1 Using contrastive inferences to learn about new words and categories

2 Claire Bergey[1] & Dan Yurovsky[2]

3 [1] The University of Chicago

4 [2] Carnegie Mellon University

5 Author Note

6 All data and code for these analyses are available at

7 https://osf.io/3f8hy/?view_only=9a196db0444c4867bc899cc70a7a1e9c.

8 Correspondence concerning this article should be addressed to Claire Bergey, 5848 S.

9 University Avenue, Chicago, IL 60637. E-mail: cbergey@uchicago.edu

10                                             Abstract

11      In the face of unfamiliar language or objects, description is one cue people can use to learn

12      about both. Beyond narrowing potential referents to those that match a descriptor (e.g.,

13      "tall"), listeners could infer that a described object is one that contrasts with other relevant

14      objects of the same type (e.g., "the tall cup" contrasts with another, shorter cup). This

15      contrast may be in relation to other present objects in the environment (this cup is tall

16      among present cups) or to the referent's category (this cup is tall for a cup in general). In

17      three experiments, we investigate whether listeners use descriptive contrast to learn new

18      word-referent mappings and learn about novel categories' feature distributions. People use

19      contrastive inferences to guide their referent choice, though size–and not color–adjectives

20      prompt them to consistently choose the contrastive target over alternatives (Experiment 1).

21      People also use color and size description to infer that a novel object is atypical of its

22      category (Experiments 2 and 3). However, these two inferences do not trade off substantially:

23      people infer a described referent is atypical even when the descriptor was necessary to

24      establish reference. We model these experiments in the Rational Speech Act (RSA)

25      framework and find it predicts both of these inferences, and a very small trade-off between

26      them consistent with the lack of trade-off we observe in people's inferences. Overall, people

27      are able to use descriptive contrast to resolve reference and make inferences about a novel

28      object's category, allowing them to learn more about new things than literal meaning alone

29      allows.

30      *Keywords:* concept learning; pragmatics; communication

31      Word count: 1385

32          Using contrastive inferences to learn about new words and categories

33          An utterance can say much more about the world than its literal interpretation might

34     suggest. For instance, if you hear a colleague say "We should hire a female professor," you

35     might infer something about the speaker's goals, the makeup of a department, or even the

36     biases of a field—none of which is literally stated. Pragmatic inferences like these are

37     pervasive in everyday conversation: by reasoning about what someone says in relation to the

38     context and what they might have said otherwise, we can glean more of their intended

39     meaning.

40          But what if you didn't know the meaning of the key words in someone's

41     utterance–could you use the same kind of pragmatic inferences to learn about new words and

42     categories? Suppose a friend asked you to "Pass the tall dax." You might look around the

43     room for two similar things that vary in height, and hand the taller one to them. Intuitively,

44     your friend must have said the word "tall" for a reason (Grice, 1975). One possibility is that

45     your friend wanted to distinguish the dax they wanted from the dax they did not. People

46     appear to make these kinds of inferences quite rapidly for objects they know; for instance, as

47     soon as they hear the word "tall," they already begin looking to a tall familiar object with a

48     short competitor nearby–even if there are other tall objects around (Sedivy, Tanenhaus,

49     Chambers, & Carlson, 1999).

50          If you only saw one object around whose name you didn't know, you might draw a

51     different inference: this dax might be a particularly tall dax. In this case, you might think

52     your friend used the word "tall" for a different reason–not to distinguish the dax they wanted

53     and other daxes around you, but to distinguish the dax they want from other daxes in the

54     world. This would be consistent with data from production studies, in which people tend to

55     describe atypical features more than they describe typical ones (Mitchell, Reiter, & Deemter,

56     2013; Rubio-Fernández, 2016; Westerbeek, Koolen, & Maes, 2015). For instance, people

57     almost always say "blue banana" to refer to a blue banana, but almost never say "yellow

banana" to refer to a yellow one.

In each of these cases, you would have used a pragmatic inference to learn something new. In the second case, you would have learned the name for a novel category "dax," and also something about the typical of size of daxes: most of them are shorter than the one you saw. In the first case, you would have also learned a new word, but would have you learned something about the typical size of daxes as well, beyond the two daxes you observed? One possibility is that you would not: You can explain your friend's use of "tall" as being motivated by the need to distinguish between the two daxes in the room, and thus you should infer nothing about the other daxes in the world. If reference is the primary motivator of speakers' word choice, as implicitly assumed in much research (e.g., Pechmann, 1989; Arts, Maes, Noordman, & Jansen, 2011; Engelhardt, Barış Demiral, & Ferreira, 2011), then people should draw no further inferences once the need for referential disambiguation can explain away a descriptor like "tall." If, on the other hand, pragmatic reasoning weighs multiple goals simultaneously–here, reference and conveying typicality–people may integrate typicality as just one factor the speaker weighs in using description, leading to graded inferences about the referent's identity and about its category's features.

In this paper, we present a series of experiments that test two ways in which people could use pragmatic inference to learn about novel categories. First, we examine whether listeners use descriptive contrast to resolve referential ambiguity. In a reference game, participants saw groups of novel objects and were asked to pick one with a referring expression, e.g., "Find the small toma." If people interpret description contrastively, they should infer that the description was necessary to identify the referent–that the small toma contrasts with some different-sized toma on the screen. We show that people can use contrastive inference–even with unfamiliar objects–to resolve reference and thus to learn the meaning of the new word "toma."

Second, we test whether people use descriptive contrast to make inferences about a

novel object's category. Participants were presented with two interlocutors who exchange objects using referring expressions, such as "Pass me the blue toma." If people interpret description as contrasting with an object's category, they should infer that in general, few tomas are blue. Crucially, we vary the object contexts such that in some contexts, the adjective is necessary to establish reference, and in others, it is superfluous. Overall, we show that people can use contrastive inferences both to establish reference and to make inferences about novel categories' feature distributions, and that they do not trade off strongly between these two inferences. We extend a version of the Rational Speech Act model to show that listeners' reasoning about speakers reflects a graded integration of informativity with respect to both reference and typicality.

In order to determine whether people can use prenominal adjective contrast to disambiguate referents, and how those inferences are affected by adjective type, we use reference games with novel objects. Novel objects provide both a useful experimental tool and an especially interesting testing ground for contrastive inferences. These objects have unknown names and feature distributions, creating the ambiguity that is necessary to test referential disambiguation and category learning. Here, we ask: can people use pragmatic inferences from description to learn about unfamiliar things in the world?

## Experiment 1

In Experiment 1, we ask whether people use descriptive contrast to identify the target of an ambiguous referring expression. Our experiment was inspired by work from Sedivy et al. (1999) showing that people interpret at least some prenominal adjective use as contrastive when the target referents are familiar objects. In their task, four objects appeared on a screen: a target (e.g., a tall cup), a contrastive pair (e.g., a short cup), a competitor that shares the target's feature but not category (e.g., a tall pitcher), and an irrelevant distractor (e.g., a key). Participants then heard a referring expression: "Pick up the tall cup." Participants looked more quickly to the correct object when the utterance referred to an

110 object with a same-category contrastive pair (tall cup vs. short cup) than when it referred to

111 an object without a contrastive pair (e.g., when there was no short cup in the display).

112     Their results suggest that listeners expect speakers to use prenominal description when

113 they are distinguishing between potential referents of the same type, and listeners use this

114 inference to rapidly allocate their attention to the target as an utterance progresses. This

115 principle does not apply equally across adjective types, however: color adjectives seem to

116 hold less contrastive weight (Sedivy, 2003), perhaps because color adjectives are often used

117 redundantly in English–that is, people describe objects' colors even when this description is

118 not necessary to establish reference (Pechmann, 1989). These experiments demonstrate that

119 listeners use contrast among familiar referents to guide their attention allocation, though not

120 their explicit referent choice, which occurs after the noun disambiguates the object.

