

# Over overinformativeness: rationally redundant referring expressions

Judith Degen, Robert X.D. Hawkins, Caroline Graf, Elisa Kreiss, & Noah D. Goodman

{jdegen}@stanford.edu

Department of Psychology, 450 Serra Mall  
Stanford, CA 94305 USA

March 23, 2018

## Abstract

Referring is one of the most basic and prevalent uses of language. How do speakers choose from the wealth of referring expressions at their disposal? Rational theories of language use have come under attack for decades for not being able to account for the seemingly irrational overinformativeness ubiquitous in referring expressions. Here we present a novel production model of referring expressions within the Rational Speech Act framework that treats speakers as agents that rationally trade off cost and informativeness of utterances. Crucially, the assumption of deterministic meanings is relaxed. This allows us to capture a large number of seemingly disparate phenomena within one unified framework: the basic asymmetry in speakers' propensity to overmodify with color rather than size; the increase in overmodification in complex scenes; the increase in overmodification with atypical features; and the preference for basic level reference in nominal reference. The findings cast a new light on the production of referring expressions: rather than being wastefully overinformative, reference is rationally redundant. This implicates a production system geared towards communicative efficiency.

**Keywords:** reference; referring expressions; informativeness; probabilistic pragmatics; experimental pragmatics

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>3</b>  |
| 1.1      | Production of referring expressions: a case against rational language use? | 4         |
| 1.2      | Modified referring expressions   | 5         |
| 1.2.1    | Asymmetry in redundant use of color and size adjectives                    | 5         |
| 1.2.2    | Scene variation  | 8         |
| 1.2.3    | Feature typicality   | 8         |
| 1.3      | Nominal referring expressions  | 9         |
| 1.4      | Summary  | 10        |
| <b>2</b> | <b>Modeling speakers' choice of referring expression</b>                   | <b>11</b> |
| 2.1      | Basic RSA  | 11        |
| 2.2      | RSA with continuous semantics – emergent color-size asymmetry              | 14        |
| 2.3      | RSA with continuous semantics – scene variation                            | 15        |

|  |           |
|--|-----------|
| <b>3 Modified referring expressions: size and color modifiers under different scene variation conditions</b> | <b>18</b> |
| 3.1 Experiment 1: scene variation in modified referring expressions . . . . .                                | 18        |
| 3.1.1 Method . . . . .   | 19        |
| 3.1.2 Data pre-processing and exclusion . . . . .  | 21        |
| 3.1.3 Results . . . . .  | 22        |
| 3.2 Model evaluation . . . . .   | 23        |
| 3.3 Discussion . . . . .   | 25        |
| <b>4 Modified referring expressions: color typicality</b>  | <b>25</b> |
| 4.1 Experiment 2: color typicality effects . . . . .   | 27        |
| 4.1.1 Method . . . . .   | 27        |
| 4.1.2 Data pre-processing and exclusion . . . . .  | 28        |
| 4.1.3 Typicality norming . . . . .   | 29        |
| 4.1.4 Results and discussion . . . . .   | 31        |
| 4.2 Model evaluation . . . . .   | 32        |
| 4.2.1 Lexicon . . . . .  | 32        |
| 4.2.2 Cost function . . . . .  | 33        |
| 4.2.3 Evaluation . . . . .   | 33        |
| 4.3 Discussion . . . . .   | 34        |
| <b>5 Unmodified referring expressions: nominal taxonomic level</b>   | <b>34</b> |
| 5.1 Experiment 3: taxonomic level of reference in nominal referring expressions . . . . .                    | 34        |
| 5.1.1 Method . . . . .   | 35        |
| 5.1.2 Data pre-processing and exclusion . . . . .  | 35        |
| 5.1.3 Results and discussion . . . . .   | 36        |
| 5.2 Model evaluation . . . . .   | 37        |
| 5.2.1 Typicality effects . . . . .   | 38        |
| 5.2.2 Cost effects . . . . .   | 39        |
| 5.3 Model evaluation: nominal choice . . . . .   | 40        |
| <b>6 General Discussion</b>  | <b>40</b> |
| 6.1 Summary . . . . .  | 40        |
| 6.2 ‘Overinformativeness’ . . . . .  | 42        |
| 6.3 Comprehension . . . . .  | 42        |
| 6.4 Fidelity . . . . .   | 43        |
| 6.5 Audience design . . . . .  | 44        |
| 6.6 Other factors affecting redundancy . . . . .   | 45        |
| 6.7 Extensions to other language production phenomena . . . . .  | 45        |
| 6.8 Conclusion . . . . .   | 46        |
| <b>A Effects of semantic value on utterance probabilities</b>  | <b>46</b> |
| <b>B Validation of interactive web-based written production paradigm</b>                                     | <b>47</b> |
| <b>C Pre-experiment quiz</b>   | <b>47</b> |

|   |           |
|---|-----------|
| <b>D Exp. 1 items</b>                                     | <b>47</b> |
| <b>E Typicality effects in Exp. 1</b>                     | <b>48</b> |
| E.1 Methods . . . . .                                     | 48        |
| E.1.1 Participants . . . . .                              | 48        |
| E.1.2 Procedure and materials . . . . .                   | 48        |
| E.2 Results and discussion . . . . .                      | 48        |
| <b>F Experiment 3a: typicality norms for Experiment 3</b> | <b>50</b> |
| F.0.1 Methods . . . . .                                   | 50        |
| F.0.2 Results and discussion . . . . .                    | 51        |
| <b>G Experiment 3 items</b>                               | <b>51</b> |
| <b>H Nominal choice model comparison</b>                  | <b>51</b> |
| <b>References</b>   | <b>53</b> |

## 1 Introduction

Reference to objects is one of the most basic and prevalent uses of language. This requires speakers to choose from amongst a wealth of referring expressions they have at their disposal. How does a speaker choose whether to refer to an object as *the animal*, *the dog*, *the dalmatian*, or *the big mostly white dalmatian*? The context within which the object occurs (other non-dogs, other dogs, other dalmatians) plays a large part in determining which features the speaker chooses to include in their utterance – speakers aim to be sufficiently informative to uniquely establish reference to the intended object. However, speakers’ utterances often exhibit what has been claimed to be *overinformativeness*: referring expressions are often more specific than necessary for establishing unique reference, and they are so in systematic ways. Providing a unified theory for speakers’ systematic patterns of overinformativeness has so far proved elusive.

This paper is concerned with accounting for these systematic patterns in overinformative referring expressions (REs). We restrict ourselves to definite descriptions of the form *the (ADJ?) + NOUN*, that is, noun phrases that minimally contain the definite determiner *the* followed by a head noun, with any number of adjectives occurring between the determiner and the noun.<sup>1</sup> A model of these REs will allow us to unify two domains in language production that have been typically treated as separate, and that have typically been treated as interesting for different reasons: the production of so-called overmodified referring expressions on the one hand, which a lot of literature in language production has been devoted to (Herrmann & Deutsch, 1976; Pechmann, 1989; Nadig & Sedivy, 2002; Sedivy, 2003; Maes, Arts, & Noordman, 2004; Engelhardt, Bailey, & Ferreira, 2006a; Arts, Maes, Noordman, & Jansen, 2011; Koolen, Gatt, Goudbeek, & Krahmer, 2011; Rubio-Fernandez, 2016); and the production of simple nominal expressions, which has so far mostly received attention in the concepts and categorization literature (Rosch, 1973; Rosch, Mervis, Gray, Johnson, & Boyes-Braem, 1976) and in the developmental literature on generalizing

---

<sup>1</sup>In contrast, we will *not* provide a treatment of pronominal referring expressions, indefinite descriptions, names, or definite descriptions with post-nominal modification, though we offer some speculative remarks on how the approach outlined here can be applied to these cases. [jd: make sure to pick this back up in the discussion]

basic level terms (? , ?). In the following, we review some of the key phenomena and puzzles in each of these literatures. We then present a model of RE production within the Rational Speech Act framework (M. C. Frank & Goodman, 2012; Goodman & Frank, 2016), which treats speakers as boundedly rational agents who optimize the tradeoff between utterance cost and informativeness. Our key innovation is to relax the assumption that semantic truth functions are deterministic. Treating speakers as boundedly rational agents operating on a continuous semantics captures that adding seemingly overinformative modifiers or using nouns that are seemingly too specific can be useful and informative, to the extent that not doing so would be too likely to lead a listener to incorrectly infer the speaker’s intention (or to invest too much processing effort in inferring the speaker’s correct intention). We thus for the first time provide a unified explanation for a number of seemingly disparate phenomena from the modified and nominal RE literature.

We spend the remainder of the paper demonstrating how the account applies to various phenomena. In Section 1 we spell out the problem and introduce the overinformativeness phenomena to capture. In Section 2 we introduce the basic (deterministic semantics) and modified (continuous semantics) Rational Speech Act framework. In Sections 3 - 5 we evaluate the continuous semantics RSA model on data from interactive online reference game experiments that exhibit the phenomena introduced in Section 1: size and color modifier choice under varying conditions of scene complexity; typicality effects in the choice of color modifier; and choice of nominal level of reference. We wrap up in Section 6 by summarizing our findings and discussing the far-reaching implications of and further challenges for this line of work.

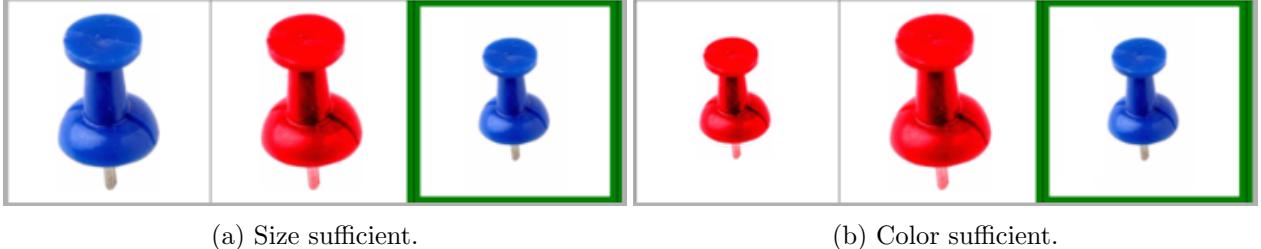
### 1.1 Production of referring expressions: a case against rational language use?

How should a cooperative choose between competing referring expressions? Grice, in his seminal work, provided some guidance by formulating his famous conversational maxims, intended as a guide to listeners’ expectations about good speaker behavior (Grice, 1975). His maxim of Quantity, consisting of two parts, requires of speakers to:

1. *Quantity-1*: Make your contribution as informative as is required (for the purposes of the exchange).
2. *Quantity-2*: Do not make your contribution more informative than is required.

That is, speakers should aim to produce neither under- nor overinformative utterances. While much support has been found for the avoidance of underinformativeness (Brennan & Clark, 1996; R. Brown, 1958; Olson, 1970; Levinson, 1983), speakers seem remarkably happy to systematically violate Quantity-2. In modified referring expressions, they routinely produce modifiers that are not necessary for uniquely establishing reference (e.g., *the small blue pin* instead of *the small pin* in contexts like Figure 1a (Gatt, van Gompel, Krahmer, & van Deemter, 2011; Gatt, Krahmer, van Deemter, & van Gompel, 2014; Arts et al., 2011; Koolen et al., 2011)). In simple nominal expressions, speakers routinely choose to refer to an object with a basic level term even when a superordinate level term would have been sufficient for establishing reference (e.g., *the dog* instead of *the animal* in contexts like Figure 16 (Rosch et al., 1976; Hoffmann & Ziessler, 1983; Tanaka & Taylor, 1991a; Johnson & Mervis, 1997; R. Brown, 1958)).

These observations have posed a challenge for theories of language production, especially those positing rational language use (including the Gricean one): why this extra expenditure of useless



(a) Size sufficient.

(b) Color sufficient.

Figure 1: Example contexts where (a) size only or (b) color only is sufficient for unique reference. A green border marks the intended referent.

effort? Why this seeming blindness to the level of informativeness requirement? Many have argued from these observations that speakers are in fact not economical (Engelhardt et al., 2006a; Pechmann, 1989). Some have derived a built-in preference for referring at the basic level from considerations of perceptual factors such as shape (Rosch et al., 1976; Rosch, 1973; Murphy & Smith, 1982). Others have argued for salience-driven effects on willingness to overmodify (Gatt et al., 2014; Westerbeek, Koolen, & Maes, 2015). In all cases, it is argued that informativeness cannot be the key factor in determining the content of speakers' referring expressions.

