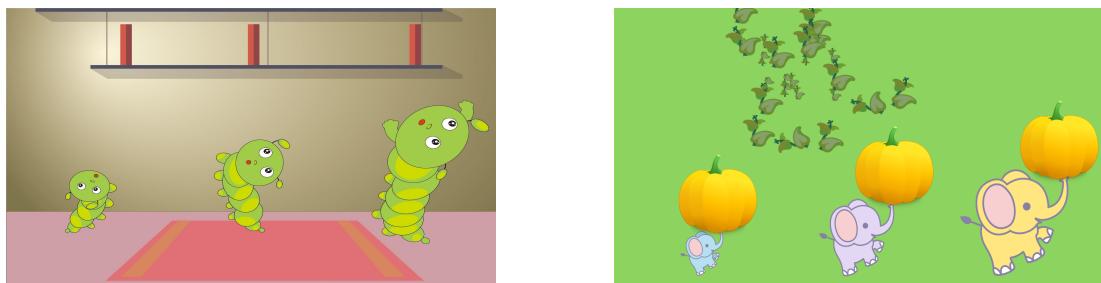


Figure 2. Four different scale types for the target stimuli: reach, lift, fit, fill.

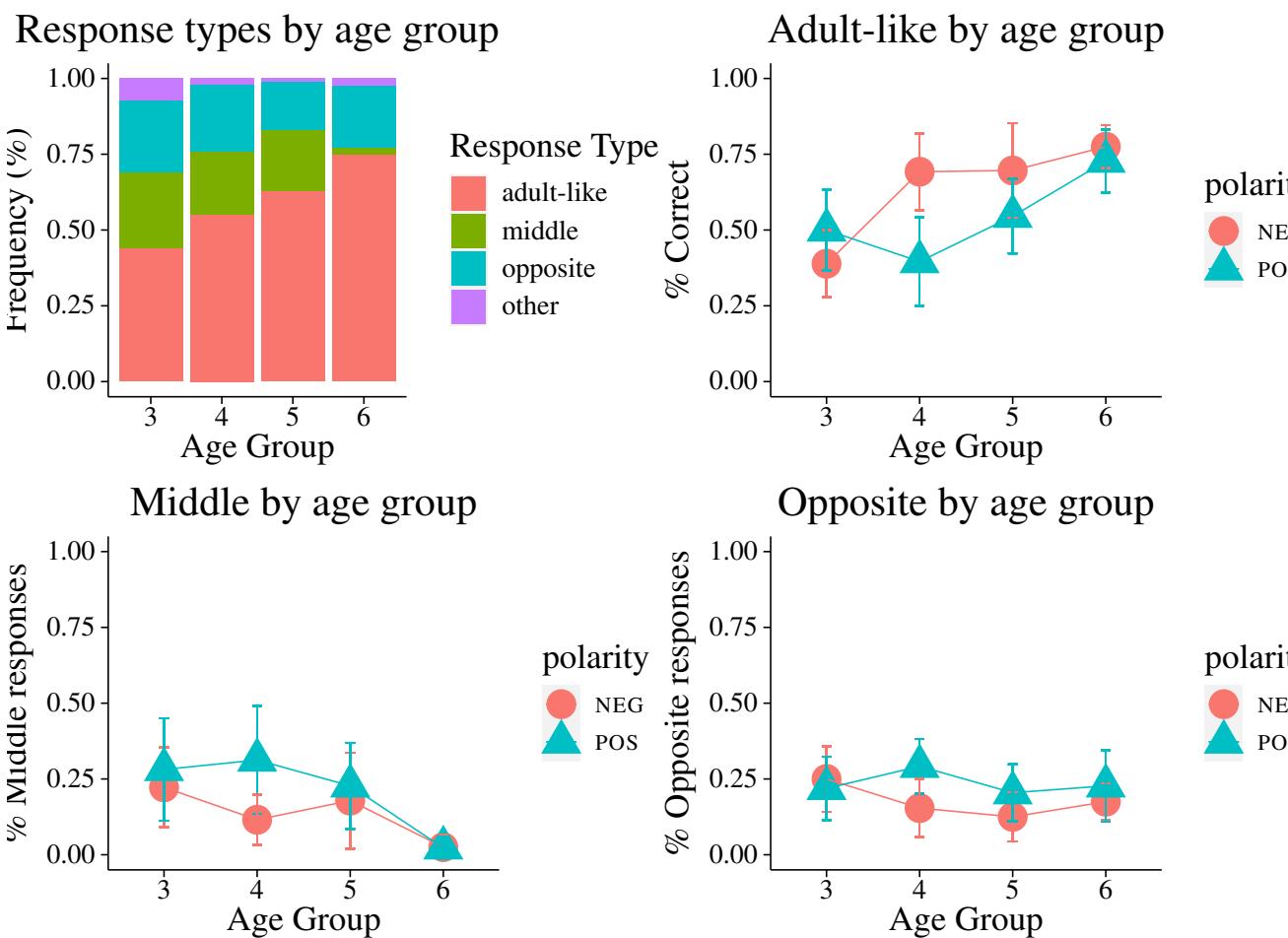
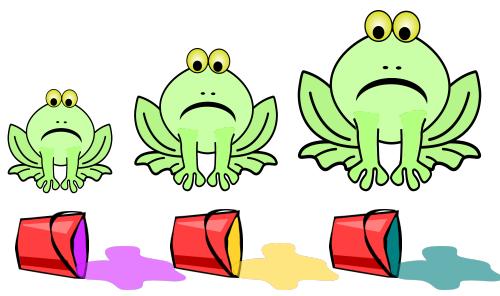
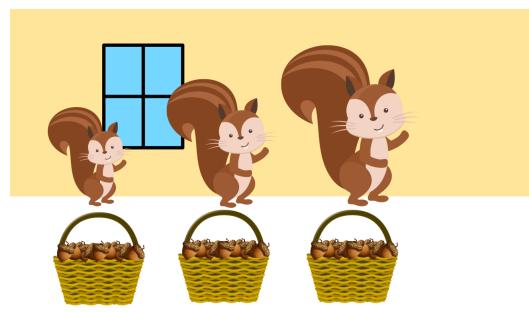