121     In a pre-registered referential disambiguation task, we presented participants with

122 arrays of novel fruit objects. On critical trials, participants saw a target object, a lure object

123 that shared the target's critical feature but not its shape, and a contrastive pair that shared

124 the target's shape but not its critical feature (Fig. 1). Participants heard an utterance,

125 sometimes denoting the critical feature: "Find the [blue/big] toma." For the target object,

126 which had a same-shaped counterpart, use of the adjective was necessary to establish

127 reference. For the lure, which was unique in shape, the adjective was relatively superfluous

128 description. If participants use contrastive inference to choose novel referents, they should

129 choose the target object more often than the lure. To examine whether contrast occurs

130 across adjective types, we test participants in two conditions: color contrast and size

131 contrast. Though we expect participants to shift toward choosing the item with a contrastive

132 pair in both conditions, we do not expect them to treat color and size equally. Because color

133 is often used redundantly in English while size is not, we expect size to hold more contrastive

134 weight, encouraging a more consistent contrastive inference (Pechmann, 1989). The

135 pre-registration of our method, recruitment plan, exclusion criteria, and analyses can be

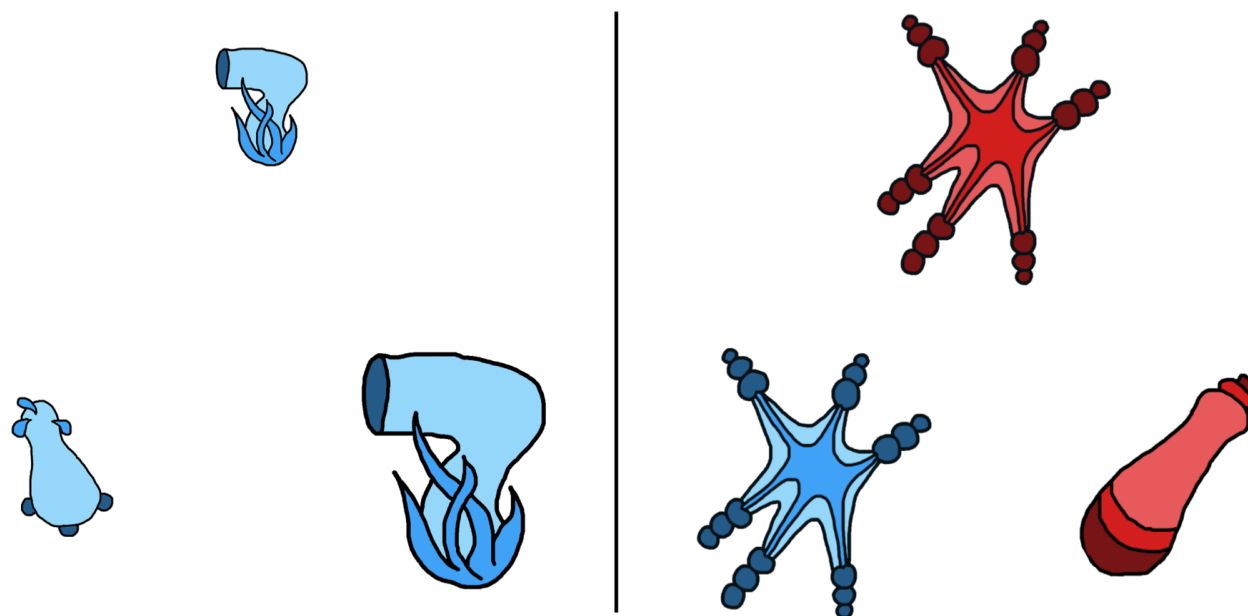found on the Open Science Framework here: https://osf.io/pqkfy .

*Figure 1*. On the left: an example of a contrastive trial in which the critical feature is size. Here, the participant would hear the instruction "Find the small toma." On the right: an example of a contrastive trial in which the critical feature is color. Here, the participant would hear the instruction "Find the red toma." In both cases, the target is the top object.

## Method

**Participants.** We recruited a pre-registered sample of 300 participants through Amazon Mechanical Turk. Half of the participants were assigned to a condition in which the critical feature was color (stimuli contrasted on color), and the other half were assigned to a condition in which the critical feature was size. Each participant gave informed consent and was paid $0.30 in exchange for their participation.

**Stimuli.** Stimulus displays were arrays of three novel fruit objects. Fruits were chosen randomly at each trial from 25 fruit kinds. Ten of the 25 fruit drawings were adapted and redrawn from Kanwisher, Woods, Iacoboni, and Mazziotta (1997); we designed the remaining 15 fruit kinds. Each fruit kind had an instance in each of four colors (red, blue,
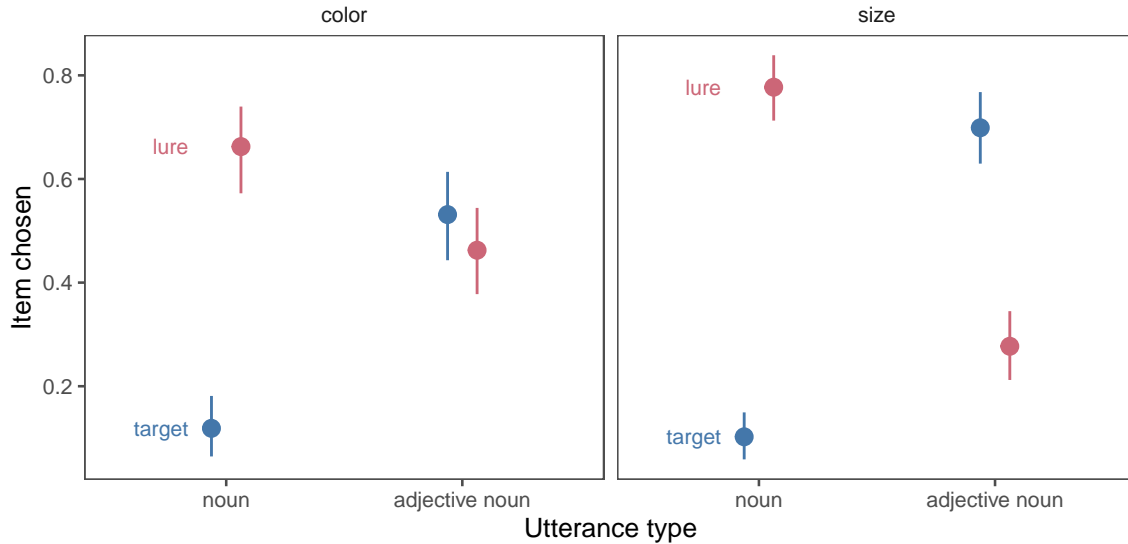
*Figure 2*. Proportion of times that participants chose the target and lure items as a function of condition and whether an adjective was provided. Points indicate group means; error bars indicate 95% confidence intervals computed by non-parametric bootstrapping.

147 green, or purple) and two sizes (big or small). Particular target colors were assigned

148 randomly at each trial and particular target sizes were counterbalanced across display types.

149 There were two display types: unique target displays and contrastive displays. Unique target

150 displays contained a target object that had a unique shape and was unique on the trial's

151 critical feature (color or size), and two distractor objects that matched each other's (but not

152 the target's) shape and critical feature. These unique target displays were included as a

153 check that participants were making reasonable referent choices and to space out contrastive

154 displays to prevent participants from dialing in on the contrastive object setup during the

155 experiment. Contrastive displays contained a target, its contrastive pair (matched the

156 target's shape but not its critical feature), and a lure (matched the target's critical feature

157 but not its shape) (Fig. 1). The positions of the target and distractor items were

158 randomized within a triad configuration.

159 **Design and Procedure.** Participants were told they would play a game in which

160 they would search for strange alien fruits. Each participant saw eight trials. Half of the trials

161  were unique target displays and half were contrastive displays. Crossed with display type,

162  half of trials had audio instructions that described the critical feature of the target (e.g.,

163  "Find the [blue/big] toma"), and half of trials had audio instructions with no adjective

164  description (e.g., "Find the toma"). A name was randomly chosen at each trial from a list of

165  eight nonce names: blicket, wug, toma, gade, sprock, koba, zorp, and lomet.

166      After completing the study, participants were asked to select which of a set of alien

167  words they had heard previously during the study. Four were words they had heard, and

168  four were novel lure words. Participants were dropped from further analysis if they did not

169  meet our pre-registered exclusion criteria of responding to at least 6 of these 8 memory check

170  questions correctly (above chance performance as indicated by a one-tailed binomial test at

171  the $p = .05$ level) and answering all four color perception check trials correctly (resulting $n =$

172  163).

**Results and Discussion**

174      We first confirmed that participants understood the task by analyzing performance on

175  unique target trials, the filler trials in which there was a target unique on both shape and

176  the relevant adjective. We asked whether participants chose the target more often than

177  expected by chance (33%) by fitting a mixed effects logistic regression with an intercept

178  term, a random effect of subject, and an offset of $logit(1/3)$ to set chance probability to the

179  correct level. The intercept term was reliably different from zero for both color ($\beta = 6.64$,

180  $t = 4.10$, $p < .001$) and size ($\beta = 2.25$ , $t = 6.91$, $p < .001$), indicating that participants

181  consistently chose the unique object on the screen when given an instruction like "Find the

182  (blue) toma." In addition, participants were more likely to select the target when an

183  adjective was provided in the audio instruction in both conditions. We confirmed this effect

184  statistically by fitting a mixed effects logistic regression predicting target selection from

185  condition, adjective use, and their interaction with random effects of participants. Use of

186  description in the audio increased target choice ($\beta = 3.85$, $t = 3.52$, $p < .001$), and adjective

187 type (color vs. size) was not statistically related to target choice ($\beta = $ -0.48, $t = $ -1.10, $p = $

188 .269). The two effects did not significantly interact ($\beta = $ -2.24, $t = $ -1.95, $p = $ .051).

189 Participants had a general tendency to choose the target in unique target trials, which was

190 strengthened if the audio instruction contained the relevant adjective. These effects did not

191 significantly differ between color and size adjectives, which suggests that participants did not

192 treat color and size differently in these baseline trials.