Here we revisit this claim and show that systematically relaxing the requirement of a deterministic semantics for referring expressions also systematically changes the informativeness of utterances. This results in a reconceptualization of what have been termed *overinformative referring expressions* as *rationally redundant referring expressions*. We begin by reviewing the phenomena of interest that a revised theory of definite referring expressions should be able to account for.

## 1.2 Modified referring expressions

Most of the literature on overinformative referring expressions has been devoted to the use of overinformative modifiers in modified referring expressions. The prevalent observation is that speakers frequently do not include only the minimal modifiers required for establishing reference, but often also include redundant modifiers (Pechmann, 1989; Nadig & Sedivy, 2002; Maes et al., 2004; Engelhardt et al., 2006a; Arts et al., 2011; Koolen et al., 2011). However, not all modifiers are created equal: there are systematic differences in the overmodification patterns observed for size adjectives (e.g., *big*, *small*), color adjectives (e.g., *blue*, *red*), material adjectives (e.g., *plastic*, *wooden*), and others (Sedivy, 2003). Here we review some key patterns of overmodification that have plagued that literature, before spelling out our account of these phenomena in Section 2.

### 1.2.1 Asymmetry in redundant use of color and size adjectives

In Figure 1a, singling out the object highlighted by the green border requires only mentioning its size (*the small pin*). But it is now well-documented that speakers routinely include redundant color adjectives (*the small blue pin*) which are not necessary for uniquely singling out the intended referent in these kinds of contexts (Pechmann, 1989; Belke & Meyer, 2002; Gatt et al., 2011). However, the same is not true for size: in contexts like Figure 1b, where color is sufficient for unique reference (*the blue pin*), speakers overmodify much more rarely. Table 1 shows proportions of color, size, and (overinformative) color-and-size mentions in conditions like those depicted in

Figure 2 across different experiments. In all cases there is a preference for overmodifying with color but not with size.<sup>2</sup>

Explanations for this asymmetry have varied. Pechmann (1989) was the first to take the asymmetry as evidence for speakers following an incremental strategy of object naming: speakers initially start to articulate an adjective denoting a feature that listeners can quickly and easily recognize (i.e., color) before they have fully inspected the display and extracted the sufficient dimension. However, this would predict that speakers routinely should produce expressions like *the blue small pin*, which violate the preference for size adjectives to occur before color adjectives in English (Bloomfield, 1933; Sproat & Shih, 1991). While Pechmann did observe such violations in his dataset, most cases of overmodification did not constitute such violations, and he himself concluded that incrementality cannot (on its own) account for the asymmetry in speakers' propensity for overmodifying with color vs. size.

Another explanation for the asymmetry is that speakers try to produce modifiers that denote features that are reasonably easy for the listener to perceive, so that, even when a feature is not fully distinguishing in context, it at least serves to restrict the number of objects that could plausibly be considered the target. Indeed, there has been some support for the idea that overmodification can be beneficial to listeners by facilitating target identification (Arts et al., 2011; Rubio-Fernandez, 2016; Paraboni, van Deemter, & Masthoff, 2007). We return to this idea in Section 2 and the General Discussion.

There have been various attempts to capture the color-size asymmetry in computational natural language generation models. The earliest contenders for models of definite referring expressions like the Full Brevity algorithm (Dale, 1989) or the Greedy algorithm (Dale, 1989) focused only on discriminatory value – that is, an utterance's informativeness – in generating referring expressions. This is equivalent to the very simple interpretation of Grice laid out above, and consequently these models demonstrated the same inability to capture the color-size asymmetry: they only produced the minimally specified expressions. Subsequently, the Incremental algorithm (Dale & Reiter, 1995) incorporated a preference order on features, with color ranked higher than size. The order is traversed and each encountered feature included in the expression if it serves to exclude at least one further distractor. This results in the production of overinformative color but not size adjectives. However, the resulting asymmetry is much greater than that evident in human speakers, and is deterministic rather than exhibiting the probabilistic production patterns that human speakers exhibit. More recently, the PRO model (Gatt, van Gompel, van Deemter, & Krahmer, 2013) has sought to integrate the observation that speakers seem to have a preference for including color terms with the observation that a preference does not imply the deterministic inclusion of said color term. The model is specifically designed to capture the color-size asymmetry: in a first step, the uniquely distinguishing property (if there is one) is first selected deterministically. In a second step, an additional property is added probabilistically, depending on both a salience parameter associated with the additional property and a parameter capturing speakers' eagerness to overmodify. If both properties are uniquely distinguishing, a property is selected probabilistically depending on its associated salience parameter. The second step proceeds as before.

However, while the PRO model – the most state-of-the-art computational model of human production of modified referring expressions – can capture the color-size asymmetry in and of itself, it is neither flexible enough to be extended straightforwardly to other modifiers beyond color

---

<sup>2</sup>There is quite a bit of variation in the actual numbers. We will discuss this variation in the Discussion of [jd: we should make a meta-analysis paper out of this instead and get rid of the table]

Table 1: Proportions of minimally informative *color* (only) or *size* (only) and overinformative *color\_size* mentions in color-sufficient vs. size-sufficient conditions across experiments.[jd: keep filling in: arts (Arts2011 examined only comprehension of overspecified expressions, not production; RTs measured (Dutch)), mitchell (Mitchell2013 investigated the effect of typicality on reference, but only along the dimensions of shape and material (English)), koolen 2011 (Koolen2011 used two domains: furniture (where attributes varied along the dimensions color, size and orientation) and faces (where attributes varied along several dimensions, not including size or color); authors did not report proportions for each domain (just overall: 53.6% overspecified, 41.4% minimally specified, 5.0% underspecified); authors reported \*number\* of redundant attributes (0.6 in furniture domain; not distinguishing between attributes) (Dutch)), Nadig & sedivy (nadig2002 report that children redundantly use modifiers (mainly color) in a privileged-ground condition in 50.0% of cases; adults redundantly use modifiers (mainly color) in a privileged-ground condition in 46.7% of cases.)]

| Study                    | Language | Color sufficient |             |                   | Size sufficient |             |                   | comments  |
|--------------------------|----------|------------------|-------------|-------------------|-----------------|-------------|-------------------|---|
|                          |          | <i>color</i>     | <i>size</i> | <i>color_size</i> | <i>color</i>    | <i>size</i> | <i>color_size</i> |   |
| Pechmann (1989)          | Dutch    | 99               | 0           | 1                 | 9               | 36          | 55                |   |
| Gatt et al. (2011)       | English  | 92               | 0           | 8                 | 3               | 17          | 80                |   |
| Gatt et al. (2011)       | Dutch    | 90               | 0           | 10                | 0               | 21          | 79                |   |
| Gatt et al. (2013)       | Dutch    | 70               | (0.5)       | (30)              | (0.5)           | 69          | (31)              | values in parentheses: proportions "overspecified"/"underspecified"                             |
| Rubio-Fernandez (2016)   | English  | NA               | NA          | NA                | (37)            | NA          | NA                | (type sufficient) Paper doll task - Monochrome  |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (5)             | NA          | NA                | (type sufficient) Paper doll task - Monochrome  |
| Rubio-Fernandez (2016)   | English  | NA               | NA          | NA                | (95)            | NA          | NA                | (type sufficient) Paper doll task - Polychrome  |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (59)            | NA          | NA                | (type sufficient) Paper doll task - Polychrome  |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (0)             | NA          | NA                | (type sufficient) Yellow pig task - Standard instructions - stereotypical                       |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (6)             | NA          | NA                | (type sufficient) Yellow pig task - Standard instructions - variable                            |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (14)            | NA          | NA                | (type sufficient) Yellow pig task - Standard instructions - atypical                            |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (0)             | NA          | NA                | (type sufficient) Yellow pig task - Cautionary instructions - stereotypical                     |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (16)            | NA          | NA                | (type sufficient) Yellow pig task - Cautionary instructions - variable                          |
| Rubio-Fernandez (2016)   | Spanish  | NA               | NA          | NA                | (67)            | NA          | NA                | (type sufficient) Yellow pig task - Cautionary instructions - atypical                          |
| Westerbeek et al. (2015) | Dutch    | NA               | NA          | NA                | (46)            | NA          | NA                | (type sufficient) color typicality scores ranging from 2 to 98                                  |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (4)             | NA          | NA                | (type sufficient) low scene variation (distractors have different t-                            |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (24)            | NA          | NA                | (type sufficient) high scene variation (distractors have different t-                           |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (9)             | NA          | NA                | (type plus (size OR orientation) sufficient) low scene variation (distractors have different t- |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (18)            | NA          | NA                | (type plus (size OR orientation) sufficient) high scene variation (distractor                   |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (10)            | NA          | NA                | ((size OR orientation) sufficient) low scene variation (distractor                              |
| Koolen et al. (2013)     | Dutch    | NA               | NA          | NA                | (27)            | NA          | NA                | ((size OR orientation) sufficient) high scene variation (distractor                             |
| Our baseline study       | English  | 94               | 0           | 6                 | 2               | 52          | 46                |   |

and size, nor can it straightforwardly be extended to capture the more subtle systematicity with which the preference to overmodify with color changes based on various features of context. We delve into these more subtle patterns in the next two sections before presenting our alternative model within the Rational Speech Act framework.

### 1.2.2 Scene variation

So far we have portrayed speakers' propensity to overmodify with color as a fixed quantity (though varying by experiment). However, this propensity is highly dependent on features of the distractor objects in the context. In particular, as the variation present in the scene increases, so does the probability of overmodifying with color (Davies & Katsos, 2013; Koolen, Goudbeek, & Kraemer, 2013). How exactly scene variation is quantified differs across experiments. One very clear demonstration of the scene variation effect was given by Koolen et al. (2013), who quantified scene variation as the number of feature dimensions along which objects in a scene vary. Over the course of three experiments, they compared a low-variation condition in which objects never differed in color with a high-variation condition in which objects differed in type, color, orientation, and size. They consistently found higher rates of overmodification with color in the high-variation (28-27%) than in the low-variation (4-10%) conditions.

The effect of scene variation on propensity to overmodify has typically been explained as the result of the demands imposed on visual search: in low-variation scenes, it is easier to discern the discriminating dimensions than in high-variation scenes, where it may be easier to simply start naming features of the target that are salient (Koolen et al., 2013).

The PRO model does not have a straightforward way of capturing the effect of scene variation on probability of overmodification. One way of doing so is to make the salience and overmodification parameters directly dependent on the amount of variation in the scene. However, this requires additional free parameters and makes the model prone to overfitting. [jd: elaborate? throw out?]

We show in Section 2 how scene variation effects fall straightforwardly out of our proposed model by capturing the intuition that color becomes an increasingly good property for speakers to mention as it becomes contextually more informative.

### 1.2.3 Feature typicality

Modifier type and amount of scene variation are not the only factors determining overmodification. Overmodification with color has been shown to be systematically related to the typicality of the color for the object. Building on work by Sedivy (2003), Westerbeek et al. (2015) (and more recently, Rubio-Fernandez (2016)) have shown that the more typical a color is for an object, the less likely it is to be mentioned when not necessary for unique reference. For example, speakers never refer to a yellow banana in the absence of other bananas as *the yellow banana* (see Figure 2a), but they sometimes refer to a brown banana as *the brown banana*, and they almost always refer to a blue banana as *the blue banana* (see Figure 2b). Similar typicality effects have been shown for other (non-color) properties. For example, Mitchell (2013) showed that speakers are more likely to include an atypical than a typical property (either shape or material) when referring to everyday objects like boxes when mentioning at least one property was necessary for unique reference.

Whether speakers are more likely to mention atypical properties over typical properties because they are more salient to *them* or because they are trying to make reference resolution easier for the listener, for whom presumably these properties are also salient, is an open question (Westerbeek



(a) Typical color, type sufficient.

(b) Atypical color, type sufficient.

Figure 2: Example contexts where type (*banana*) is sufficient for unique reference and color is (a) typical or (b) atypical. A green border marks the intended referent.

et al., 2015). Some support for the audience design account comes from a study by Huettig and Altmann (2011), who found that listeners, after hearing a noun with a diagnostic color (e.g., *frog*), are more likely to fixate objects of that diagnostic color (green), indicating that typical object features are rapidly activated and aid visual search. Similarly, Arts et al. (2011) showed that overspecified expressions result in faster referent identification. Nevertheless, the benefit for listeners and the salience for speakers might simply be a happy coincidence and speakers might not, in fact, be designing their utterances for their addressees. We will remain agnostic about the underlying reason for typicality effects for the time being and will return to this issue in the General Discussion.

Irrespective of the source of typicality effects, it is unclear how the PRO model could accommodate them. In addition, one is left with the task of explaining how scene variation and typicality should interact. We show in Section 4 that the production model we propose straightforwardly accounts for these typicality effects in a principled way. In addition, we provide a principled way in which distractor typicality can be taken into account in modifier choice.