Figure 3. Results for Experiment 1. **(A)** Rate of response types by age group. **(B)** Rate of adult-like responses by age group and polarity. **(C)** Rate of middle responses by age group and polarity. **(D)** Rate of opposite responses by age group and polarity.



“Even Frida_F has spilled her paint!”



“Even Sammy_F hasn’t spilled her basket!”

Figure 4. Spilling stimuli

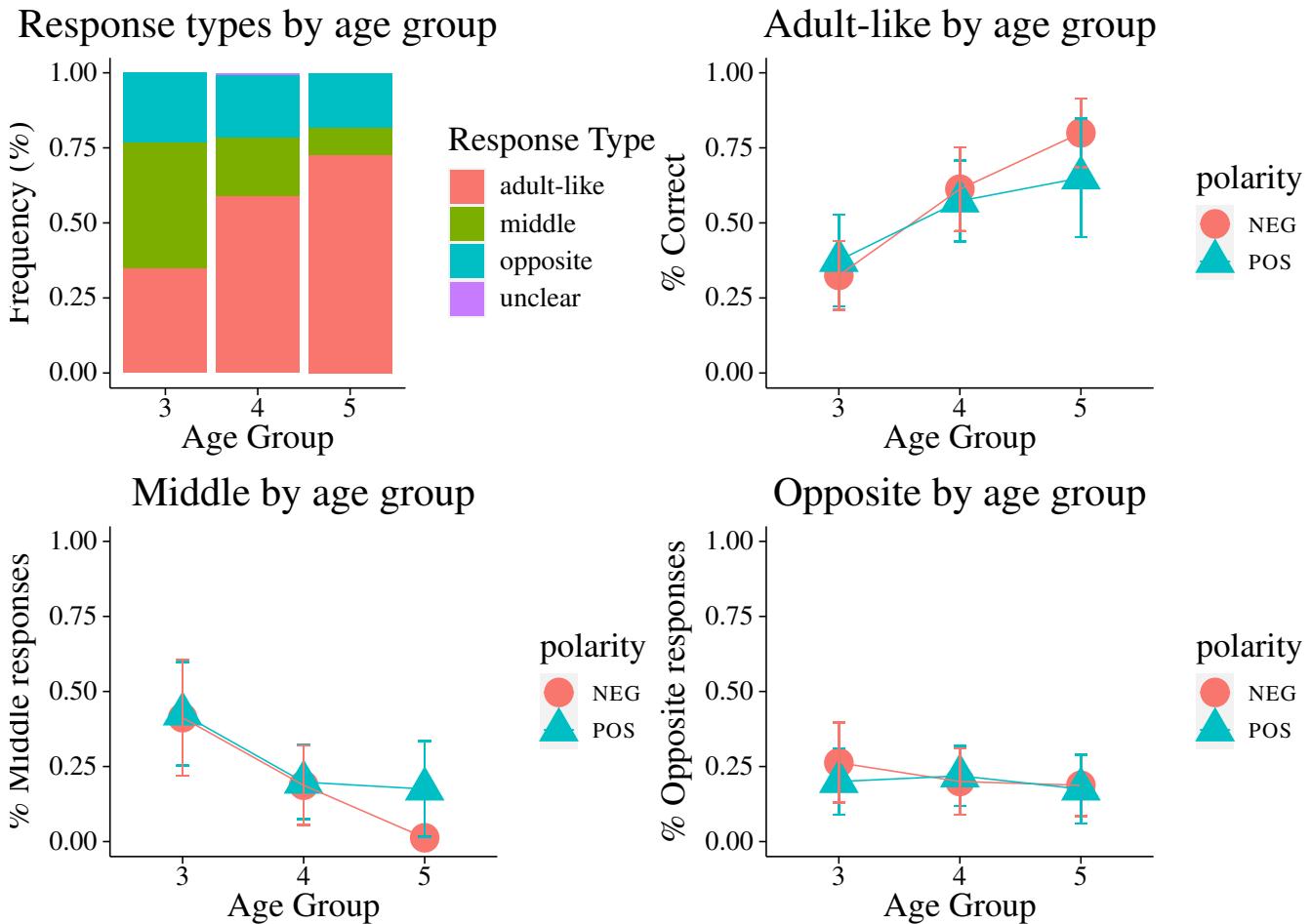


Figure 5. Results for Experiment 2. **(A)** Rate of response types by age group. **(B)** Rate of adult-like responses by age group and polarity. **(C)** Rate of middle responses by age group and polarity. **(D)** Rate of opposite responses by age group and polarity.

TABLES

		Likelihood inference	
		Least-likely	Most-likely
Sentence polarity	POS	✓	✓
	NEG	✓	✓

(A) Ambiguity theory of *even*

		Likelihood inference	
		Least-likely	Most-likely
Sentence polarity	POS	✓	
	NEG	✓	✓

(B) Scope theory of *even*

Table 1. Space of in principle available inferences associated with *even* on each theory.

		Selection pattern			
		Target characters for both sentence types	Opposite characters for both sentence types	Always rightmost or leftmost character	Any characters
Sentence type	Test sentences	33.3% (36/90)	38.9% (35/90)	27.8% (25/90) (22.2% for rightmost, 5.6% for leftmost)	