193       Our key pre-registered analysis was whether participants would choose the target

194 object on contrastive trials–when they heard an adjective in the referential expression. To

195 perform this test, we compared participants' rate of choosing the target to their rate of

196 choosing the lure, which shares the relevant critical feature with the target, when they heard

197 the adjective. Overall, participants chose the target with a contrasting pair more often than

198 the unique lure, indicating that they used contrastive inferences to resolve reference ($\beta = $

199 0.53, $t = $ 3.83, $p = \; < .001$). To test whether the strength of the contrastive inference differed

200 between color and size conditions, we pre-registered a version of this regression with a term

201 for adjective type, and found that people were more likely to choose the target over the lure

202 in the size condition than the color condition ($\beta = $ 0.87, $t = $ 3.12, $p = $ .002). Given this

203 result, we tested whether people consistently chose the target over the lure on the color and

204 size data separately, as a stricter check of whether the effect was present in both conditions.

205 Considering color and size separately, participants chose the target significantly more often

206 than the lure in the size condition ($\beta = $ 0.86, $t = $ 4.41, $p = \; < .001$), but not in the color

207 condition ($\beta = $ 0.15, $t = $ 0.75, $p = $ .455). On contrastive trials in which a descriptor was not

208 given, participants dispreferred the target, instead choosing the lure object, which matched

209 the target on the descriptor but had a unique shape ($\beta = $ -2.65, $t = $ -5.44, $p = \; < .001$).

210 Participants' choice of the target in the size condition was therefore not due to a prior

211 preference for the target in contrastive displays, but relied on contrastive interpretation of

212 the adjective. In the supplemental materials, we report an additional pre-registered analysis

213 of all Experiment 1 data with maximal terms and random effects; those results are consistent

with the more focused tests reported here.

When faced with unfamiliar objects referred to by unfamiliar words, people can use pragmatic inference to resolve referential ambiguity and learn the meanings of these new words. In Experiment 1, we found that participants have a general tendency to choose objects that are unique in shape when reference is ambiguous. However, when they hear an utterance with description (e.g., "blue toma", "small toma"), they shift away from choosing unique objects and toward choosing objects that have a similar contrasting counterpart. Furthermore, use of size adjectives–but not color adjectives–prompts people to choose the target object with a contrasting counterpart more often than the unique lure object. We find that people are able to use contrastive inferences about size to successfully resolve which unfamiliar object an unfamiliar word refers to.

## Model

To formalize the inference that participants were asked to make, we developed a model in the Rational Speech Act Framework (RSA, Frank & Goodman, 2012). In this framework, pragmatic listeners ($L$) are modeled as drawing inferences about speakers' ($S$) communicative intentions in talking to a hypothetical literal listener ($L_0$). This literal listener makes no pragmatic inferences at all, evaluating the literal truth of statements (e.g., it is true that a red toma can be called "toma" and "red toma" but not "blue toma"), and chooses randomly among all referents consistent with a statement. In planning their referring expressions, speakers choose utterances that are successful at accomplishing two goals: (1) making the listener as likely as possible to select the correct object, and (2) minimizing their communicative cost (i.e., producing as few words as possible). Pragmatic listeners use Bayes' rule to invert the speaker's utility function, essentially inferring what the speaker's intention was likely to be given the utterance they produced.

$$Literal : P_{Lit} = \delta\left(u, r\right) P\left(r\right)$$

$$Speaker : P_S\left(u|r\right) \propto \alpha\left(P_{Lit}\left(r|u\right) - C\right)$$

$$Listener : P_{Learn}\left(r|u\right) \propto P_s\left(u|r\right) P\left(r\right)$$

For this experiment, we build on a Rational Speech Act model developed by Frank and Goodman (2014) to jointly resolve reference and learn new words. The primary extension of RSA is that the pragmatic learner is a pragmatic listener who has uncertainty about the meanings of words in their language, and thus cannot directly compute the speaker's utility as written. Instead, the speaker's utility is conditioned on the set of mappings, and the learners must also infer which set of mappings is correct:

$$Learner : P_L\left(r|u\right) \propto P_s\left(u|r; m\right) P\left(r\right) P\left(m\right)$$

In these experiments, we assume that the prior probability to refer to each object $\left(P\left(r\right)\right)$ is equal, and similarly that all mappings $\left(P\left(m\right)\right)$ are equally likely, so they cancel out in computations. We further assume that the cost of producing any word is identical, and so the cost of an utterance is equal to its length. All that remains is to specify the possible mappings, and literal meanings, and alternative utterances possible on each trial of the experiment. We describe the size condition here, but the computation for the color condition is analogous.

On the trial shown in the left panel of Figure 1 people see two objects that look something like a hair dryer and one that looks like a pear and they are asked to "Find the toma." Here, in the experiment design and the model, we take advantage of the fact that English speakers tend to assume that nouns generally correspond to differences in shape rather than other features (Landau, Smith, & Jones, 1992). Given this, the two possible mappings are $\{m_1 : hairdryer - \text{``}toma\text{''}, pear - \text{``}?\text{''}\}$, and

259 $\{m_2 : hairdryer - "?", pear - "toma"\}$. The literal semantics of each object allow them to

260 be referred to by their shape label (e.g. "toma"), or by a descriptor that is true of them

261 (e.g. "small"), but not names for other shapes or untrue descriptors.

262      Having heard "Find the toma," the model must now choose a referent. If the true

263 mapping for "toma" is the hair dryer ($m_1$), this utterance is ambiguous to the literal listener,

264 as there are two referents consistent with the literal meaning toma. Consequently, whichever

265 of the two referents the speaker intends to point out to the learner, the speaker's utility will

266 be relatively low. Alternatively, if the true mapping for "toma" is the pear ($m_1$), then the

267 utterance will be unambiguous to the literal listener, and thus the speaker's utterance will

268 have higher utility. As a result, the model can infer that the more likely mapping is $m_2$ and

269 choose the pear, simultaneously resolving reference and learning the meaning of "toma."

270      If instead the speaker produced "Find the small toma," the model will make a different

271 inference. If the true mapping for "toma" is hair dryer ($m_2$), this utterance now uniquely

272 identifies one referent for the literal listener and thus has high utility. It also uniquely

273 identifies the target if "toma" means pear ($m_1$). However, if "toma" means pear, the

274 speaker's utterance was inefficient because the single word utterance "toma" would have

275 identified the target to the literal listener and incurred less cost. Thus, the model can infer

276 that "toma" is more likely to mean hair dryer and choose the small hair dryer appropriately.

277      While these descriptions use deterministic language for clarity, the model's

278 computation is probabilistic and thus reflects tendencies to choose those objects rather than

279 fixed rules. Figure 3 shows model predictions alongside people's behavior for the size and

280 color contrast conditions in Experiment 1. In line with the intuition above, the model

281 predicts that hearing a bare noun (e.g. "toma") should lead people to infer that the intended

282 referent is the unique object (lure), whereas hearing a modified noun (e.g. "small toma")

283 should lead people to infer that the speaker's intended referent has a same-shaped

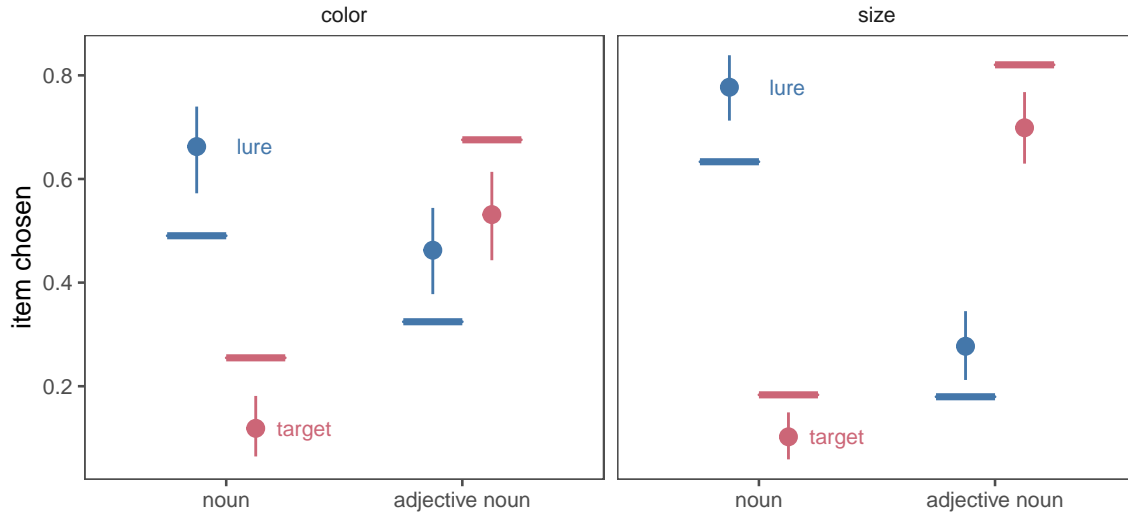284 counterpart without the described feature (i.e., is the target object).