### 1.3 Nominal referring expressions

A problem related to the issue of how many additional features to include in a modified referring expression, but which has received much less attention in the language production literature, is that of deciding at which taxonomic level to refer to an object in a simple nominal expression. That is, even in the absence of adjectives, a referring expression can be more or less informative: *the dalmatian* communicates more information about the object in question than *the dog* (being a dalmatian entails being a dog), which in turn is globally more informative than *the animal*. Thus, this choice can be considered analogous to the choice of adding more modifiers – in both cases, the speaker has a choice of being more or less specific about the intended referent. However, the choice of reference level in simple nominal referring expressions is also interestingly different from that of adding modifiers in that there is no additional word-level cost associated with being more specific – the choice is between different one-word utterances, not between utterances that differ in word length.

Nevertheless, cognitive cost affects the choice of reference level: in particular, speakers prefer more frequent words over less frequent ones (Oldfield & Wingfield, 1965), and they prefer shorter ones over longer ones (Degen, Franke, & Jäger, 2013; Rohde, Seyfarth, Clark, Jäger, & Kaufmann, 2012). This may go part of the way towards explaining the well-documented effect from the concepts and categorization literature that speakers prefer to refer at the *basic level* (Rosch et



(a) Subordinate level term necessary.

(b) Superordinate level term sufficient.

Figure 3: Example contexts in which different levels of reference are necessary for establishing unique reference to the target marked with a green border. (a) subordinate (*dalmatian*) necessary; (b) superordinate (*animal*) sufficient, but basic (*dog*) or subordinate (*dalmatian*) possible.

al., 1976; Tanaka & Taylor, 1991b). That is, in the absence of other constraints, even when a superordinate level term would be sufficient for establishing reference (as in Figure 3b), speakers prefer to say *the dog* rather than *the animal*.

Contextual informativeness is another factor that has been shown to affect speakers' nominal production choices (e.g., Brennan & Clark, 1996). For instance, in a context like Figure 3a, speakers should use the subordinate level term *dalmatian* to refer to the target marked with a green border, because a higher-level term (*dog*, *animal*) would be contextually underinformative. However, there are nevertheless cases of contexts where either the superordinate *animal* or the basic level *dog* term would be sufficient for unique reference, as in Figure 3b, in which speakers nevertheless prefer to use the subordinate level term *the dalmatian*. This is the case when the object is a particularly good instance of the subordinate level term or a particularly bad instance of the basic level term, compared to the other objects in the context. For example, penguins, which are rated as particularly atypical birds, are often referred to at the subordinate level *penguin* rather than at the basic level *bird*, despite the general preference for the basic level (Jolicoeur, Gluck, & Kosslyn, 1984).

#### 1.4 Summary

In sum, the production of modified and simple nominal referring expressions is governed by a rich interplay of many factors, including an utterance's informativeness, its cost relative to alternative utterances, and the typicality of an object or its features. In the next section, we provide an explicit computational account of how these different factors interact, with a focus on cases where speakers appear to be overinformative – either by adding more modifiers or by referring at a more specific level than necessary for establishing unique reference. A summary of the effects we will focus on in the remainder of the paper is provided in Table 2.

To date, there is no theory to account for all of these different phenomena; and no model has attempted to unify overinformativeness in the domain of modified and nominal referring expressions. We touched on some of the explanations that have been proposed for these phenomena. We also highlighted where computational models have been proposed for individual phenomena, and how they fall short. In the next section, we present the Rational Speech Act modeling framework,

<sup>3</sup>Reported by many (e.g., Pechmann, 1989; Engelhardt et al., 2006a; Gatt et al., 2011; Rubio-Fernandez, 2016)

<sup>4</sup>Multiple replications reported (e.g., Davies & Katsos, 2013; Koolen et al., 2013)

<sup>5</sup>Multiple replications reported (e.g. Sedivy, 2003; Westerbeek et al., 2015; Rubio-Fernandez, 2016)

<sup>6</sup>Originally reported by Rosch et al. (1976), dozens of replications.

<sup>7</sup>Reported by Jolicoeur et al. (1984)

Table 2: List of effects a theory of referring expression production should account for and paper section(s) in which they are treated.

| Section | Effect                 | Description  |
|---------|------------------------|--|
| 2 & 3   | Color/size asymmetry   | More redundant use of color than size <sup>3</sup>                         |
| 2 & 3   | Scene variation        | More redundant use of color with increasing scene variation <sup>4</sup>   |
| 4       | Color typicality       | More redundant use of color with decreasing color typicality <sup>5</sup>  |
| 5       | Basic level preference | Preference for basic level term when superordinate sufficient <sup>6</sup> |
| 5       | Subordinate level use  | Unnecessary use of subordinate level term <sup>7</sup>                     |

within which we will provide precisely the kind of theory that can account for at least all of the phenomena listed here and holds great promise for scaling up to many other overinformativeness phenomena.

## 2 Modeling speakers' choice of referring expression

Here we propose a computational model of referring expression production that accounts for the phenomena introduced above. The model is formulated within the Rational Speech Act (RSA) framework (M. C. Frank & Goodman, 2012; Goodman & Frank, 2016).<sup>8</sup> It provides a principled explanation for the phenomena reviewed in the previous section and holds promise for being generalizable to many further production phenomena related to overinformativeness, which we discuss in Section 6. We proceed by first presenting the general framework in Section 2.1, and show why the most basic model, as formulated by M. C. Frank & Goodman, 2012, does not produce the phenomena outlined above due to its strong focus on speakers maximizing the informativeness of one-word expressions under a deterministic semantics. In Section 2.2 we introduce the crucial innovation: relaxing the assumption of a deterministic semantics. We show that the model can qualitatively account both for speakers' asymmetric propensity to overmodify with color rather than with size and (in Section 2.3) for speakers' propensity to overmodify more with increasing scene variation. In Section 3 we report an interactive reference game experiment which functions as a quantitative test of the model. In Section 4 we explore how the model captures color typicality effects. In Section 5 we apply the model to the choice of simple nominal referring expressions.

### 2.1 Basic RSA

As has been pointed out by Gatt et al. (2013), the basic Rational Speech Act model as formulated by M. C. Frank and Goodman (2012) does not generate overinformative referring expressions for two reasons: first, it trivially cannot do so because it is limited to one-word utterances (see also Baumann, Clark, & Kaufmann, 2014). But even when allowing two-word (or  $n$ -word) utterances, the speaker's utility function does not allow for producing redundant referring expressions as long as additional words contribute non-negative costs to the overall utterance cost. To see this, and as

---

<sup>8</sup>All RSA models and Bayesian Data Analyses reported in this paper were implemented in the probabilistic programming language WebPPL (Goodman & Stuhlmüller, electronic) and can be viewed at XXX.

Table 3: Row-wise literal listener distributions  $P_{L_0}(o|u)$  for each utterance  $u$  in the size-sufficient context depicted in Figure 1a, allowing only simple one-word utterances (left) or one- and two-word utterances (middle, right) under a deterministic semantics (left, middle) or under a continuous semantics (right) with  $\alpha = 1$ ,  $x_{\text{size}} = .8$ ,  $x_{\text{color}} = .99$ ,  $\beta_c = 0$ . Bolded numbers indicate crucial comparisons between literal listener probabilities in correctly selecting the intended referent  $o_{\text{small\_blue}}$  in response to observing the sufficient *small* and the redundant *small blue* utterances.

|                   | deterministic (simple) |                       |                          | deterministic (complex) |                       |                          | non-deterministic      |                       |                          |
|-------------------|------------------------|-----------------------|--------------------------|-------------------------|-----------------------|--------------------------|------------------------|-----------------------|--------------------------|
|                   | $o_{\text{big\_blue}}$ | $o_{\text{big\_red}}$ | $o_{\text{small\_blue}}$ | $o_{\text{big\_blue}}$  | $o_{\text{big\_red}}$ | $o_{\text{small\_blue}}$ | $o_{\text{big\_blue}}$ | $o_{\text{big\_red}}$ | $o_{\text{small\_blue}}$ |
| <i>big</i>        | .5                     | .5                    | 0                        | .5                      | .5                    | 0                        | .44                    | .44                   | .11                      |
| <i>small</i>      | 0                      | 0                     | 1                        | 0                       | 0                     | <b>1</b>                 | .17                    | .17                   | <b>.67</b>               |
| <i>blue</i>       | .5                     | 0                     | .5                       | .5                      | 0                     | .5                       | .50                    | .01                   | .50                      |
| <i>red</i>        | 0                      | 1                     | 0                        | 0                       | 1                     | 0                        | .01                    | .99                   | .01                      |
| <i>big blue</i>   | NA                     | NA                    | NA                       | 1                       | 0                     | 0                        | .79                    | .01                   | .20                      |
| <i>big red</i>    | NA                     | NA                    | NA                       | 0                       | 1                     | 0                        | .01                    | .99                   | .00                      |
| <i>small blue</i> | NA                     | NA                    | NA                       | 0                       | 0                     | <b>1</b>                 | .20                    | .00                   | <b>.80</b>               |

a basis for the innovation introduced in Section 2.2 it is useful to reiterate the basic form of the model.

The production component of RSA aims to soft-maximize the utility of utterances, where utility is defined in terms of the contextual informativeness of an utterance, given each utterance’s literal semantics. Formally, this is treated as a pragmatic speaker  $S_1$  reasoning about a literal listener  $L_0$ , who can be described by the following formula:

$$P_{L_0}(o|u) \propto \mathcal{L}(u, o). \quad (1)$$

The literal listener  $L_0$  observes an utterance  $u$  from the set of utterances  $U$ , consisting of single adjectives denoting features available in the context of a set of objects  $O$ , and returns a distribution over objects  $o \in O$ . Here,  $\mathcal{L}(u, o)$  is the lexicon that encodes deterministic lexical meanings such that:

$$\mathcal{L}(u, o) = \begin{cases} 1 & \text{if } u \text{ is true of } o \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Thus,  $P_{L_0}(o|u)$  returns a uniform distribution over all contextually available  $o$  in the extension of  $u$ . For example, in the size-sufficient context shown in Figure 1a,  $U = \{\text{big}, \text{small}, \text{blue}, \text{red}\}$  and  $O = \{o_{\text{big\_blue}}, o_{\text{big\_red}}, o_{\text{small\_blue}}\}$ . Upon observing *blue*, the literal listener therefore assigns equal probability to  $o_{\text{big\_blue}}$  and  $o_{\text{small\_blue}}$ . Values of  $P_{L_0}(o|u)$  for each  $u$  are shown on the left in Table 3.

The pragmatic speaker in turn produces an utterance with probability proportional to the utility of that utterance:

$$P_{S_1}(u|o) \propto e^{\alpha U(u, o)} \quad (3)$$

The speaker’s utility  $U(u, o)$  is a function of both the utterance’s *informativeness* with respect to the literal listener  $P_{L_0}(o|u)$  and the utterance’s *cost*  $c(u)$ :

$$U(u, o) = \ln P_{L_0}(o|u) - \beta_c c(u) \quad (4)$$

Two free parameters enter the computation: the speaker’s overall utility is weighted by parameter  $\alpha$  and utterance cost is weighted by parameter  $\beta_c$ .<sup>9</sup> In order to understand the effect of  $\alpha$ , it is useful to explore its effect when utterances are cost-free. In this case, as  $\alpha$  approaches infinity, the speaker increasingly only chooses utterances that maximize informativeness; if  $\alpha$  is 0, informativeness is disregarded and the speaker chooses randomly from the set of all available utterances; if  $\alpha$  is 1, the speaker probability-matches, i.e., chooses utterances proportional to their informativeness (equivalent to Luce’s choice rule, Luce, 1959). Applied to the example in Table 3, if the speaker wants to refer to  $o_{\text{small\_blue}}$  they have two semantically possible utterances, *small* and *blue*, where *small* is twice as informative as *blue*. They produce *small* with probability 1 when  $\alpha \rightarrow \infty$ , probability 2/3 when  $\alpha = 1$  and probability 1/4 when  $\alpha = 0$ .

Conversely, disregarding informativeness and focusing only on cost, any asymmetry in costs will be exaggerated with increasing  $\beta_c$ , such that the speaker will choose the least costly utterance with higher and higher probability as  $\beta_c$  increases.