	Control sentences	20% (18/90)	40% (36/90)	36.7% (33/90) (26.7% for rightmost, 10% for leftmost)	3.3% (3/90)

Table 2. Rate of responses out of different types of pragmatics for test and control sentences in children's group. From [Kim \(2011\)](#).

Response Type	Justification Type	Age Group:	3	4	5	6	Total
adult-like	none		17	16	9	12	54
	random		0	4	7	6	17
	scalar		13	35	47	45	140
middle	none		15	9	5	1	30
	random		2	9	11	1	23
	scalar		0	3	4	0	7
opposite	none		10	10	3	1	24
	random		3	3	2	1	9
	scalar		3	9	11	15	38

Table 3. Experiment 1: Justifications for each response type by age.

Response Type	Polarity	% of responses	Std.dev
adult-like	POS	88%	32%
	NEG	82%	38%
middle	POS	0%	0%
	NEG	2%	14%
opposite	POS	12%	32%
	NEG	16%	37%

Table 4. Experiment 1, Adult control: Rates of response types by polarity.

Age Group vs. Age Group	4yo			5yo			6yo		
		Mean	CI		Mean	CI		Mean	CI
3yo	Adult-like	13	-11;39	Adult-like	33	-2;67	Adult-like	37	10;66
	Middle	-6	-34;16	Middle	-12	-47;21	Middle	-17	-43;0
	Opposite	-6	-26;13	Opposite	-21	-44;-1	Opposite	-20	-44;2
4yo				Adult-like	20	-10;48	Adult-like	24	2;47
				Middle	-5	-30;22	Middle	-11	-28;0
				Opposite	-15	-33;0	Opposite	-13	-33;4
5yo							Adult-like	4	-13;29
							Middle	-6	-28;3
							Opposite	1	-9;13

(A) All data.