*Figure 3*. Proportion of times that people (and our model) chose the target and lure items as a function of adjective type and whether an adjective was provided. Points indicate empirical means; error bars indicate 95% confidence intervals computed by non-parametric bootstrapping. Solid horizontal lines show model predictions.

²⁸⁵     Our empirical data suggest that people treat color and size adjectives differently,

²⁸⁶ making a stronger contrastive inference with size than with color. One potential explanation

²⁸⁷ for this difference is that people are aware of production asymmetries between color and size.

²⁸⁸ As mentioned, speakers tend to over-describe color, providing more color adjectives than

²⁸⁹ necessary to establish reference, while describing size more minimally (Nadig & Sedivy, 2002;

²⁹⁰ Pechmann, 1989). Listeners may be aware of this production asymmetry and discount the

²⁹¹ contrastive weight of color adjectives with respect to reference.

²⁹²     In the Rational Speech Act model, this kind of difference is captured neatly by a

²⁹³ difference in the listener's beliefs about the speaker's rationality (i.e. how sensitive the

²⁹⁴ speaker is to differences in utility of different utterances). To determine the value of the

²⁹⁵ rationality parameter in each condition, we used Empirical Bayesian inference to estimate

²⁹⁶ the likely range of parameter values. These estimates varied substantially across conditions,

²⁹⁷ with the rationality parameter in the color condition estimated to be 2.00 with a 95%

298 credible interval of [1.37, 2.63], and the rationality parameter in the size condition estimated
299 to be 3.98 [3.22, 4.74].

300     Figure 3 shows the model predictions along with the empirical data from Experiment 1.
301 The model broadly captures the contrastive inference–when speakers produce an adjective
302 noun combination like "red toma," the model selects the target object more often than the
303 lure object. The extent to which the model makes this inference varies as predicted between
304 the color and size adjective conditions in line with the different estimated rationality values.
305 In both conditions, despite estimating the value of rationality that makes the observed data
306 more likely, the model overpredicts the extent of the contrastive inference that people make.
307 Intuitively, it appears that in over the strength of their contrastive inferences, people have an
308 especially strong tendency to choose a unique object when they hear an unmodified noun
309 (e.g. "toma"). In an attempt to capture this uniqueness tendency, the model overpredicts the
310 extent of the contrastive inference.

311     The model captures the difference between color and size in a difference in the
312 rationality parameter, but leaves open the ultimate source of this difference in rationality.
313 Why there is a production asymmetry in the first place? For now, we bracket this question
314 and note that listeners in our task appropriately discount color's contrastive weight given
315 production norms.

316     An alternative way to capture this preference would be to locate it in a different part
317 of the model. One possibility is that the semantics of color and size work differently. A
318 recent model from Degen, Hawkins, Graf, Kreiss, and Goodman (2020) does predict a
319 color–size asymmetry based on different semantic exactness. In this model, literal semantics
320 are treated as continuous rather than discrete, so "blue" is neither 100% true nor 100% false
321 of a particular object, but can instead be 90% true. They successfully model a number of
322 color–size asymmetries by treating color as having stronger literal semantics (i.e. "blue toma"
323 is a better description of a small blue toma than "small toma" is). However, this model

324 predicts the opposite asymmetry of what we found. Because color has stronger semantics

325 than size, the listener in this model shows a stronger contrast effect for color than size. We

326 show this effect in appendix A. Thus, though a continuous semantics can explain our

327 asymmetry, this explanation is unlikely given the continuous semantics that predicts other

328 empirical color–size asymmetries does not predict our findings.

329      Overall, we found that people can use contrastive inferences from description to map

330 an unknown word to an unknown object. This inference is captured by an extension of the

331 Rational Speech Act model using a pragmatic learner, who is simultaneously making

332 inferences over possible referents and possible lexicons. This model can also capture people's

333 tendency to make stronger contrastive inferences from color description than size description

334 through differences in the rationality parameter, though the origin of these differences cannot

335 be pinned down with this experiment alone. Our experiment and model results suggest that

336 people can resolve a request like "Give me the small dax" by reasoning that the speaker must

337 have been making a useful distinction by mentioning size, and therefore looking for multiple

338 similar objects that differ in size and choosing the smaller one. Immediately available objects

339 are not the only ones worth making a distinction from, though. Next, we turn to another

340 salient set of objects a speaker might want to set a referent apart from: the referent's

341 category.

## **Experiment 2**

343      When referring to a *big red dog* or a *hot-air balloon*, we often take care to describe

344 them–even when there are no other dogs or balloons around. Speakers use more description

345 when referring to objects with atypical features (e.g., a yellow tomato) than typical ones

346 (e.g., a red tomato; Mitchell et al., 2013; Bergey, Morris, & Yurovsky, 2020; Rubio-Fernández,

347 2016; Westerbeek et al., 2015). This selective marking of atypical objects potentially supplies

348 useful information to listeners: they have the opportunity to not only learn about the object

349 at hand, but also about its broader category. Further, this kind of contrast may help make

350 sense of the asymmetry between color and size adjectives we found in Experiment 1. Color

351 adjectives that are redundant with respect to reference are not necessarily redundant in

352 general. Rubio-Fernández (2016) demonstrates that speakers often use 'redundant' color

353 adjectives to describe colors when they are central to the category's meaning (e.g., colorful

354 t-shirts) or when they are atypical (e.g., a purple banana). Therefore, color and size may

355 hold similar contrastive weight with respect to the category's feature distribution. In

356 Experiment 2, we test whether listeners use descriptive contrast with a novel object's

357 category to learn about the category's feature distribution.

358      If listeners do make contrastive inferences about typicality, it may not be as simple as

359 judging that an over-described referent is atypical. Description can serve many purposes. In

360 the prior experiment, we investigated its use in contrasting between present objects. If a

361 descriptor was needed to distinguish between two present objects, it may not have been used

362 to mark atypicality. For instance, in the context of a bin of heirloom tomatoes, a speaker

363 who wanted a red one in particular might specify that they want a "red tomato" rather than

364 just asking for a "tomato." In this case, the adjective "red" is being used contrastively with

365 respect to reference (as in Experiment 1), and not to mark atypicality. Thus, a listener who

366 does not know much about tomatoes may attribute the use of "red" to referential

367 disambiguation given the context and not infer that red is an unusual color for tomatoes.

368      In Experiment 2, we used an artificial language task to set up just this kind of learning

369 situation. We manipulated the contexts in which listeners hear adjectives modifying novel

370 names of novel referents. We asked whether listeners infer that these adjectives identify

371 atypical features of the named objects, and whether the strength of this inference depends on

372 the referential ambiguity of the context in which adjectives are used.
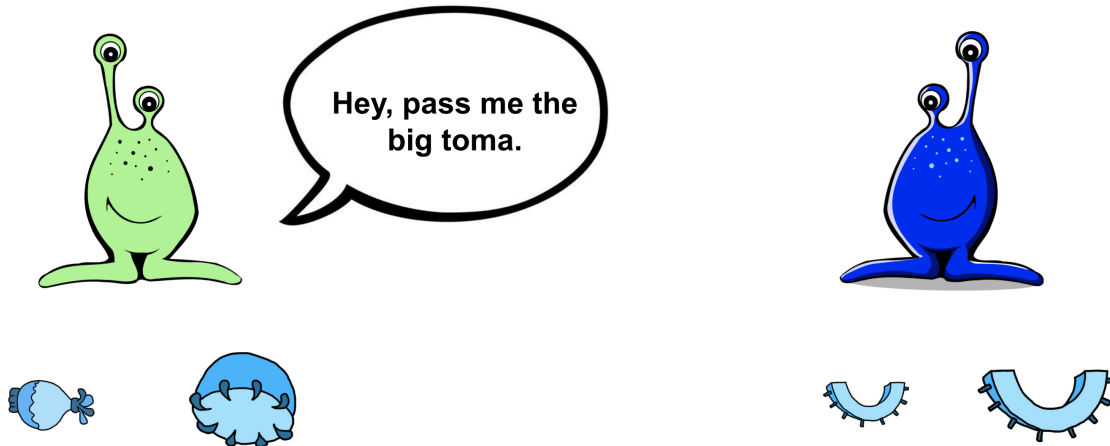
*Figure 4*. Experiment 2 stimuli. In the above example, the critical feature is size and the object context is a within-category contrast: the alien on the right has two same-shaped objects that differ in size.

## Method

**Participants.**    240 participants were recruited from Amazon Mechanical Turk. Half of the participants were assigned to a condition in which the critical feature was color (red, blue, purple, or green), and the other half of participants were assigned to a condition in which the critical feature was size (small or big).

**Stimuli & Procedure.**    Stimulus displays showed two alien interlocutors, one on the left side (Alien A) and one on the right side (Alien B) of the screen, each with two novel fruit objects beneath them (Figure 4). Alien A, in a speech bubble, asked Alien B for one of its fruits (e.g., "Hey, pass me the big toma.") Alien B replied, "Here you go!" and the referent disappeared from Alien B's side and reappeared on Alien A's side.