As noted above, this model does not generate redundant referring expressions for multiple reasons. One of these is trivial:  $U$  only contains one-word utterances. We can ameliorate this easily by allowing complex two-word utterances. We assume an intersective semantics for complex utterances  $u_{\text{complex}}$  that consist of a two adjective sequence  $u_{\text{size}} \in \{\text{big}, \text{small}\}$  and  $u_{\text{color}} \in \{\text{blue}, \text{red}\}$ , such that the meaning of a complex two-word utterance is defined as

$$\mathcal{L}(u_{\text{complex}}, o) = \mathcal{L}(u_{\text{size}}, o) \times \mathcal{L}(u_{\text{color}}, o). \quad (5)$$

The resulting renormalized literal listener distributions for our example size-sufficient context in Figure 1a are shown in the middle columns in Table 3.

Unfortunately, simply including complex utterances in the set of alternatives does not solve the problem. Let’s turn again to the case where the speaker wants to communicate the small blue object. There are now two utterances, *small* and *small blue*, which are both more informative than *blue* and equally informative as each other, for referring to the small blue object. Because they are equally contextually informative, the only way for the complex utterance to be chosen with greater probability than the simple utterance is if it was the *cheaper* one. While this would achieve the desired mathematical effect, the cognitive plausibility of complex utterances being cheaper than simple utterances is highly dubious. Even if it wasn’t dubious: as mentioned previously, proportions of overinformative referring expressions are variable across experiments. The only way to achieve that variability under the basic model is to assume that the costs of utterances vary from task to task. This also seems to us an implausible assumption. Thus, unless we want to introduce dubious cost assumptions, we must look elsewhere to account for overinformativeness. We propose that the place to look is the computation of informativeness itself. This is what we turn to next.

---

<sup>9</sup>M. C. Frank and Goodman (2012) did not include cost in their formulation because they assumed equal costs for all utterances. Subsequent work has demonstrated the importance of taking into account utterance cost in modeling interpretation phenomena like cost-based quantity implicatures (Degen, Franke, & Jäger, 2013) and M-implicature (Bergen, Levy, & Goodman, 2016). We include it here because of the importance that cost has played in explanations of overinformative referring expressions, where it typically surfaces as the idea that speakers have different overall preferences for mentioning color vs. size modifiers (Dale & Reiter, 1995; Koolen et al., 2011; ?, ?). At this point we remain agnostic about the factors that contribute to an utterance’s cost  $c(u)$ . In later sections we treat cost as a function of both an utterance’s length and frequency.

## 2.2 RSA with continuous semantics – emergent color-size asymmetry

Here we introduce the crucial innovation: rather than assuming a deterministic truth-conditional semantics that returns true (1) or false (0) for any combination of expression and object, we assume a non-deterministic, continuous, semantics that returns real values in the interval [0, 1]. Formally, the only change is in the values that the lexicon returns:

$$\mathcal{L}(u, o) = x \in \mathbb{R} : x \in [0, 1] \quad (6)$$

That is, rather than assuming that an object is unambiguously big (or not) or unambiguously blue (or not), this continuous semantics captures that objects count as big or blue to varying degrees (similar to approaches in fuzzy logic and prototype theory, Zadeh, 1965; ?, ?). In principle, we could allow the semantic value of any utterance-object combination to vary. However, to see the basic effect of switching to a continuous semantics and to see how far we can get in capturing overinformativeness patterns with just this simple change, we will only distinguish between the semantic value of color and size adjectives in this most basic formulation of the continuous semantics model. When a size adjective is ‘true’ of an object under a deterministic semantics, we call that adjective’s semantic value as applied to the object  $x_{\text{size}}$ . When it is ‘false’ of the object, the semantic value is  $1 - x_{\text{size}}$ . Similarly for color adjectives. This results in two free model parameters,  $x_{\text{size}}$  and  $x_{\text{color}}$ , that can take on different values, capturing that size and color adjectives may on average apply more or less to objects.

To understand the motivation for this rather drastic move, consider some of the notable differences between color and size adjectives: color adjectives are typically treated as *absolute adjectives* while size adjectives are inherently *relative* (Kennedy & McNally, 2005). That is, while both size and color adjectives are vague, size adjectives are arguably context-dependent in a way that color adjectives are not – whether an object is big depends inherently on its comparison class; whether an object is red does not.<sup>10</sup> In addition, color as a property has been claimed to be inherently salient in a way that size is not (Arts et al., 2011; Gatt et al., 2013). Finally, we have shown in recent work that color adjectives are rated as less subjective than size adjectives (Scontras, Degen, & Goodman, 2017). We use these observations as motivation for exploring the effects of the assumption that the semantic value of size adjectives is inherently lower (i.e., that size adjectives are inherently noisier) than that of color adjectives.

The more extreme (closer to 0 and 1) an utterance type’s semantic value, the less uncertainty there is on the literal listener’s end about whether an object exhibits the property denoted by the observed expression; conversely, the less extreme (closer to 0.5) an utterance type’s semantic value, the more uncertainty there is in the literal listener. We defer a discussion of the important potential psychological and linguistic interpretation of these semantic values to the General Discussion in Section 6.

As an example, the resulting renormalized literal listener distributions for the size-sufficient example context in Figure 1a are shown for values  $x_{\text{size}} = .8$  and  $x_{\text{color}} = .99$  on the right in Table 3. Recall that in this context, the speaker intends for the listener to select the small blue pin. To see which would be the best utterance to produce for this purpose, we can compare the literal listener probabilities in the  $o_{\text{small}, \text{blue}}$  column. The two best utterances under both the deterministic

---

<sup>10</sup>This is not entirely true, as has been repeatedly pointed out (e.g., Cohen & Murphy, 1984): red hair has a very different color than red wine, which in turn has a different color from a red bell pepper. If presented out of context, only the last red is likely to be judged as red. For our purposes, it suffices that one can give a color judgment but not a size judgment for an object presented in isolation.

and the continuous semantics are bolded in the table: under the deterministic semantics, the two best utterances are *small* and *small blue*, with no difference in listener probability. In contrast, under the continuous semantics *small* has a smaller literal listener probability (.67) of retrieving the intended referent than the redundant *small blue*. Consequently, the pragmatic speaker will be more likely to produce *small blue* than *small*, though the precise probabilities depend on the cost parameter  $\beta_c$  and the rationality parameter  $\alpha$ .

Crucially, the reverse is not the case when color is the distinguishing dimension. Imagine the speaker in the same context wanted to communicate the big red pin. The two best utterances for this purpose are *red* (.99) and *big red* (.99). In contrast to the results for the small blue pin, these utterances do not differ in their capacity to direct the literal listener to the intended referent. The reason for this is that we defined color to be almost noiseless, with the result that the literal listener distributions in response to utterances containing color terms are more similar to those obtained via a deterministic semantics than the distributions obtained in response to utterances containing size terms. The reader is encouraged to verify this by comparing the row-wise distributions under the deterministic and continuous semantics in Table 3.

To gain a wider understanding of the effects of assuming non-deterministic meanings in contexts like that depicted in Figure 1a, we visualize the results of varying  $x_{\text{size}}$  and  $x_{\text{color}}$  in Figure 4. To orient the reader to the graph: the deterministic semantics of utterances is approximated where the semantic values of both size and color utterances are close to 1 (.999, top right-most point in graph). In this case, the simple sufficient (*small pin*) and complex redundant utterance (*small blue pin*) are equally likely around .5, because they are both equally informative and utterances are assumed to have 0 cost. All other utterances are highly unlikely. The interesting question is under which circumstances, if any, the standard color-size asymmetry emerges. This is the yellow/orange/red space in the ‘small blue’ facet, characterized by values of  $x_{\text{size}}$  that are lower than  $x_{\text{color}}$ , with high values for  $x_{\text{color}}$ . That is, redundant utterances are more likely than sufficient utterances when the redundant dimension (in this case color) is less noisy than the sufficient dimension (in this case size) and overall is close to noiseless.

Thus, when size adjectives are noisier than color adjectives, the model produces overinformative referring expressions with color, but not with size – precisely the pattern observed in the literature (Pechmann, 1989; Gatt et al., 2011). Note also that no difference in adjective *cost* is *necessary* for obtaining the overinformativeness asymmetry. However, assuming a greater cost for size than for color does further increase the observed asymmetry. We defer a discussion of costs to Section 3.1, where we infer the best parameter values for both the costs and the semantic values of size and color, given data from a reference game experiment.

To summarize, we have thus far shown that RSA with non-deterministic adjective semantics can give rise to the well-documented color-size asymmetry in the production of overinformative referring expressions when size adjectives are noisier than color adjectives. The crucial mechanism is this: when modifiers are noisy, adding additional, less noisy modifiers adds information. From this perspective, these redundant modifiers are not *overinformative*; they are rationally redundant, or sufficiently informative, given the needs of the listener. We spend the remainder of the paper demonstrating the far-reaching effects of assuming non-deterministic semantic values.

### 2.3 RSA with continuous semantics – scene variation

We begin by demonstrating the qualitative effect of a continuous semantics on redundant referring expressions in contexts that vary in the amount of visual complexity. As discussed in Section 1,

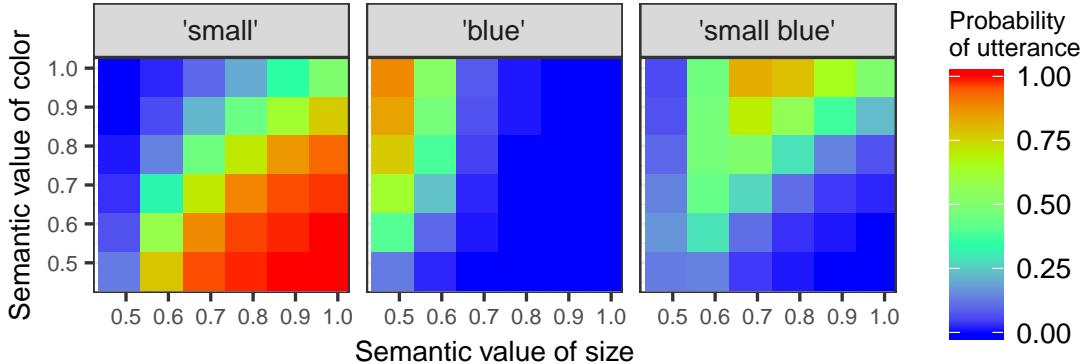


Figure 4: Probability of producing sufficient *small pin*, insufficient *blue pin*, and redundant *small blue pin* in contexts as depicted in Figure 1a, as a function of semantic value of color and size utterances (for  $\alpha = 30$  and  $\beta_c = 0$ ). For a visualization of model behavior under varying  $\alpha$ s, see Appendix A.

increased scene variation has been shown to increase the probability of referring expressions that are overmodified with color. Here we simulate the experimental conditions reported by Koolen et al. (2013) and explore continuous semantics RSA’s predictions for these situations. Koolen et al. (2013) quantified scene variation as the number of feature dimensions along which pieces of furniture in a scene varied: type (e.g., chair, fan), size (big, small), and color (e.g., red, blue).<sup>11</sup> Here, we simulate the high and low variation conditions from their Experiments 1 and 2, reproduced in Figure 5a.

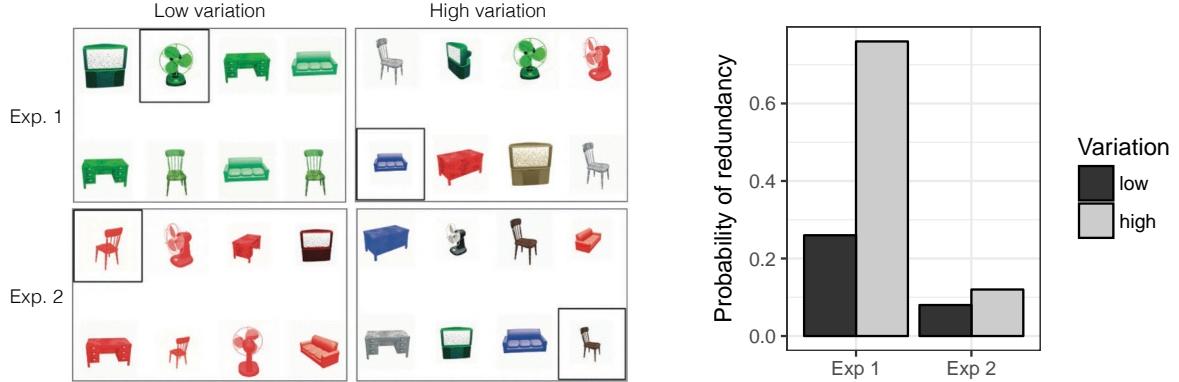
In both conditions in both experiments, color was not necessary for establishing reference; that is, color mentions were always redundant. The two experiments differed in the dimension necessary for unique reference. In Exp. 1, only type was necessary (*fan* and *couch* in the low and high variation conditions in Figure 5a, respectively). In Exp. 2, size and type were necessary (*big chair* and *small chair* in Figure 5a, respectively). Koolen et al. (2013) found lower rates of redundant color use in the low variation conditions (4% and 9%) than in the high variation conditions (24% and 18%).