Age Group vs. Age Group	4yo			5yo			6yo		
		Mean	CI		Mean	CI		Mean	CI
3yo	Adult-like	29	0;59	Adult-like	29	0;60	Adult-like	34	5;65
	Middle	-13	-48;11	Middle	-10	-47;16	Middle	-17	-50;1
	Opposite	-17	-44;7	Opposite	-19	-46;4	Opposite	-17	-45;8
4yo				Adult-like	0	-22;22	Adult-like	5	-15;27
				Middle	2	-14;22	Middle	-5	-19;3
				Opposite	-3	-20;13	Opposite	0	-18;18
5yo							Adult-like	5	-17;28
							Middle	-7	-25;2
							Opposite	3	-15;21

(B) Only data from negative stimuli.

Response Type vs. Response Type	Adult-like			Opposite		
		Mean	CI		Mean	CI
Middle	None	-57	-89;-19	None	-57	-90;-18
	Random	-8	-29;2	Random	-6	-28;6
	Scalar	63	26;94	Scalar	62	23;93
Opposite	None	1	-24;25			
	Random	-2	-8;2			
	Scalar	1	-24;27			

(C) Mean difference in rate of justification types by response type.

Table 5. Summary of statistical analysis for experiment 1: (A) Mean differences between age groups in rate of adult-like, middle, and opposite responses, for all data (B) Mean differences between age groups in rate of adult-like, middle, and opposite responses, for just negative stimuli. (C) Mean differences in rate of justification types by response type. Cells highlighted in green indicate a significant difference (i.e., not centered around 0).

Response Type	Justification Type	Age Group:	3	4	5	Total
adult-like	none		27	27	8	62
	random		6	13	4	23
	scalar		23	64	104	191
middle	none		47	15	1	63
	random		16	13	8	37
	scalar		4	6	6	16
opposite	none		18	10	6	34
	random		7	6	2	15
	scalar		12	21	21	54

Table 6. Experiment 2: Justifications offered for each response type by age.

Response Type	Polarity	% of responses	Std.dev
adult-like	POS	95%	21%
	NEG	91%	29%
middle	POS	3%	17%
	NEG	1%	11%
opposite	POS	2%	13%
	NEG	8%	27%

Table 7. Experiment 2, Adult control: Rates of response types by polarity.

Age Group vs. Age Group	4yo			5yo		
		Mean	CI		Mean	CI
3yo	Adult-like	39	6;71	Adult-like	56	27;86
	Middle	-37	-77;-1	Middle	-45	-80;10
	Opposite	-2	-23;20	Opposite	-12	-33;7
4yo				Adult-like	18	-2;41
				Middle	-7	-23;2
				Opposite	-10	-27;6

(A) Rates of different response types.

Response Type vs. Response Type	Adult-like			Opposite		
		Mean	CI		Mean	CI
Middle	None	-18	-53;61	None	-25	-63;1
	Random	-13	-36;0	Random	-12	-35;2
	Scalar	31	2;61	Scalar	37	5;68
Opposite	None	7	-10;33			
	Random	-1	-6;2			
	Scalar	-6	-31;12			

(B) Rates of different justification types by response type.

Table 8. Experiment 2: Mean differences in (A): the rate of adult-like, middle, and opposite responses, and (B): the rate of justification types (none, random, and scalar) by response type. Cells highlighted in green indicate a significant difference (i.e., not centered around 0).

Age group	Sentence Polarity	Inference:	Least-likely	Most-likely	Unclear	Total
3	POS		47	1	3	51
	NEG		1	17		18
4	POS		99	1	2	104
	NEG		1	161	5	167
	Unclear		1	1	1	3
5	POS		25	(1)	1	(29) 28
	NEG			42		42
	Unclear		1	1		2
6	POS		6		2	8
	NEG			23	1	24
Adults	POS		754	(57) 0	31	(842) 785
	NEG		16	1143	27	1186
	Unclear		1	2	7	10

Table 9. Child and child-directed (adult) production of *even* in negative and positive environments between 3-6 years old. Figures in gray show numbers before post-hoc reanalysis (cf. footnote 27).