We manipulated the critical feature type (color or size) between subjects. Two factors (presence of the critical adjective in the referring expression and object context) were fully crossed within subjects. Object context had three levels: within-category contrast, between-category contrast, and same feature (Figure 5). In the within-category contrast

condition, Alien B possessed the target object and another object of the same shape, but with a different value of the critical feature (e.g., a big toma and a small toma). In the between-category contrast condition, Alien B possessed the target object and another object of a different shape, and with a different value of the critical feature (e.g., a big toma and a small dax). In the same feature condition, Alien B possessed the target object and another object of a different shape but with the same value of the critical feature as the target (e.g., a big toma and a big dax). Thus, in the within-category contrast condition, the descriptor was necessary to distinguish the referent; in the between-category contrast condition it was unnecessary but potentially helpful; and in the same feature condition it was unnecessary and unhelpful.

Note that in all context conditions, the set of objects onscreen was the same in terms of the experiment design: there was a target (e.g., big toma), an object with the same shape as the target and a different critical feature (e.g., small toma), an object with a different shape from the target and the same critical feature (e.g., big dax), and an object with a different shape from the target and a different critical feature (e.g., small dax). Context was manipulated by rearranging these objects such that the relevant referents (the objects under Alien B) differed and the remaining objects were under Alien A. Thus, in each case, participants saw the target object and one other object that shared the target object's shape but not its critical feature–they observed the same kind of feature distribution of the target object's category in each trial type. The particular values of the features were randomly chosen at each trial.

Participants completed six trials. After each exchange between the alien interlocutors, they made a judgment about the prevalence of the target's critical feature in the target object's category. For instance, after seeing a red blicket being exchanged, participants would be asked, "On this planet, what percentage of blickets do you think are red?" and answer on a sliding scale between zero and 100. In the size condition, participants were

asked, "On this planet, what percentage of blickets do you think are the size shown below?" with an image of the target object they just saw available on the screen.

After completing the study, participants were asked to select which of a set of alien words they had seen previously during the study. Four were words they had seen, and four were novel lure words. Participants were dropped from further analysis if they did not respond to at least 6 of these 8 correctly (above chance performance as indicated by a one-tailed binomial test at the $p = .05$ level). This resulted in excluding 47 participants, leaving 193 for further analysis.
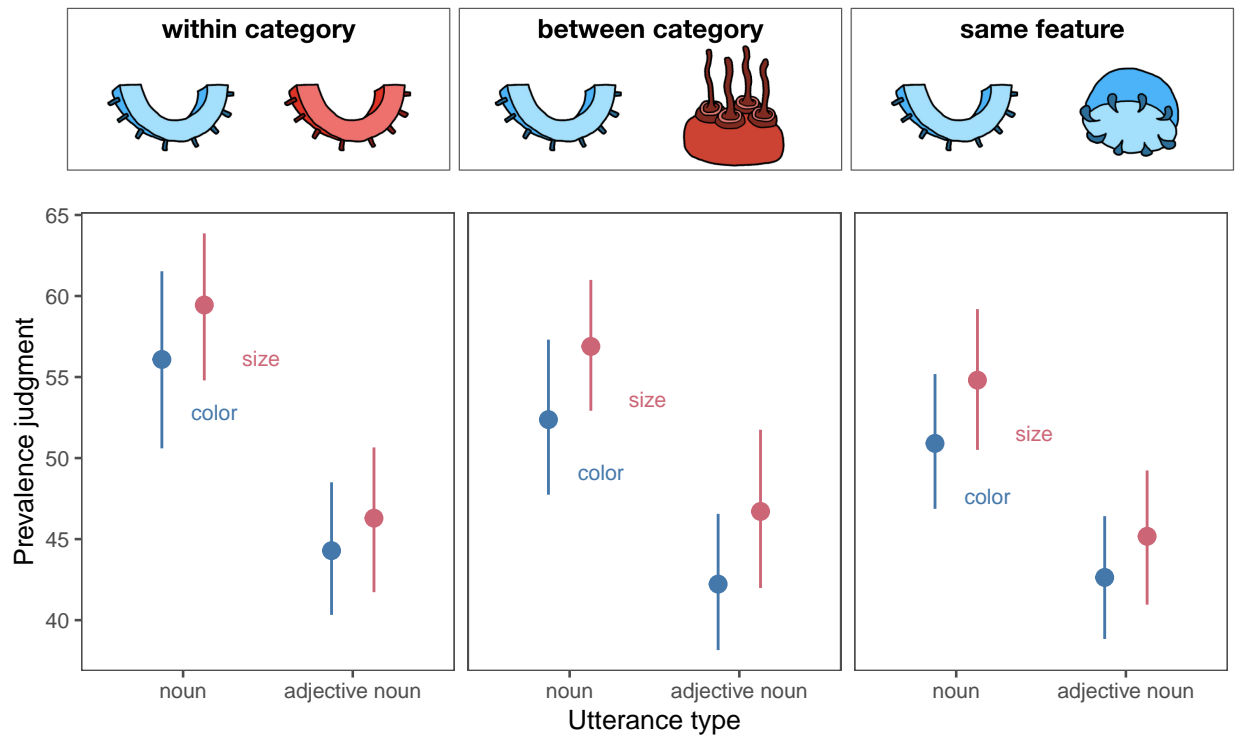
**Results**



*Figure 5*. Prevalence judgments from Experiment 2. Participants consistently judged the target object as less typical of its category when the referent was described with an adjective (e.g., "Pass me the blue toma") than when it was not (e.g., "Pass me the toma"). This inference was not significantly modulated by object context (examples shown above each figure panel).

422    We analyzed participants' judgments of the prevalence of the target object's critical

423 feature in its category. We began by fitting a maximum mixed-effects linear model with

424 effects utterance type (adjective or no adjective), context type (within category, between

425 category, or same feature), and critical feature (color or size) as well as all interactions and

426 random slopes of utterance type and context type nested within subject. Random effects

427 were removed until the model converged. The final model revealed a significant effect of

428 utterance type ($\beta_{adjective}$ = -11.80, $t$ = -3.90, $p < .001$), such that prevalence judgments were

429 lower when an adjective was used than when it was not. Participants also made lower

430 prevalence judgments in the same-feature context type relative to within-category context

431 type ($\beta_{same}$ = -5.41, $t$ = -2.25, $p = .025$), but there was no significant effect of

432 between-category relative to within-category contexts ($\beta_{between}$ = -3.92, $t$ = -1.63, $p = .104$).

433 There was not a significant interaction between context and presence of an adjective in the

434 utterance ($\beta_{same*adjective}$ = 3.71, $t$ = 1.09, $p = .277$; $\beta_{between*adjective}$ = 1.58, $t$ = 0.46, $p =$

435 .644). That is, participants slightly adjusted their inferences according to the object context,

436 though not in a way that depended on whether an adjective was used in the utterance.

437 However, they robustly inferred that described features were less prevalent in the target's

438 category than unmentioned features.

**Discussion**

440    Description is often used not to distinguish among present objects, but to pick out an

441 object's feature as atypical of its category. In Experiment 2, we asked whether people would

442 infer that a described feature is atypical of a novel category after hearing it mentioned in an

443 exchange. We found that people robustly inferred that a mentioned feature was atypical of

444 its category, across both size and color description. Further, participants did not use object

445 context to substantially explain away description. That is, even when description was

446 necessary to distinguish among present objects (e.g., there were two same-shaped objects

447 that differed only in the mentioned feature), participants still inferred that the feature was

atypical of its category. This suggests that, in the case of hearing someone ask for a "red tomato" from a bin of many-colored heirloom tomatoes, a person naive about tomatoes would infer that tomatoes are relatively unlikely to be red.

Unlike Experiment 1, in which people made stronger contrastive inferences for size than color, there were not substantial differences between people's inferences about color and size in Experiment 2. If an account based on production norms is correct, this suggests that people do not only track how often people use color compared to size description but also for what purpose–contrasting with present objects or with the referent's category. That is, color description may be more likely to be used superfluously with respect to present objects but informatively with respect to the category. Indeed, color description that seems overdescriptive with respect to object context often occurs when the category has many-colored members (e.g., t-shirts) or when the object's color is atypical (Rubio-Fernández, 2016). However, our results are consistent with several potential explanations of the color-size asymmetry (or lack thereof). Future work addressing the source of the color-size asymmetry will need to explain differences in its extent when distinguishing among present objects compared to the referent's category.

## Model

To allow the Rational Speech Act Framework to capture inferences about typicality, we modified the Speaker's utility function to have an additional term: the listener's expected processing difficulty. Speakers may be motivated to help listeners to select the correct referent not just eventually but as quickly as possible. People are both slower and less accurate at identifying atypical members of a category as members of that category (Dale, Kehoe, & Spivey, 2007; Rosch, Simpson, & Miller, 1976). If speakers account for listeners' processing difficulties, they should be unlikely to produce bare nouns to refer to low typicality exemplars (e.g. unlikely to call a purple carrot "carrot"). This is roughly the kind of inference encoded in Degen et al. (2020)'s continuous semantics Rational Speech Act model.