We generated model predictions for precisely these four conditions. Note that by adding the type dimension as a distinguishing dimension, we must allow for an additional semantic value  $x_{\text{type}}$ , which encodes how noisy nouns are.

Koolen et al. (2013) counted any mention of color as a redundant mention. In Exp. 1, this includes the simple redundant utterances like *blue couch* as well as complex redundant utterances like *small blue couch*. In Exp. 2, where size was necessary for unique reference, only the complex redundant utterance *small brown chair* was truly redundant. The results of simulating these conditions for  $\alpha = 30$ ,  $\beta_c = c(u_{\text{size}}) = c(u_{\text{color}}) = 1$ ,  $x_{\text{size}} = .8$ ,  $x_{\text{color}} = .999$ , and  $x_{\text{type}} = .9$  are shown in Figure 5b.

For both experiments, the model retrieves the empirically observed effect of variation on the probability of redundant color mention: when variation is greater, redundant color mention is more likely. While the absolute values predicted by the model ( $\approx 8\%$  to  $\approx 75\%$ ) are different from

<sup>11</sup>They also included orientation (left-facing, right-facing) as a dimension along which objects could vary in certain cases. We ignore this dimension here for the sake of simplicity.



(a) Contexts from Koolen et al.’s low variation (left column) and high variation (right column) conditions in Koolen conditions for  $\alpha = 30$ ,  $\beta_c = c(u_{\text{size}}) = c(u_{\text{color}}) = 1$ ,  $x_{\text{size}} = .8$ ,  $x_{\text{color}} = .999$ ,  $x_{\text{type}} = .9$ .

Figure 5: Koolen et al. contexts and RSA model predictions.

the values observed by Koolen et al. (2013) ( $\approx 4\%$  to  $\approx 24\%$ ), it is note-worthy that continuous semantics RSA captures the qualitative scene variation effect with no further modifications.

Differences in exact values may stem from various sources. First, the best  $\alpha$  value to assume may differ from experiment to experiment. Second, semantic values may differ between experiments. Indeed, assuming a lower  $x_{\text{color}}$  of .9 maintains the qualitative effects but lowers the highest probability of redundancy to .26 (which is much closer to the 24% observed by Koolen et al.). Importantly, the basic requirements to yield the empirical scene variation effect are that semantic values for size, type, and color are ordered as follows:  $x_{\text{size}} \leq x_{\text{type}} < x_{\text{color}}$ . If  $x_{\text{type}}$  is greater than  $x_{\text{color}}$ , the probability of redundantly mentioning color is close to zero and does not differ between variation conditions. This is because in those cases, color mention reduces, rather than adds, information about the target. Third, the values reported by Koolen et al. (2013) were averaged over many different items – here, we only reported model predictions for the example items they reported.

These results are encouraging: RSA with a continuous semantics not only predicts a systematic color-size asymmetry in propensity to redundantly produce adjectives when size is noisier than color; it also predicts that there should be more redundant color mention as the number of dimensions along which objects in the scene vary increases. However, thus far we have only probed the model for qualitative effects from very few data points previously reported in the literature. Independently evaluating the utility of the model requires testing it on large datasets. This is what we turn to next. In Sections 3, 4, and 5 we quantitatively evaluate continuous semantics RSA on datasets capturing the phenomena described in the Introduction (for a summary see Table 2): modifier type and scene variation effects on modified referring expressions, typicality effects on color mention, and the choice of taxonomic level of reference in nominal choice, respectively.

### 3 Modified referring expressions: size and color modifiers under different scene variation conditions

Adequately assessing the explanatory value of RSA with non-deterministic truth functions requires evaluating how well it does at predicting the probability of various types of utterances occurring in large datasets of naturally produced referring expressions. To this end we proceed in two steps. First we report the results of a web-based interactive reference game in which we systematically manipulate scene variation (in a somewhat different way than Koolen et al. (2013) did). We then perform Bayesian data analysis to generate model predictions, conditioning on the observed production data. This allows us to both a) assess how likely the model is to generate the actually observed data – i.e., to obtain a measure of model quality – and b) infer the posterior probability of parameter values – i.e., to understand whether the assumed asymmetries in the adjectives' semantic values and/or cost discussed in the previous section are warranted.

#### 3.1 Experiment 1: scene variation in modified referring expressions

We saw in Section 2.3 that continuous semantics RSA correctly predicts qualitative effects of scene variation on redundant adjective use. In particular, we saw that color is more likely to be used redundantly as the number of dimensions along which objects in a scene vary increases. However, we would like to a) go beyond a qualitative investigation of scene variation effects and also b) ask whether redundant size mention is also affected by scene variation. The notion of scene variation we employ is the proportion of distractor items that do not share the value of the insufficient feature with the target, that is, as the number of distractors  $n_{\text{diff}}$  that differ in the value of the insufficient feature divided by the total number of distractors  $n_{\text{total}}$ :

$$\text{scene variation} = \frac{n_{\text{diff}}}{n_{\text{total}}}$$

To explain, let's turn again to Figure 1a. Here, the target item is the small blue pin and there are two distractor items: a big blue pin and a big red pin. Thus, for the purpose of establishing unique reference, size is the sufficient dimension and color the insufficient dimension. There is one distractor that differs from the target in color (the big red pin) and there are two distractors in total. That is,  $\text{scenevar} = \frac{1}{2} = .5$ . Scene variation is minimal when all distractors are of the same color as the target, in which case it is 0. Scene variation is maximal when all distractors except for one (in order for the dimension to remain insufficient for establishing reference) are of a different color than the target. That is, scene variation may take on values between 0 and  $\frac{n_{\text{total}}-1}{n_{\text{total}}}$ , i.e., approaching but never reaching 1.<sup>12</sup>

Using the same parameter values as in the previous two model explorations ( $\alpha = 30$ ,  $\beta_c = c(u_{\text{size}}) = c(u_{\text{color}}) = 1$ ,  $x_{\text{size}} = .8$ ,  $x_{\text{color}} = .999$ ), we generate model predictions for size-sufficient and color-sufficient contexts, varying scene variation by varying number of distractors (2, 3, or 4)

<sup>12</sup>Some readers might find this unintuitive: shouldn't scene variation be maximal when there is an equal number of same and different colors? Or when the different colors are also all different from one another? As discussed in the Introduction, there are many ways of quantifying (different aspects of) scene variation. Here we explore just one such measure; it is an interesting question whether RSA accounts equally well for different ways of quantifying scene variation. Fortunately, it is very straightforward to implement such different measures by manipulating features of distractor items and exploring the model's behavior in these contexts.

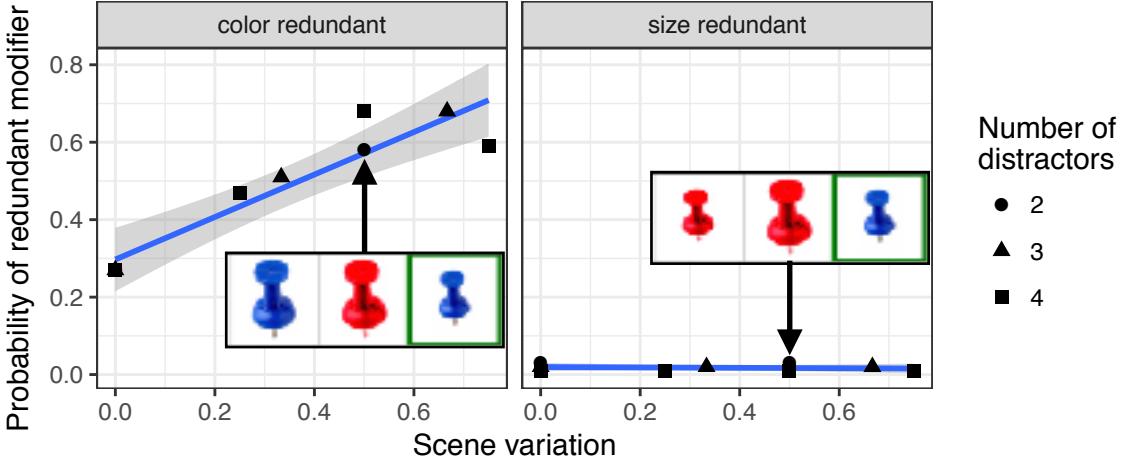


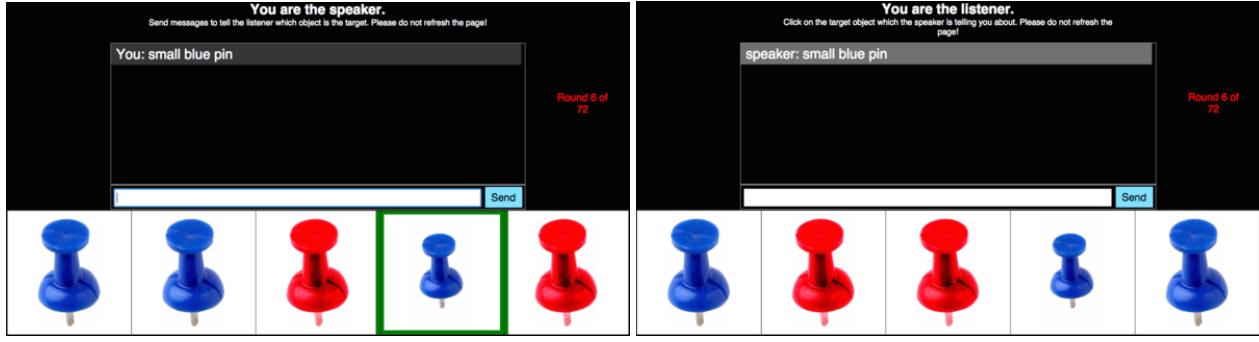
Figure 6: Probability of redundant utterance (*small blue pin*) as a function of scene variation when size is sufficient (and color redundant, left) and when color is sufficient (and size redundant, right), for  $\alpha = 30$ ,  $\beta_c = c(u_{\text{size}}) = c(u_{\text{color}}) = 1$ ,  $x_{\text{size}} = .8$ ,  $x_{\text{color}} = .999$ . Linear smoothers overlaid.

and number of distractors that don't share the insufficient feature value. The resulting model predictions are shown in Figure 6: the probability of redundant adjective use increases with increasing scene variation when size is sufficient (and color redundant), but not when color is sufficient (and size redundant). This can be explained by noise distributions in the literal listener across contexts: in size-sufficient contexts, as the number of distractors of a different color than the target increases, using the relatively noiseless color term in addition to the more noisy size term reduces uncertainty about the target object more and more. However, the same is not true of the color-sufficient contexts: there is very little uncertainty about the target upon observing the minimal color utterance – adding the size term only introduces more uncertainty about the target, regardless of the amount of scene variation. Note that this is highly dependent on the actual semantic value of color, with slightly lower semantic values for color, the model predicts small increases in redundant size use. This will be important for the interpretation of the empirical results. In general: increased scene variation is predicted to lead to a greater increase in redundant adjective use for less noisy adjectives.

To test continuous semantics RSA predictions, we conducted an interactive web-based written production study within a reference game setting. Speakers and listeners were shown arrays of objects that varied in color and size. Speakers were asked to produce a referring expression to allow the listener to identify a target object. We manipulated the number of distractor objects in the grid, as well as the variation in color and size among distractor objects.

### 3.1.1 Method

**Participants** We recruited 58 pairs of participants (116 participants total) over Amazon's Mechanical Turk who were each paid \$1.75 for their participation. Data from another 7 pairs who prematurely dropped out of the experiment and who could therefore not be compensated for their work, were also included. Here and in all other experiments reported in this paper, participants' IP address was limited to US addresses and only participants with a past work approval rate of at



(a) Speaker’s perspective.

(b) Listener’s perspective.

Figure 7: Example displays from the (a) speaker’s and the (b) listener’s perspective on a *size-sufficient 4-2* trial.

least 95% were accepted.

**Procedure** Participants were paired up through a real-time multi-player interface (Hawkins, 2015). For each pair, one participant was assigned the speaker role and one the listener role. They initially received written instructions that informed participants that one of them would be the Speaker and the other the Listener. They were further told that they would see some number of objects on each round and that the speaker’s task is to communicate one of those objects, marked by a green border, to the listener. They were explicitly told that using locative modifiers (like *left* or *right*) would be useless because the order of objects on their partner’s screen would be different than on their own screen. Before continuing to the experiment, participants were required to correctly answer a series of questions about the experimental procedure. These questions are listed in Appendix C.