Age	NEG:POS <i>even</i> ratio	NEG count	POS count
3 yrs	0.35	18	51
4 yrs	1.6	167	104
5 yrs	1.5	42	28
6 yrs	3	24	8
Adult	(1.5) 1.6	(1186) 1243	785

Table 10. The ratio of negative *even* sentences by age. In adults, we include downward-entailing sentences which lack sentential negation in the NEG figures, with ratios/counts excluding these sentences in gray.

Supplementary Material

0.1 Figures

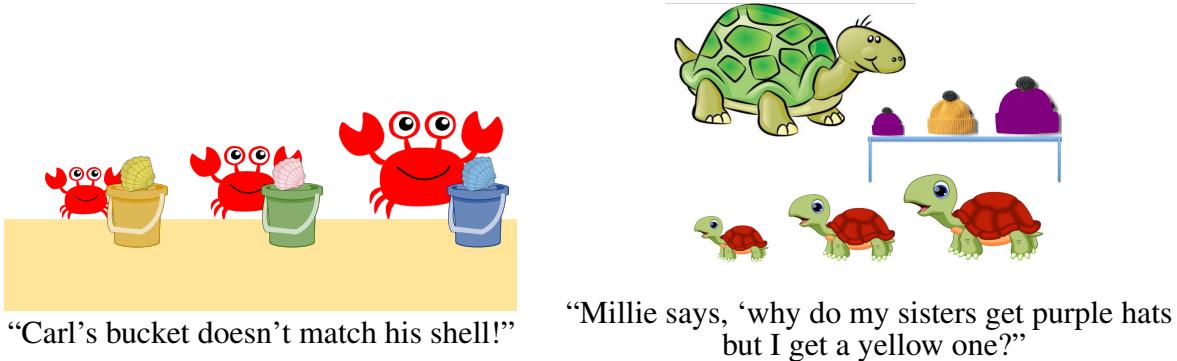


Figure S1. Two different filler items: matching by color/size.

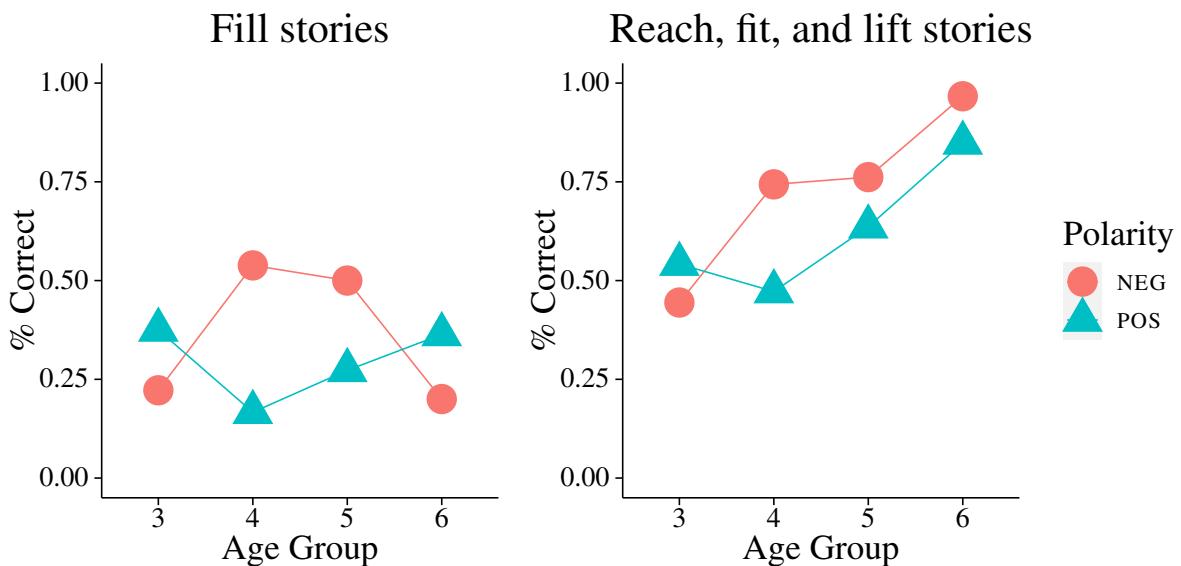


Figure S2. Adult-like behavior in Experiment 1, by story type.