We model the speaker as reasoning about the listener's label verification process. Because the speed of verification scales with the typicality of a referent, a natural way of modeling it is as a process of searching for that particular referent in the set of all exemplars of the named category, or alternatively of sampling that particular referent from the set of all exemplars in that category, $P\left(r|Cat\right)$. On this account, speakers want to provide a modifying adjective for atypical referents because the probability of sampling them from their category is low, but the probability of sampling them from the modified category is much higher.[1] Typicality is just one term in the speaker's utility, and thus is directly weighed with the literal listener's judgment and against cost.

If speakers use this utility function, a listener who does not know the feature distribution for a category can use a speaker's utterance to infer it. Intuitively, speakers should prefer not to modify nouns with adjectives because they incur a cost for producing that adjective. If they did, it must be because they thought the learner would have a difficult time finding the referent from a bare noun alone because of typicality, competing referents, or both. To infer the true prevalence of the target feature in the category, learners combine the speaker's utterance with their prior beliefs about the feature distribution. We model the learner's prior about the prevalance of features in any category as a Beta distribution with two parameters $\alpha$ and $\beta$ that encode the number of hypothesized prior psuedo-exemplars with the feature and without feature that the learner has previously observed (e.g. one red dax and one blue dax). We assume that the learner believes they have previously observed one hypothetical psuedo-examplar of each type, which is a weak symmetric prior indicating that the learner expects features to occur in half of all members of a category on average, but would find many levels of prevalence unsurprising. To model the learner's direct experience with the category, we add the observed instances in the experiment to these

———

[1] This is a generalization of the size principle (Xu & Tenenbaum, 2007) to categories where exemplars are not equally likely.

hypothesized prior instances. After observing one member of the target category with the

relevant feature and one without, the listeners prior is thus updated to be Beta (2, 2).

As in Experiment 1, we used Empirical Bayesian methods to estimate the rationality

parameter that participants are using to draw inferences about speakers in both the color

and size conditions. In contrast to Experiment 1, the absolute values of these parameters are

driven largely by the number of pseudo-exemplars assumed by the listener prior to exposure.

Thus, the rationality parameters inferred in the two experiments are not directly comparable.

However, differences between color and size within each model are interpretable. As in

Experiment 1, we found that listeners inferred speakers to be more rational when using size

adjectives 0.89 [0.63, 1.13] than color adjectives 0.60 [0.37, 0.83], but the two inferred

confidence intervals were overlapping, suggesting that people treated the adjective types as

more similar to each other when making inferences about typicality than when making

inferences about reference.

Figure 6 shows the predictions of our Rational Speech Act model compared to

empirical data from participants. The model captures the trends in the data correctly,

inferring that the critical feature was less prevalent in the category if it is referred to with an

adjective (e.g., "red dax") than if it was not mentioned (e.g., "dax"). The model also infers

the prevalence of the critical feature to be numerically more likely in the within-category

condition, like people do. That is, in the within-category condition when an adjective is used

to distinguish between referents, the model thinks that the target color is slightly less

atypical. When an adjective would be useful to distinguish between two objects of the same

shape but one is not used, the model infers that the color of the target object is more typical.

**Discussion**

In contrast to the reference-first view that these two kinds of inferences trade off

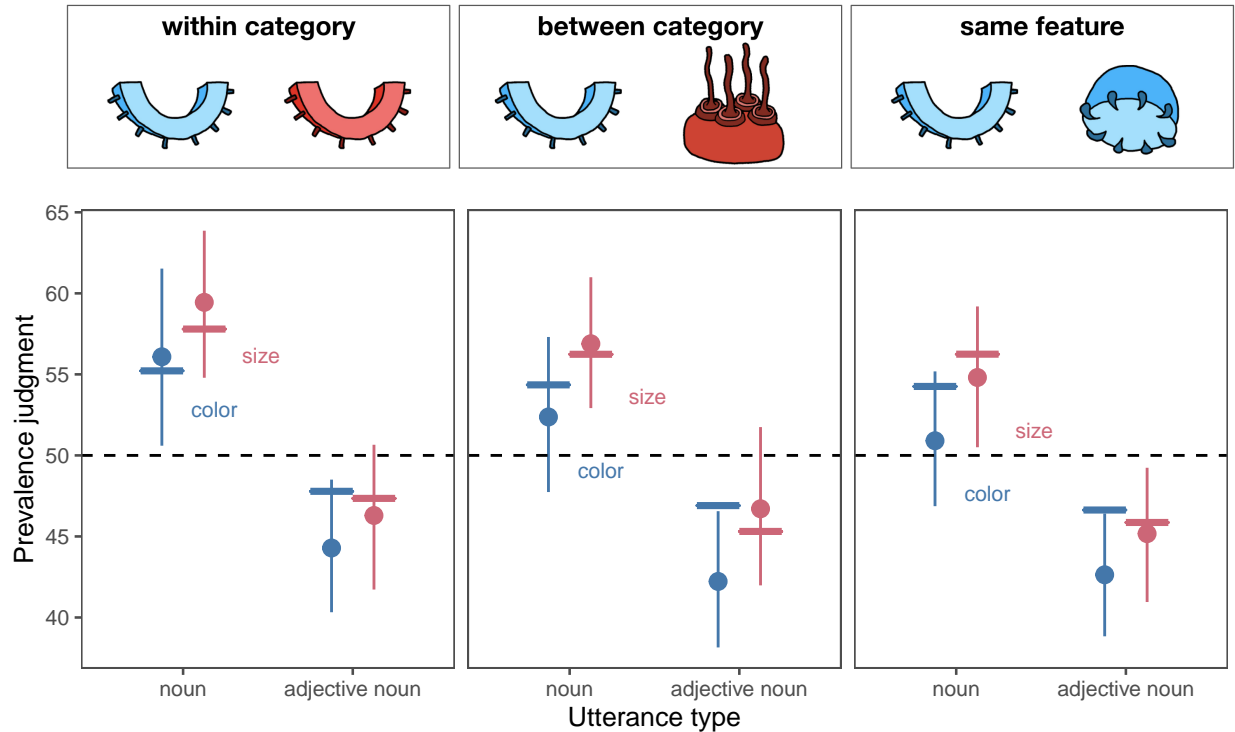strongly–that is, adjectives are used primarily for reference, and such use blocks the inference

*Figure 6*. Participants' prevalence judgments from Experiment 2, compared to model predictions (horizontal lines).

that they are marking typicality–the model captures the graded way in which people interpolate between them. When an adjective is helpful for reference, whether it is used or not makes both the model and people give it slightly less weight in inferring the typical features of the target object, but the weight is still significant. Our model's explanation for this is that while people choose their language in order to refer successfully, their choices also reflect their knowledge of features of those objects. In the model as constructed, we cannot distinguish between listener and speaker design explanations for the impact of feature knowledge. One possibility is that the pressure from this feature knowledge is communicative: speakers could be intentionally transmitting information to the listener about the typical features of their intended referent. Alternatively, the influence of this feature knowledge could be unintentional, driven by pressures from the speaker's semantic representation. We consider these implications more fully in the General Discussion. In

either case, listeners can leverage the impact of speakers' feature knowledge on their productions in order to infer the typical features of the objects they are talking about, even if this is their first exposure to these novel objects.

## Experiment 3

In Experiments 1 and 2, we established that people can use contrastive inferences to resolve referential ambiguity and to make inferences about the feature distribution of a novel category. Additionally, in Experiment 2, we found that these two inferences do not seem to trade off substantially: even if an adjective is necessary to establish reference, people infer that it also marks atypicality. We also found that inferences of atypicality about color and size adjectives pattern very similarly, though their baseline typicality is shifted, while color and size are not equally contrastive with respect to referential disambiguation.

To strengthen our findings in a way that would allow us to better detect potential trade-offs between these two types of inference, we conducted a pre-registered replication of Experiment 2 with a larger sample of participants. In addition, we test how people's prevalence judgments from utterances with and without an adjective compare to their null inference about feature prevalence by adding a control utterance condition: an alien utterance, which the participants cannot understand. This also tests the model assumption we made in Experiment 2: that after seeing two exemplars of the target object with two values of the feature (e.g., one green and one blue), people's prevalence judgments would be around 50%. In addition to validating this model assumption, we more strongly test the model here by comparing predictions from same model, with parameters inferred from the Experiment 2 data, to data from Experiment 3. Our pre-registration of the method, recruitment plan, exclusion criteria, and analyses can be found on the Open Science Framework here: https://osf.io/s8gre .

## Method

**Participants.**    A pre-registered sample of four hundred participants were recruited from Amazon Mechanical Turk. Half of the participants were assigned to a condition in which the critical feature was color (red, blue, purple, or green), and half of the participants were assigned to a condition in which the critical feature was size (small or big).