On each trial participants saw an array of objects. The array contained the same objects for both speaker and listener, but the order of objects was randomized and was typically different for speaker and listener. In the speaker’s display, one of the objects – henceforth the *target* – was highlighted with a green border. See Figure 7 for an example of the listener’s and speaker’s view on a particular trial.

The speaker produced a referring expression to communicate the target to the listener by typing into an unrestricted chat window. After pressing Enter or clicking the ‘Send’ button, the speaker’s message was shown to the listener. The listener then clicked on the object they thought was the target, given the speaker’s message. Once the listener clicked on an object, a red border appeared around that object in both the listener and the speaker’s display for 1 second before advancing to the next trial. That is, both participants received feedback about the speaker’s intended referent and the listener’s inference.

Both speakers and listeners could write in the chat window, allowing listeners to request clarification if necessary. Listeners could only click on an object and advance to the next trial once the speaker sent a message.

**Materials** Participants proceeded through 72 trials. Of these, half were critical trials of interest and half were filler trials. On critical trials, we varied the feature that was sufficient to mention for

uniquely establishing reference, the total number of objects in the array, and the number of objects that shared the insufficient feature with the target.

Objects varied in color and size. On 18 trials, color was sufficient for establishing reference. On the other 18 trials, size was sufficient. Figure 7 shows an example of a size-sufficient trial. We further varied the amount of variation in the scene by varying the number of distractor objects in each array (2, 3, or 4) and the number of distractors that did share the redundant feature value with the target. That is, when size was sufficient, we varied the number of distractors that shared the same color as the target. This number had to be at least one, since otherwise the redundant property would have been sufficient for uniquely establishing reference, i.e. mentioning it would not have been redundant. Each total number of distractors was crossed with each possible number of distractors that shared the redundant property, leading to the following nine conditions: 2-1, 2-2, 3-1, 3-2, 3-3, 4-1, 4-2, 4-3, and 4-4, where the first number indicates the total number and the second number the shared number of distractors. Each condition occurred twice with each sufficient dimension. Objects never differed in type within one array (e.g., all objects are pins in Figure 7 but always differed in type across trials. Each object type could occur in two different sizes and two different colors. We deliberately chose photo-realistic objects of intuitively fairly typical colors. The 36 different object types and the colors they could occur with are listed in Appendix D.

Fillers were target trials from Exp. 2, a replication of Graf, Degen, Hawkins, and Goodman (2016). Each filler item contained a three-object grid. None of the filler objects occurred on target trials. Objects stood in various taxonomic relations to each other and required neither size nor color mention for unique reference. See Section 5.1 for a description of these materials.

### 3.1.2 Data pre-processing and exclusion

We collected data from 2171 critical trials. Because we did not restrict participants' utterances in any way, they produced many different kinds of referring expressions. Testing the model's predictions required, for each trial, classifying the produced utterance as an instance of a *color-only* mention, a *size-only* mention, or a *color-and-size* mention (or excluding the trial if no classification was possible). To this end we conducted the following semi-automatic data pre-processing.

First, 33 trials on which the listener selected the wrong referent were excluded, leading to the elimination of 1.5% of trials. Then, an R script automatically checked whether the speaker's utterance contained a precoded color (i.e. *black*, *blue*, *brown*, *gold*, *green*, *orange*, *pink*, *purple*, *red*, *silver*, *violet*, *white*, *yellow*) or size (i.e. *big*, *bigger*, *biggest*, *huge*, *large*, *larger*, *largest*, *little*, *small*, *smaller*, *smallest*, *tiny*) term. In this way, 95.7 % of cases were classified as mentioning size and/or color. However, this did not capture that sometimes, participants produced meaning-equivalent modifications of color/size terms for instance by adding suffixes (e.g., *bluish*), using abbreviations (e.g., *lg* for *large* or *purp* for *purple*), or using non-precoded color labels (e.g., *lime* or *lavender*). Expressions containing a typo (e.g., *pruple* instead of *purple*) could also not be classified automatically. In the next step, one of the authors (CG) therefore manually checked the automatic coding to include these kinds of modifications in the analysis. This covered another 1.5% of trials. Most of the time, participants converged on a convention of producing only the target's size and/or color, e.g., *purple* or *big blue*, but not an article (e.g., *the*) or the noun corresponding to the object's type (e.g., *comb*). Articles were omitted in 93.1 % of cases and nouns were omitted in 71.5 % of cases. We did not analyze this any further.

There were 50 cases (2.3%) in which the speaker made reference to the distinguishing dimension in an abstract way, e.g. *different color*, *unique one*, *ripest*, *very girly*, or *guitar closest to viewer*.

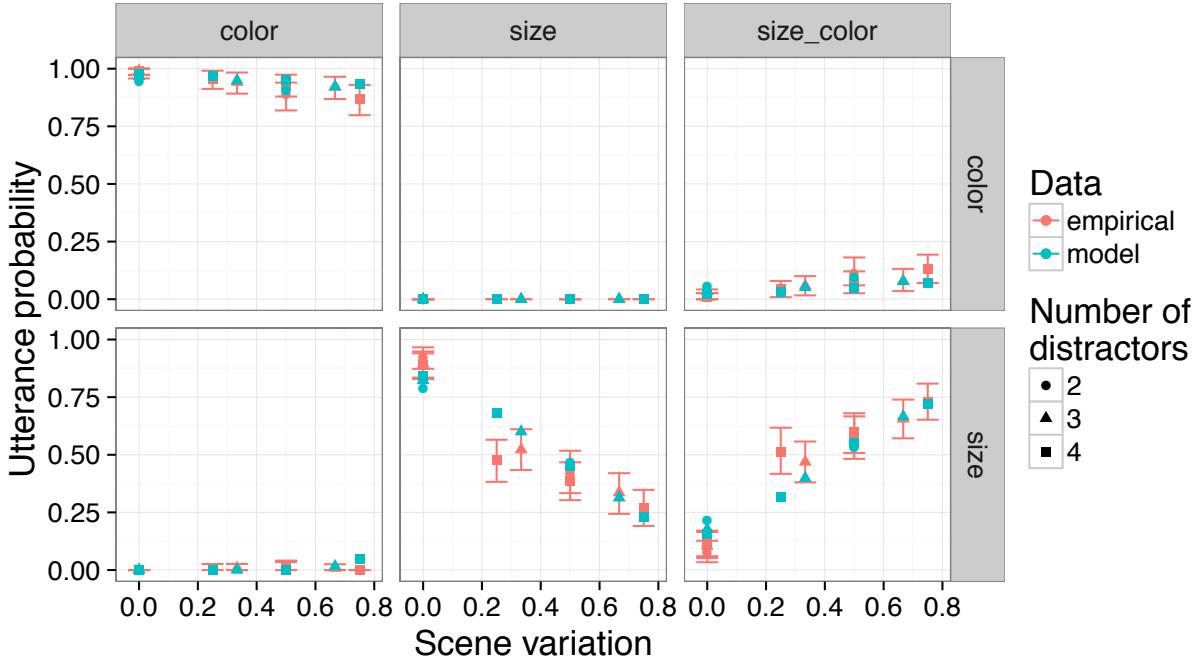


Figure 8: Empirical utterance proportions (red) alongside point-wise maximum a posteriori (MAP) estimates of the RSA model’s posterior predictives for utterance probability (blue) as a function of scene variation. Rows indicate the sufficient dimension, columns the produced utterance. Here and in all following plots, error bars indicate 95% bootstrapped confidence intervals.

While interesting as utterance choices,<sup>13</sup> these cases were excluded from the analysis. There were 3 cases that were nonsensical, e.g. *bigger off a shade*, which were also excluded. Finally, there were 6 cases where only the insufficient dimension was mentioned – these were excluded from the analysis reported in the next section, where we are only interested in minimal or redundant utterances, not underinformative ones, but were included in the Bayesian data analysis reported in Section 3.2. After the exclusion, 2079 cases classified as one of *color*, *size*, or *color-and-size* entered the analysis.

### 3.1.3 Results

Proportions of redundant *color-and-size* and minimal *color* or *size* utterances are shown in Figure 8 alongside model predictions (to be explained further in Section 3.2). There are three main questions of interest: first, do we replicate the color/size asymmetry in probability of redundant adjective use? Second, do we replicate the previously established effect of increased redundant color use with increasing scene variation? Third, is there an effect of scene variation on redundant size use and if so, is it smaller compared to that on color use, as is predicted under asymmetric semantic values for color and size adjectives?

We addressed all of these questions in one fell swoop by conducting a mixed effects logistic

<sup>13</sup>Certain participants seemed to have deliberately used this as a strategy even though simply mentioning the distinguishing property would have been shorter in most cases. In all, only 12 participants produced these kinds of utterances: one 18 times, one 8 times, one 6 times, two 3 times, one 2 times, and the remaining six only once each.

regression analysis predicting redundant over minimal adjective use from fixed effects of sufficient property (color vs. size), scene variation (proportion of distractors that does not share the insufficient property value with the target), and the interaction between the two.<sup>14</sup> The model included the maximal random effects structure that allowed the model to converge: by-speaker and by-item random intercepts as well as by-speaker random slopes for scene variation.

We observed a main effect of sufficient property, such that speakers were more likely to redundantly use color than size adjectives ( $\beta = 3.61$ ,  $SE = .23$ ,  $p < .0001$ ), replicating the much-documented color-size asymmetry. We further observed a main effect of scene variation, such that redundant adjective use increased with increasing scene variation ( $\beta = 4.11$ ,  $SE = .49$ ,  $p < .0001$ ). Finally, we also observed a significant interaction between sufficient property and scene variation ( $\beta = 3.03$ ,  $SE = .81$ ,  $p < .0002$ ). Simple effects analysis revealed that the interaction was driven by the scene variation effect being much smaller in the *color-sufficient* condition ( $\beta = 2.59$ ,  $SE = .78$ ,  $p < .0009$ ) than in the *size-sufficient* condition ( $\beta = 5.63$ ,  $SE = .45$ ,  $p < .0001$ ), as predicted if size modifiers are noisier than color modifiers.

### 3.2 Model evaluation

In order to evaluate RSA with non-deterministic truth functions, we asked how well it captures the empirical data. To this end we conducted a Bayesian data analysis. This allowed us to simultaneously generate model predictions and infer likely parameter values, by conditioning on the observed production data (coded into *size*, *color*, and *size-and-color* utterances as described above) and integrating over the following free parameters: semantic value for color  $x_{\text{color}}$ , semantic value for size  $x_{\text{size}}$ , color cost  $c(u_{\text{color}})$ , size cost  $c(u_{\text{size}})$ , cost weight  $\beta_c$ , and speaker rationality parameter  $\alpha$ . We assumed uniform priors for each parameter:  $x_{\text{color}} \sim \mathcal{U}(0, 1)$ ,  $x_{\text{size}} \sim \mathcal{U}(0, 1)$ ,  $c(u_{\text{color}}) \sim \mathcal{U}(0, 2)$ ,  $c(u_{\text{size}}) \sim \mathcal{U}(0, 2)$ ,  $\beta_c \sim \mathcal{U}(0, 10)$ ,  $\alpha \sim \mathcal{U}(0, 40)$ . Inference for the cognitive model was exact. We used Markov Chain Monte Carlo (MCMC) to infer posteriors for the six free parameters.

Point-wise maximum a posteriori (MAP) estimates of the model’s posterior predictives for each combination of utterance, sufficient dimension, number of distractors, and number of different distractors (collapsing across different items) are compared to empirical data in Figure 9. At this level, the model achieves a correlation of  $r = .99$ . Looking at results additionally on the by-item level yields a correlation of  $r = .85$ . The model thus does a very good job of capturing the quantitative patterns in the data. This can also be seen in Figure 8, where model predictions are plotted alongside the empirical proportions by condition. The only clear flaw is that the model predicts greater redundant adjective use than empirically observed when there is no scene variation at all. [jd: Though this is a pretty minor thing. Noah, can you add a sentence on why, given that you’ve thought about this for the negation case as well? Or should we just leave it out?].

Posteriors over parameters are shown in Figure 10. Crucially, the semantic value of color is inferred to be higher than that of size – there is no overlap between the 95% highest density intervals (HDIs) for the two parameters. That is, size modifiers are inferred to be noisier than color modifiers. The relatively high inferred  $\alpha$  suggests that this difference in semantic value contributes substantially to the observed color-size asymmetries in redundant adjective use. As for cost, there is a lot of overlap in the inferred cost of size and color modifiers and a very low weight on cost,

---

<sup>14</sup>All mixed effects analyses reported in this paper were conducted with the `lme4` package (Bates, Mächler, Bolker, & Walker, 2015) in R (R Core Team, 2017).