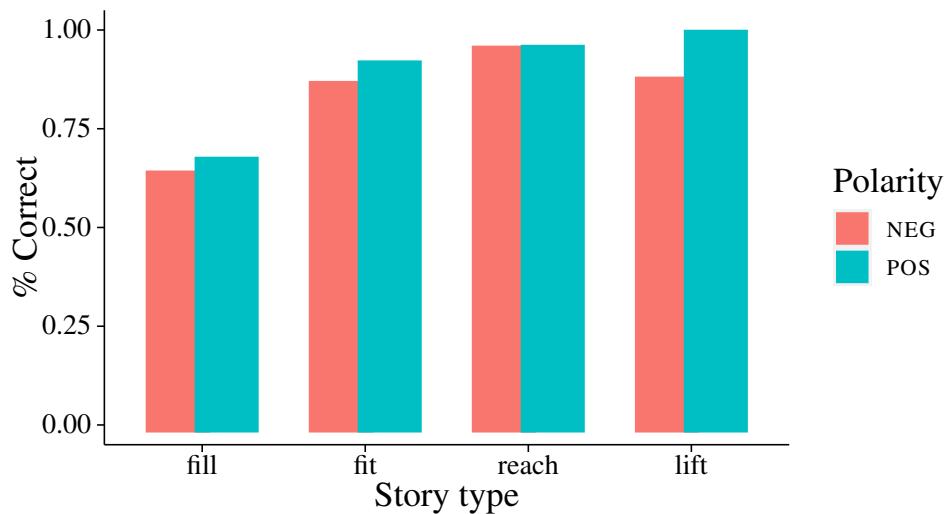


Figure S3. Adult study 1: adult performance by story type.

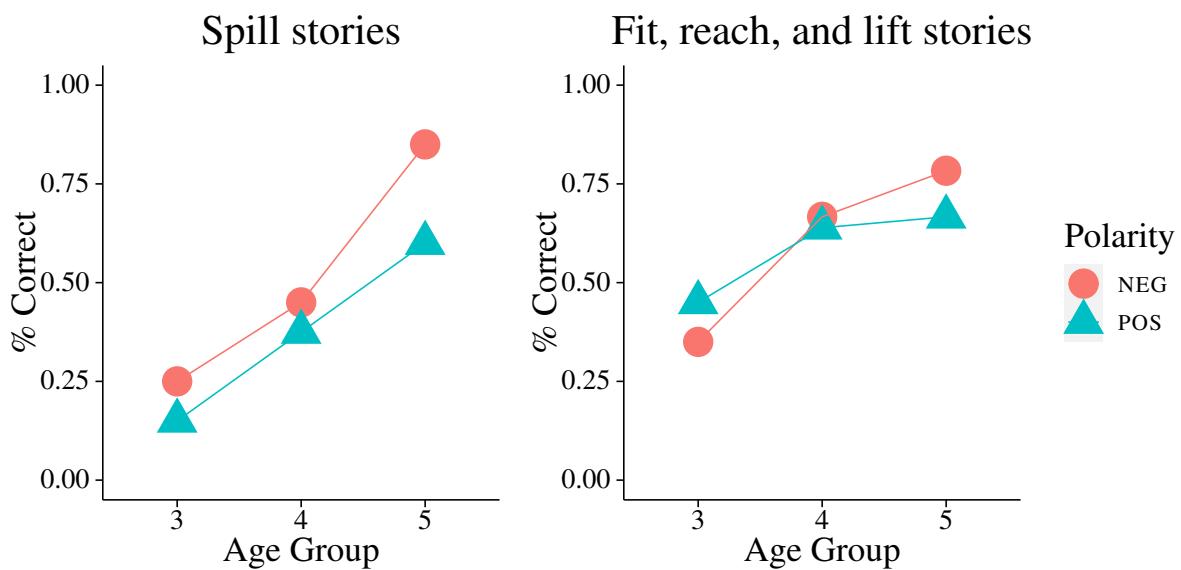


Figure S4. Adult-like behavior in Experiment 2, by story type.