**Stimuli & Procedure.**    The stimuli and procedure were identical to those of Experiment 2, with the following modifications. Two factors, utterance type and object context, were fully crossed within subjects. Object context had two levels: within-category contrast and between-category contrast. In the within-category context condition, Alien B possessed the target object and another object of the same shape, but with a different value of the critical feature (color or size). In the between-category contrast condition, Alien B possessed the target object and another object of a different shape, and with a different value of the critical feature. Thus, in the within-category contrast condition, the descriptor is necessary to distinguish the referent; in the between-category contrast condition it is unnecessary but potentially helpful. There were three utterance types: adjective, no adjective, and alien utterance. In the two alien utterance trials, the aliens spoke using completely unfamiliar utterances (e.g., "Zem, noba bi yix blicket"). Participants were told in the task instructions that sometimes the aliens would talk in a completely alien language, and sometimes their language will be partly translated into English. To keep participants from making inferences about the content of the alien utterances using the utterance content of other trials, both alien language trials were first; other than this constraint, trial order was random. We manipulated the critical feature type (color or size) between subjects.

After completing the study, participants were asked to select which of a set of alien words they had seen previously during the study. Four were words they had seen, and four were novel lure words. Participants were dropped from further analysis if they did not meet our pre-registered criteria of responding to at least 6 of these 8 correctly (above chance

performance as indicated by a one-tailed binomial test at the $p = .05$ level) and answering all
four color perception check questions correctly. Additionally, six participants were excluded
because their trial conditions were not balanced due to an error in the run of the experiment.
This resulted in excluding 203 participants, leaving 197 for further analysis.

**Results**

We began by fitting a pre-registered maximum mixed-effects linear model: effects
utterance type (alien utterance, adjective, or no adjective; alien utterance as reference level),
context type (within category or between category), and critical feature (color or size) as
well as all interactions and random slopes of utterance type and context type nested within
subject. Random effects were removed until the model converged, which resulted in a model
with all fixed effects, all interactions and a random slope of utterance type by subject. The
final model revealed a significant effect of the no adjective utterance type compared to the
alien utterance type ($\beta = 13.05$, $t = 4.88$, $p = < .001$) and a marginal effect of the adjective
utterance type compared to the alien utterance type ($\beta = 5.13$, $t = 1.95$, $p = .052$). The
effects of context type (within-category or between-category) and adjective type (color or
size) were not significant ($\beta_{between} = 2.70$, $t_{between} = 1.23$, $p_{between} = .220$; $\beta_{size} = 5.68$, $t_{size} = 1.70$, $p_{size} = .090$). There was a significant interaction between the adjective utterance type
and the size condition ($\beta = -8.78$, $t = -2.31$, $p = .022$). Thus, participants inferred that an
object referred to in an intelligible utterance with no description was more typical of its
category on the target feature than an object referred to with an alien utterance. They also
inferred that an object referred to in an intelligible utterance with description was marginally
less typical than an object referred to with an alien utterance, and this effect was slightly
stronger in the size condition. Participants did not substantially adjust their inferences
based on the object context.

Given that interpretation of these results with respect to the alien utterance condition
can be difficult, we pre-registered a version of the same full model excluding alien utterance

611  trials with the no adjective utterance type as the reference level. This model revealed a

612  significant effect of utterance type: participants' prevalence judgments were lower when an

613  adjective was used than when it was not ($\beta$ = -7.92, $t$ = -3.38, $p$ = .001). No other effects

614  were significant. This replicates the main effect of interest in Experiment 2: that when an

615  adjective is used in referring to the object, participants infer that the described feature is less

616  typical of that object's category than when the feature goes unmentioned. In the

617  supplemental materials, we report two more pre-registered tests of the effect of utterance

618  type alone on prevalence judgments, whose results are consistent with the fuller models
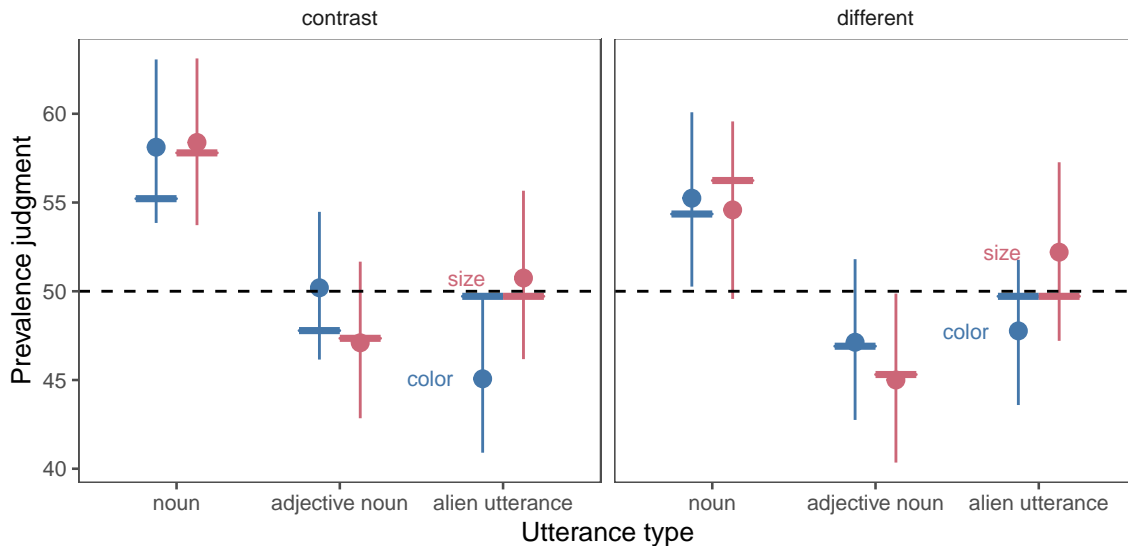
619  reported here.



*Figure 7*. Participants' prevalence judgments in Experiment 3, with model predictions using
the parameters estimated in Experiment 2 (horizontal lines).

620        To validate the model we developed for Experiment 2, we compared its estimates using

621  the previously fit parameters to the new data for Experiment 3. As show in Figure 7, the

622  model predictions were well aligned with peoples' prevalence judgments. In addition, in

623  Experiment 2, we fixed the model's prior beliefs about the prevalence of the target object's

624  color or size to be centered at 50% because the model had seen one pseudo-exemplar of the

625  target color/size, and on psuedo-exemplar of the non-target color/size. In Experiment 3, we

aimed to estimate this prior empirically in the alien utterance condition, reasoning that

people could only use their prior to make a prevalence judgment (as we asked the model to

do). In both the color and size conditions, peoples' judgments were indeed around 50%,

although in the color condition they were directionally lower. This small effect may arise

from a fundamental difference between polar adjectives like size (where objects can be big or

small) and adjectives like color where there may be many nameable alternatives (e.g. red,

blue, green, etc.). Thus, the results of Experiment 3 confirm the modeling assumptions we

made in estimating peoples' prior beliefs, and further validate the model we developed as a

good candidate model for how people simultaneously draw inferences about speakers'

intended referents and the typicality of these referents. That is, when people think about

why a speaker chose their referring expression, they think about not only the set of present

objects as providing the context of referents, but also the broader set of categories that they

belong to.

**Discussion**

In Experiment 3, we replicated the main finding of interest in Experiment 2: when a

novel object's feature is described, people infer that the feature is rarer of its category than

when it goes unmentioned. Again, this effect was consistent across both size and color

adjectives, and people did not substantially adjust this inference based on how necessary the

description was to distinguish among potential referents. We also added an alien language

condition, in which the entire referring expression was unintelligible to participants, to probe

people's priors on feature typicality. We found that in the alien language condition, people

judged features to be roughly between the adjective utterance and no adjective utterance

conditions, and significantly different from the no adjective utterance condition. In the alien

language condition, people's prevalence judgments were roughly around our model's

prevalence judgments (50%) after observing the objects on each trial and before any

inferences about the utterance.

⁶⁵² The similarity of people's prevalence judgments in the alien language condition and the

⁶⁵³ adjective condition raises the question: is this effect driven by an atypicality inference in the

⁶⁵⁴ adjective conditions, or a *typicality* effect when the feature is unmentioned? Our results

⁶⁵⁵ suggest that it is a bit of both. When someone mentions an object without extra description,

⁶⁵⁶ the listener can infer that its features are likely more typical than their prior; when they use

⁶⁵⁷ description, they can infer that its features are likely less typical. Because using an extra

⁶⁵⁸ word–an adjective–is generally not thought of as the default way to refer to something, this

⁶⁵⁹ effect is still best described as a contrastive inference of *atypicality* when people use

⁶⁶⁰ description. However, the fact that people infer high typicality when an object is referred to

⁶⁶¹ without description suggests that, in some sense, there is no neutral way to refer: people will

⁶⁶² make broader inferences about a category from even simple mentions of an object.


⁶⁶³ **General Discussion**


⁶⁶⁴ When we think about what someone is trying to communicate to us, we go far beyond

⁶⁶⁵ the literal meanings of the words they say. Instead, we make pragmatic inferences about why

⁶⁶⁶ they chose those particular words rather than other words they could have used instead.