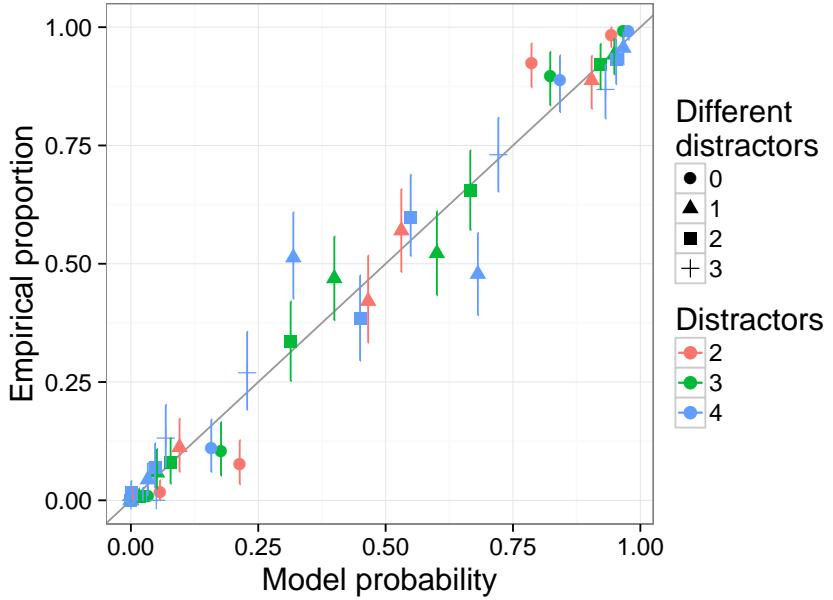


Figure 9: Scatterplot of point-wise maximum a posteriori (MAP) estimates of the RSA model’s posterior predictives against empirical proportions ( $r = .85$ ). [jd: update this plot once bda run]

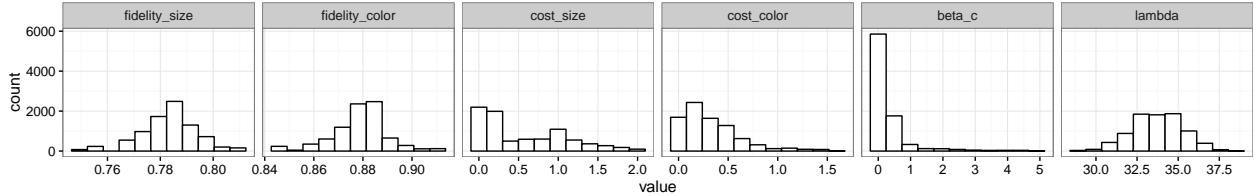


Figure 10: Posterior distribution over model parameters. Maximum a posteriori (MAP)  $f_s = 0.79$ , 95% highest density interval (HDI) = [0.76,0.80]; MAP  $f_c = 0.88$ , HDI = [0.86,0.91]; MAP  $c_{size} = .08$ , HDI = [0, 1.5]; MAP  $c_{color} = 0.07$ , HDI = [0,0.9]; MAP  $\beta_c = 0.04$ , HDI = [0,1.6]; MAP  $\alpha = 34.0$ , HDI = [30.8,36.5] [jd: update this plot once new bda run]

suggesting that no or very minimal cost difference is necessary to obtain the color-size asymmetry and the scene variation effects. While these results are compatible with part of the explanation for the color-size asymmetry stemming from the low cognitive cost involved in producing color modifiers compared to size modifiers, they also suggest that a cost asymmetry is not the driving force behind the asymmetry in redundant adjective use. Note further that the asymmetry cannot be reduced to cost differences: in Section 2.2 we showed that the color-size asymmetry in redundant adjective use requires an asymmetry in semantic value. An asymmetry in cost only serves to further enhance the asymmetry brought about by the asymmetry in semantic value, but cannot carry the redundant use asymmetry on its own.

### 3.3 Discussion

In this section, we reported the results of a dataset of freely collected referring expressions that replicated the well-documented color-size asymmetry in redundant adjective, the effect of scene variation on redundant color use, and showed a novel effect of scene variation on redundant size use. We also showed that continuous semantics RSA provides an excellent fit to these data. In particular, the crucial element in obtaining the color-size asymmetry in overmodification is that size adjectives be noisier than color adjectives, captured in RSA via a lower semantic value for size compared to color. The effect is that color adjectives are more informative than size adjectives when controlling for the number of distractors each would rule out under a deterministic semantics. Asymmetries in the cost of the adjectives only serve to further enhance the modification asymmetry resulting from the asymmetry in semantic value. In addition, we showed that asymmetric effects of scene variation on overmodification straightforwardly fall out of continuous semantics RSA: scene variation leads to a greater increase in overmodification with less noisy than with more noisy modifiers because the less noisy modifiers (colors) on average provide more information about the target.

Some readers may find themselves wondering about the status of these semantic values: are we claiming that color modifiers have inherently higher semantic values than size modifiers? Is the difference constant? What if the color modifier is a less well known one like *mauve*? The way we have set the model up thus far, there would indeed be no difference in semantic value between *red* and *mauve*. Moreover, the model is not equipped to handle potential object-level idiosyncracies such as the typicality effects discussed in Section 1.2.3. We defer a fuller discussion of the status of the semantic value term to the General Discussion and turn first to continuous semantics RSA’s potential for capturing these typicality effects.

## 4 Modified referring expressions: color typicality

In Section 3 we showed that continuous semantics RSA successfully captures both the basic asymmetry in overmodification with color vs. size as well as effects of scene variation, quantified in various different ways. But in Section 1.2.3 we discussed a further characteristic of speakers’ overmodification behavior: speakers are more likely to redundantly produce modifiers that denote atypical rather than typical object features, i.e., they are more likely to refer to a blue banana as a *blue banana* rather than as a *banana*, and they are more likely to refer to a yellow banana as a *banana* than as a *yellow banana* (Sedivy, 2003; Westerbeek et al., 2015). Continuous semantics RSA as we have set it up thus far does not capture this asymmetry: it knows that a particular modifier is a color modifier with a particular semantic value; it does not know anything about the typicality of the denoted properties for the referent.

We would like to warn and disillusion the reader upfront: we will not solve the problem of how to get overmodification behavior from the typicality of features compositionally. This is a problem for all theories of modification (Kamp & Partee, 1995). However, we would like to offer a proof of concept showing that, if the non-determinism in the RSA semantics is not at the adjective type (color, size) level, but instead at the level of combinations of referring expressions and objects, the model produces precisely the sorts of typicality effects reported in the literature.

Let us elaborate. Where before we took a semantic value to be a number between 0 and 1 indicating how likely a type of modifier (size, color) was to correctly apply to an object, we now treat it as indicating how good an instance of a particular referring expression the object in question



(a) Typical color.

(b) Mid-typical color.

(c) Atypical color.

Figure 11: Three hypothetical contexts where color is redundant for referring to the target banana. Banana varies in typicality from left to right. Each context contains one distractor of the same color as the target, and one of a different color.

Table 4: Hypothetical semantic values for utterances (rows) as applied to objects (columns).

|                      | yellow banana | brown banana | blue banana | other |
|----------------------|---------------|--------------|-------------|-------|
| <i>banana</i>        | .9            | .35          | .1          | .015  |
| <i>yellow banana</i> | .99           | .015         | .015        | .015  |
| <i>brown banana</i>  | .015          | .99          | .015        | .015  |
| <i>blue banana</i>   | .015          | .015         | .99         | .015  |
| other                | .015          | .015         | .015        | .99   |

is. For example, take the banana case: assume the three contexts in Figure 11. The target object in each is the banana, which varies in how typical its color is. The banana is the only object of its type, making type mention (*banana*) sufficient for unique reference and color redundant. Additionally, each context contains a distractor of the same color as the target, further making color redundant, as well as a distractor of a different color. Assume further the hypothetical semantic values shown in Table 4. These values should be read as follows: a yellow banana is a very good or typical instance of a *banana – banana* applied to yellow bananas has a high semantic value of .9. In contrast, brown bananas are less typical instances of *bananas* (.35), and blue bananas are highly atypical *bananas* (.1) but still better than objects of an other non-banana type (.015). Going along the diagonal, we assume for each remaining utterance that its semantic value is very high (.99) when applied to an object in its (deterministic truth-conditional) extension and very low otherwise (.015).

Inputting these contexts and semantic values into the RSA model used thus far, with  $\alpha = 12$  and  $\beta_c = 5$  (that is, both informativeness and utterance cost receive a substantial weight), the resulting speaker probabilities for the (minimal) *banana* are .99, .37, and .05, to refer to the yellow banana, the brown banana, and the blue banana, respectively. In contrast, the resulting speaker probabilities for the redundant *yellow banana*, *brown banana*, and *blue banana* are .01, .63, and .95, respectively. That is, redundant color mention increases with decreasing semantic value of the simple *banana* utterance.

So far we have shown that continuous semantics RSA can capture typicality effects in principle if we assume that semantic values do not operate at the adjective type level but instead captures the typicality of an object for the (minimal and redundant) referring expressions. If an object is more typical for the redundant expression than for the minimal expression, then the bigger the difference in typicality, the greater the relative informativeness of the redundant expression, and the greater the probability of it being produced.

This example is somewhat oversimplified. In practice, speakers sometimes just mention an

object’s color, without mentioning the noun. In the contexts presented in Figure 11 this does not make much sense because there is always a competitor of the same color present. In contrast, in the contexts in Figure 13a and Figure 13c, color alone disambiguates the target. This suggests that we should consider among the set of utterance alternatives not just the simple type mentions (e.g., *banana*) and color-and-type mentions (e.g., *yellow banana*), but also simple color mentions (e.g., *yellow*). The dynamics of the model proceed as before.

We can now ask whether taking into account this more fine-grained notion of a continuous semantics affects the probability of redundantly mentioning color. Because the stimuli for Exp. 1 were specifically designed to be realistic objects with low color-diagnosticity, they do not include objects with low typicality values or large degrees of variation in typicality. This makes the dataset from Exp. 1 not well-suited for investigating typicality effects.<sup>15</sup> We therefore conducted a separate production experiment in the same paradigm but with two broad changes: first, objects’ color varied in typicality; and second, we did not manipulate object size, focusing only on color mention. This allows us to ask two questions: first, do we replicate the typicality effects reported in the literature – that is, are less color-typical objects more likely to lead to redundant color use than more color-typical objects? Second, does RSA with empirically elicited typicality values as proxy for a continuous semantics capture speakers’ behavior better than just typicality alone? [jd: do we actually evaluate the latter?]

Additionally, the proposed RSA model allows us to ask questions that have not been addressed previously. In particular, it allows us to test for complex interactions of contextual informativeness of color and type, by manipulating both whether there is a distractor of the same type present and whether there is a distractor of the same color present.

## 4.1 Experiment 2: color typicality effects

### 4.1.1 Method

**Participants** We recruited 61 pairs of participants (122 participants total) over Amazon’s Mechanical Turk who were each paid \$1.80 for their participation. Two participant-pairs were excluded because they did not finish the experiment and therefore could not receive payment. Trials in which the speaker did not produce any utterances were excluded as well, which resulted in the exclusion of two additional participant-pairs. Finally, there were 10 speakers who consistently used round-about descriptions instead of direct referring expressions (e.g., *monkeys love...* to refer to banana). These pairs were also excluded because they were playing a game more akin to Taboo.

**Procedure** The procedure was identical to that of Exp. 1. See Figure 12 for an example speaker and listener perspective.

**Materials** Each participant completed 42 trials. In this experiment, there were no filler trials, since pilot studies with and without fillers delivered very similar results. Each array presented to the participants consisted of three objects that could differ in type and color. One of the three objects functioned as a target and the other two as its distractors.

The stimuli were selected from seven color-diagnostic food items (apple, avocado, banana, carrot, pear, pepper, tomato), which all occurred in a typical, mid-typical and atypical color for

---

<sup>15</sup>We did elicit typicality norms for the items in Exp. 1 and replicated the previously documented typicality effects on the four items that did exhibit variation in typicality. See Appendix E for details.