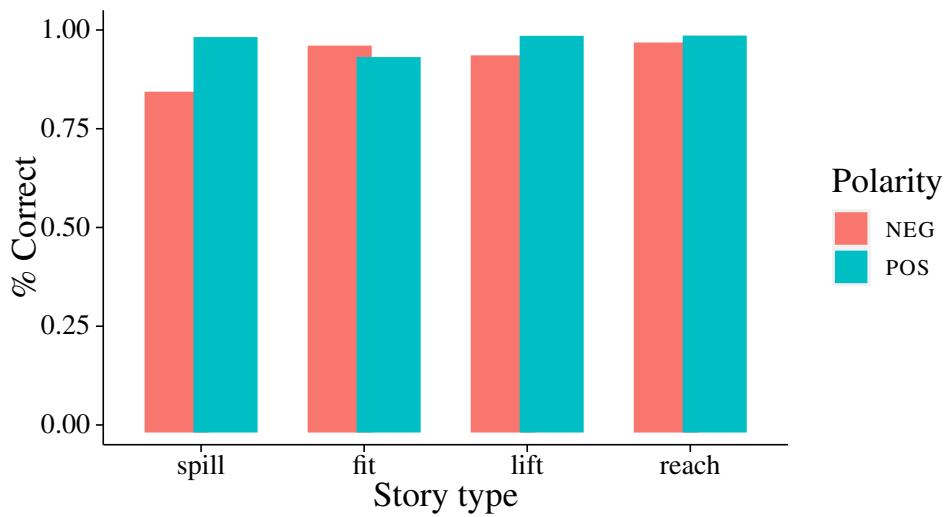


Figure S5. Adult study 2: adult performance by story type.

CHILDES Coding

Enterer
Elise

Window Size:
1 3 5 7 9 11 13 15

Token ID: 4840302 (108 / 1034)

Irrelevant/Repetition

Embedded Even

Embedded Associate

Comparative

Negation Present?
Yes

Neg-Even

Sentential Negation?
Yes

Syntactic Position
Post-Subject

Focus Associate
Verb

Likelihood
Most-likely

Comments

< Previous Next >

Transcript Responses All Tokens SQL Code
See full CHILDES transcript (<https://childestalkbank.org/browser/index.php?url=Eng-NA/Kuczaj/040301.cha#>)

Abe (Target_Child, CHI): Dad xxx through this xxx

Abe (Target_Child, CHI): that's better

NA (Father, FAT): did they all make it Abe

Abe (Target_Child, CHI): not yet

Abe (Target_Child, CHI): they all are going to make it for sure

Abe (Target_Child, CHI): this one has_to climb up

Abe (Target_Child, CHI): that one just stays there

Abe (Target_Child, CHI): you know why

NA (Father, FAT): why

Abe (Target_Child, CHI): nobody could even move it not even the water not even the wind

NA (Father, FAT): that boat's so strong that it can just stay where it wants to

Abe (Target_Child, CHI): yep and it's even stronger than a tidal wave

NA (Father, FAT): wow

Abe (Target_Child, CHI): whoops

Abe (Target_Child, CHI): that's not even stronger than a tidal wave

Abe (Target_Child, CHI): except for sure if you want to get back on you'd have_to jump really high

Abe (Target_Child, CHI): that wasn't even a xxx

NA (Father, FAT): what happened Abe

Abe (Target_Child, CHI): a tidal wave hitted this one

Figure S6. An example instance of *even* from the CHILDES corpus, along with its coding.

	<i>Even's</i> position	Sentence Polarity	Focus associate:	Subject	Other	Total
Children	Pre-subject	POS		41	27	68
		NEG		15	3	18
	Post-subject	POS		13	207	220
		NEG		10	371	381
		Unclear		5	5	5
Adults	Pre-subject	POS		89	89	178
		NEG		24	23	47
		Unclear		2	2	2
	Post-subject	POS		11	637	648
		NEG		8	1123	1131
		Unclear		5	5	5

Table S1. Both children and adults most often associate pre-subject *even* with subject focus and post-subject *even* with VP-internal focus. Very few instances of pre-subject *even*. Highlighted cells are those which instantiate the form of our comprehension study stimuli.

<i>Even-NEG</i> Order	Sentential Negation	Child Total	Adult Total
NEG- <i>Even</i>	Yes	457	1156
	No	6	6
	Unclear	5	4
<i>Even-NEG</i>	Yes	16	23
	No	26	58
	Unclear	4	1

Table S2. In negative polarity sentences (Sentential Negation=Yes), both children and adults have a strong tendency to place *even* after the negative element. Highlighted cells are those which instantiate the form of our comprehension stimuli.