⁶⁶⁷ This kind of reasoning allows us to draw scalar implicatures (e.g. "some" means "some but

⁶⁶⁸ not all"), identify negative beliefs from the absence of positive language in recommendation

⁶⁶⁹ letters, and to make the kind of typicality inferences we studied here. In most work on

⁶⁷⁰ pragmatic reasoning, speakers and listeners share the same knowledge of language, and the

⁶⁷¹ question of interest is whether listeners can use their knowledge of language to learn

⁶⁷² something about the unknown state of the world. Here we focus on an even more challenging

⁶⁷³ problem: Can pragmatic inference be used to learn about language and the world

⁶⁷⁴ simultaneously?


⁶⁷⁵ In three studies we showed that people can use pragmatic inference to (1) learn the

⁶⁷⁶ meaning of a novel word, (2) learn the typical features of the category described by this

⁶⁷⁷ novel word, and (3) rationally integrate these two kinds of reasoning processes. In

Experiment 1, we show that people can use descriptive contrast implied by adjectives like "big" or "blue" to resolve referential ambiguity to learn a new word; in the case of color, they shift substantially in the direction of the correct mapping, and in the case of size, they choose the correct mapping significantly more often than the incorrect one. In Experiments 2 and 3, we show that people use the presence of the same kind of descriptor to infer that the noted feature is atypical of the object being referred to. Critically, people infer that the described feature is atypical even when the descriptor is helpful for referential disambiguation—although the size of the atypicality inference is slightly reduced.

Why do people think that the mentioned feature is atypical even when its mention is helpful for referential disambiguation? If people use language for multiple goals–for example, both for reference and for description– then listeners should reason jointly about all of the possible reasons why speakers could have used a word when they hear it. To determine what rational listeners would do in this circumstance, we developed an extension of the Rational Speech Act Framework that reasons both about reference and about the typical features of categories to which objects belong. The behavior of this model was closely aligned to the behavior we observed from people. Because rational inference is probabilistic rather than deterministic, descriptors still lead to atypicality inferences even when they are helpful for referential disambiguation. This work thus adds to the growing body of work extending the Rational Speech Act framework from reasoning about just reference to reasoning about other goals as well, such as inferring that speech is hyperbolic (e.g. waiting "a million years" means waiting a long time), inferring when speakers are being polite rather than truthful, and learning new words in ambiguous contexts (Frank & Goodman, 2014; Goodman & Frank, 2016; Kao, Wu, Bergen, & Goodman, 2014; Yoon, Tessler, Goodman, & Frank, 2020).

Though the participants in our experiments were adults, the ability to disambiguate novel referents using contrast most obviously serves budding language learners: children. Contrastive use of adjectives is a pragmatic regularity in language that children could

704 potentially exploit to establish word–referent mappings. Use of adjectives has been shown to

705 allow children to make contrastive inferences among familiar present objects (Davies,

706 Lingwood, Ivanova, & Arunachalam, 2021; Huang & Snedeker, 2008). When paired with

707 other contrastive cues such as prosody, preschoolers can make inferences about novel object

708 typicality (Horowitz & Frank, 2016), and can use novel adjectives and nouns to restrict

709 reference (Diesendruck, Hall, & Graham, 2006; Gelman & Markman, 1985). Future work

710 should explore whether adjective contrast that is less scaffolded by other cues is a viable way

711 for children to learn about novel concepts.

712      The core computation in pragmatic inference is reasoning about alternatives–things the

713 speaker could have said and did not. Given that others are reasoning about these

714 alternatives, no choice is neutral. In the studies in this paper, for instance, using an adjective

715 in referring to an object led people to infer that the feature described by that adjective was

716 less typical than if it had not been mentioned. But, conversely, *not* using an adjective led

717 them to think that the feature was more typical than if they could not understand the

718 meaning of the utterance at all–all communicative choices leak one's beliefs about the world.

719 This has implications not only for learning about novel concrete objects, as people did here,

720 but for learning about less directly accessible entities such as abstract concepts and social

721 groups. These inferences can be framed positively, as ways for learners to extract additional

722 knowledge that was not directly conveyed, but can also spread beliefs that the speaker does

723 not intend. A core challenge will be to understand how people reason about the many

724 potential meanings a speaker might convey in naturalistic contexts to learn about others'

725 words for and beliefs about the world.

## Acknowledgements

# References

Arts, A., Maes, A., Noordman, L. G. M., & Jansen, C. (2011). Overspecification in written instruction. *Linguistics*, *49*(3), 555–574.

Bergey, C., Morris, B., & Yurovsky, D. (2020). *Children hear more about what is atypical than what is typical.* PsyArXiv. https://doi.org/10.31234/osf.io/5wvu8

Dale, R., Kehoe, C., & Spivey, M. J. (2007). Graded motor responses in the time course of categorizing atypical exemplars. *Memory & Cognition*, *35*(1), 15–28.

Davies, C., Lingwood, J., Ivanova, B., & Arunachalam, S. (2021). Three-year-olds' comprehension of contrastive and descriptive adjectives: Evidence for contrastive inference. *Cognition*, *212*, 104707. https://doi.org/10.1016/j.cognition.2021.104707

Degen, J., Hawkins, R. D., Graf, C., Kreiss, E., & Goodman, N. D. (2020). When redundancy is useful: A bayesian approach to "overinformative" referring expressions. *Psychological Review.*

Diesendruck, G., Hall, D. G., & Graham, S. A. (2006). Children's Use of Syntactic and Pragmatic Knowledge in the Interpretation of Novel Adjectives. *Child Development*, *77*(1), 16–30.

Engelhardt, P. E., Barış Demiral, Ş., & Ferreira, F. (2011). Over-specified referring expressions impair comprehension: An ERP study. *Brain and Cognition*, *77*(2), 304–314. https://doi.org/10.1016/j.bandc.2011.07.004

Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, *336*(6084), 998–998.

Frank, M. C., & Goodman, N. D. (2014). Inferring word meanings by assuming that speakers are informative. *Cognitive Psychology*, *75*, 80–96.

Gelman, S. A., & Markman, E. M. (1985). Implicit contrast in adjectives vs. Nouns: Implications for word-learning in preschoolers*. *Journal of Child Language*, *12*(1), 125–143.

Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, *20*(11), 818–829.

Grice, H. P. (1975). Logic and conversation. *1975*, 41–58.

Horowitz, A. C., & Frank, M. C. (2016). Children's Pragmatic Inferences as a Route for Learning About the World. *Child Development*, *87*(3), 807–819.

Huang, Y. T., & Snedeker, J. (2008). Use of referential context in children's language processing. *Proceedings of the 30th Annual Meeting of the Cognitive Science Society.*

Kanwisher, N., Woods, R. P., Iacoboni, M., & Mazziotta, J. C. (1997). A locus in human extrastriate cortex for visual shape analysis. *Journal of Cognitive Neuroscience*, *9*(1), 133–142.

Kao, J. T., Wu, J. Y., Bergen, L., & Goodman, N. D. (2014). Nonliteral understanding of number words. *Proceedings of the National Academy of Sciences*, *111*(33), 12002–12007.

Landau, B., Smith, L. B., & Jones, S. (1992). Syntactic context and the shape bias in children's and adults' lexical learning. *Journal of Memory and Language*, *31*(6), 807–825.

Mitchell, M., Reiter, E., & Deemter, K. van. (2013). Typicality and Object Reference, 7.

Nadig, A. S., & Sedivy, J. C. (2002). Evidence of Perspective-Taking Constraints in Children's On-Line Reference Resolution. *Psychological Science*, *13*(4), 329–336.

Pechmann, T. (1989). Incremental speech production and referential overspecification. *Linguistics*, *27*(1), 89–110.

Rosch, E., Simpson, C., & Miller, R. S. (1976). Structural bases of typicality effects. *Journal of Experimental Psychology: Human Perception and Performance*, *2*(4), 491.

Rubio-Fernández, P. (2016). How Redundant Are Redundant Color Adjectives? An Efficiency-Based Analysis of Color Overspecification. *Frontiers in Psychology*, *7*.

Sedivy, J. C. (2003). Pragmatic Versus Form-Based Accounts of Referential Contrast: Evidence for Effects of Informativity Expectations. *Journal of Psycholinguistic Research*, *32*(1), 3–23.

Sedivy, J. C., Tanenhaus, M. K., Chambers, C. G., & Carlson, G. N. (1999). Achieving incremental semantic interpretation through contextual representation. *Cognition*, *71*(2), 109–147.

Westerbeek, H., Koolen, R., & Maes, A. (2015). Stored object knowledge and the production of referring expressions: The case of color typicality. *Frontiers in Psychology*, *6*. https://doi.org/10.3389/fpsyg.2015.00935

Xu, F., & Tenenbaum, J. B. (2007). Word learning as bayesian inference. *Psychological Review*, *114*(2), 245.

Yoon, E. J., Tessler, M. H., Goodman, N. D., & Frank, M. C. (2020). Polite speech emerges from competing social goals. *Open Mind*, *4*, 71–87.