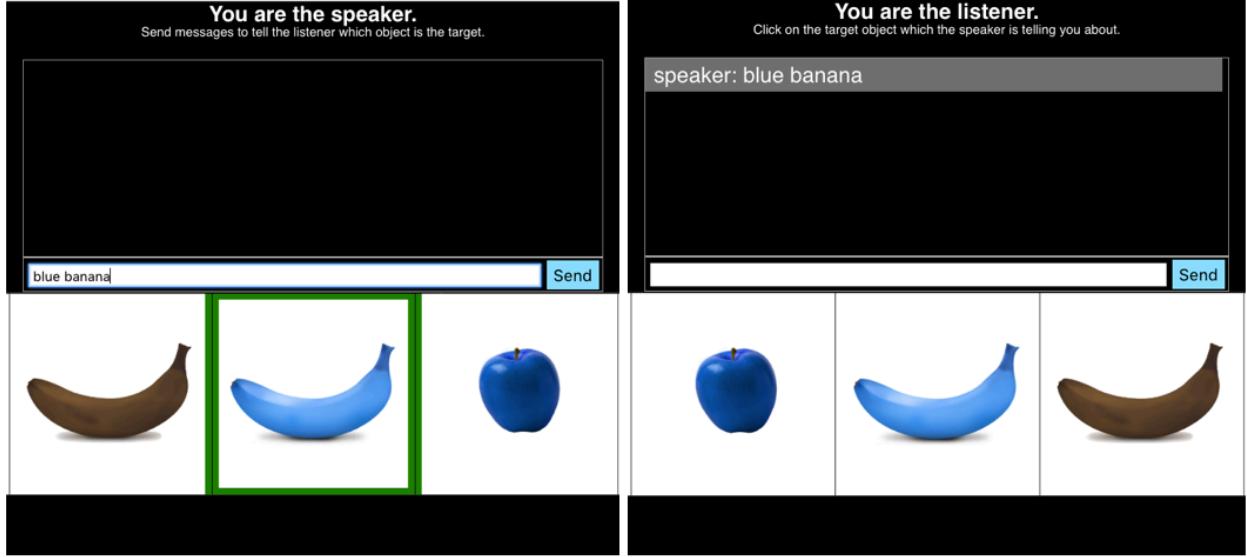


Figure 12: Example displays from the speaker’s (left) and listener’s (right) perspective in an informative-cc (i.e., presence of another object of the same type and one with the same color) condition.

that object. For example, the banana appeared in the colors yellow (typical), brown (midtypical), and blue (atypical). All items were presented as targets and as distractors. Pepper additionally occurred in a fourth color, which only functioned as a distractor due to the need for a green color competitor (as explained in the following paragraph).

We refer to the different context conditions as “informative”, “informative-cc”, “overinformative”, and “overinformative-cc” (see Figure 13). A context was “overinformative” (Figure 13c) when mentioning the type of the item, e.g., banana, was sufficient for unambiguously identifying the target. In this condition, the target never had a color competitor. This means that mentioning color alone (without a noun) was also unambiguously identifying. In contrast, in the overinformative condition with a color competitor (“overinformative-cc”, Figure 13d), color alone was not sufficient. In the informative conditions, color and type mention were necessary for unambiguous reference. Again, one context type did (Figure 13a) and one did not (Figure 13d) include a color competitor among its distractors.

In the end, each participant saw 42 different contexts. Each of the 21 items (color-type combinations) was the target exactly twice, but the context in which they occurred was drawn randomly from the four possible conditions mentioned above. In total, there were 84 different possible configurations (seven target food items, each of them in three colors, where each could occur in four contexts). Trial order was randomized.

#### 4.1.2 Data pre-processing and exclusion

We analyzed data from 1974 trials. Just as in Exp. 1, participants communicated freely, which led to a vast amount of different referring expressions. To test the model’s predictions, the utterance produced for each trial was classified as belonging to one of the following categories: *type-only* (“banana”), *color-and-type* (“yellow banana”), and *color-only* (“yellow”) utterances. Refer-

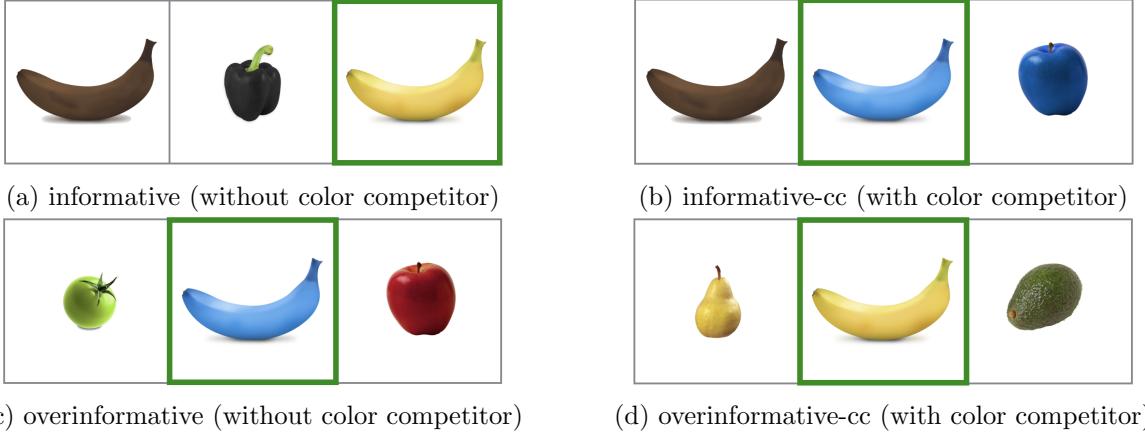


Figure 13: The four different context conditions in Exp. 2. They differed in the presence of an object of the same type (informative vs. overinformative) and in the presence of another object of the same color as the target (with color competitor vs. without color competitor). The green border marks the intended referent.

ring expressions that included categories (“yellow fruit”), descriptions (“has green stem”), color-circumscriptions (“funky carrot”), and negations (“yellow but not banana”) were regarded as *other* and excluded. To this end we conducted the following semi-automatic data pre-processing.

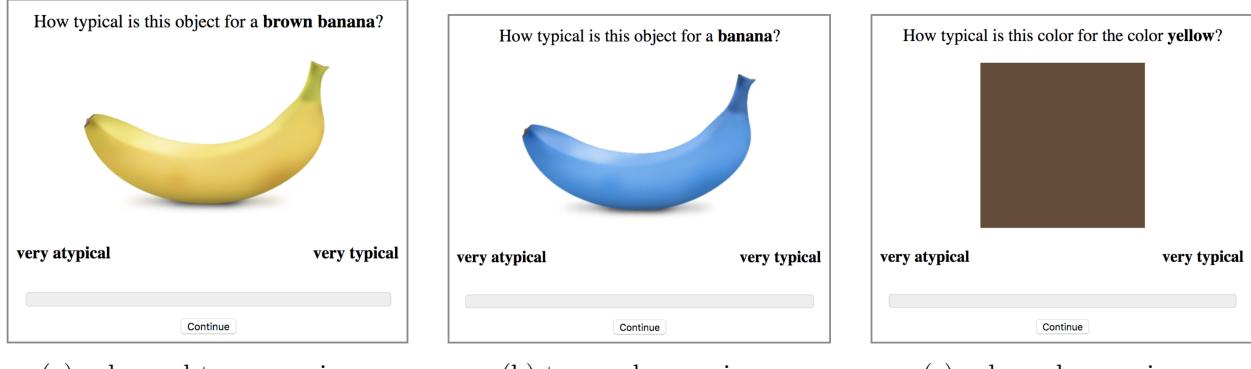
The referring expressions were analyzed similarly to Exp. 1. First, 32 trials (1.6%) were excluded because the listener selected the wrong referent. 109 trials (5.6%) were excluded because the referring expressions included one of the exceptional cases described above (e.g., using negations). An R script then automatically checked the remaining 1833 utterances for whether they contained a pre-coded color term (i.e. *green*, *purple*, *white*, *black*, *brown*, *yellow*, *orange*, *blue*, *pink*, *red*, *grey*) or type (i.e. *apple*, *banana*, *carrot*, *tomato*, *pear*, *pepper*, *avocado*). This way, 96.5% of the remaining cases were classified as mentioning type and/or color.

However, this did not capture that sometimes, participants produced meaning-equivalent modifications of color/type terms for instance by adding suffixes (e.g., *pinkish*), using abbreviations (e.g., *yel* for *yellow*), or using non-pre-coded color and type labels (e.g., *lavender* or *jalapeno*). In addition, expressions that contained a typo (e.g., *blake* instead of *black*) could also not be classified automatically. One of the authors (EK) therefore manually hand-coded these cases.

There were 6 cases (0.3%) that could not be categorized. Those were mostly greetings (e.g., *Hi*), other comments (e.g., *I have instructions to follow sometimes*) and not certainly identifiable utterances (e.g., *re*). These were excluded. After exclusion, 1827 cases classified as one of *color*, *type*, or *color-and-type* entered the analysis.

#### 4.1.3 Typicality norming

In order to test for typicality effects on the production data and to evaluate RSA’s performance, we collected empirical typicality values for each utterance/object pair across three separate studies. The first study collected typicalities for *color-and-type*/object pairs (e.g., *yellow banana* as applied to a yellow banana, a blue banana, an orange pear, etc., see Figure 14a). The second study collected typicalities for *type-only*/object pairs (e.g., *banana* as applied to a yellow banana, a blue banana,



(a) color-and-type norming.

(b) type-only norming.

(c) color-only norming.

Figure 14: Example stimuli exemplifying the three different typicality norming studies.

| Utterances | Example              | Images      | Participants | Trials | Items     | Excluded |
|------------|----------------------|-------------|--------------|--------|-----------|----------|
| Adj Noun   | <i>yellow banana</i> | object      | 174          | 110    | 484       | 14       |
| Noun       | <i>banana</i>        | object      | 75           | 90     | 154 (198) | 1        |
| Adj        | <i>yellow</i>        | color patch | 110          | 90     | 176       | None     |

Table 5: Overview of typicality norming studies; the value in brackets shows the number of items including *fruit*, *vegetable* and *cup*.

an orange pear, etc., Figure 14b). The third study collected typicalities for *color/color* pairs (e.g., *yellow* as applied to a color patch of the average yellow from the yellow banana stimulus or to a color patch of the average orange from the orange pear stimulus, and so on, for all other colors, Figure 14c).

On each trial, participants saw one of the stimuli used in the production experiment in isolation and were asked: “How typical is this object for a *utterance*”, where *utterance* was replaced by an utterance of interest. In the color typicality study, they were asked “How typical is this color for the color *color*?”, where *color* was replaced by one of the relevant color terms. They then adjusted a continuous sliding scale with endpoints labeled “very atypical” and “very typical” to indicate their response. An overview of the differences between the three typicality norming studies is shown in Table 5.

Slider values were coded as falling between 0 (‘very atypical’) and 1 (‘very typical’). For each utterance-object combination, we computed mean typicality ratings. The means for the banana items are shown in Table 6. The values are very similar to those hypothesized for the purpose of the example in Table 4. The means for all items are displayed in XXX [jd: should we include all these, or just say the banana case is typical and put the data online?]

The typicality elicitation procedure we employed here is somewhat different from that employed by Westerbeek et al. (2015), who asked their participants “How typical is this color for this object?” We did this for conceptual reasons: the values that go into the semantics of the RSA model are most easily conceptualized as the typicality of an object as an instance of an utterance, rather than as the degree to which an object’s color is representative of that object. While the typicality of a feature for an object type no doubt plays into how good of an instance of the utterance the

| Utterance       | Banana items |       |      | Other |
|-----------------|--------------|-------|------|-------|
|                 | yellow       | brown | blue |       |
| “banana”        | .98          | .66   | .42  | .05   |
| “yellow banana” | .97          | .30   | .15  | .05   |
| “brown banana”  | .22          | .91   | .15  | .04   |
| “blue banana”   | .16          | .15   | .92  | .06   |

Table 6: Mean typicalities for banana items. Combinations where deterministic semantics would return TRUE are marked in boldface. [jd: can you add the color patch data included below in another three rows, and with all the formatting analogous to table 4, and an “other” utterance row that just has the mean of all the other non-banana options?]

object is, deriving our typicalities from the statistical properties of the distributions of features in objects is beyond the scope of this paper. We expect, however, that the simple TYPE-object typicalities most closely approximates the Westerbeek question because the employed objects are color-diagnostic – asking whether a blue or a yellow banana is a typical *banana* is similar to asking whether or not the bananas’ most salient property – their color – is typical.<sup>16</sup>

Mean typicality values obtained in the norming studies are used in the analyses and visualizations in the following.

#### 4.1.4 Results and discussion

Proportions of type-only (*banana*), color-and-type (*yellow banana*), color-only (*yellow*), and other (*funky carrot*) utterances are shown in Figure 15 as a function of the described item’s mean type-only (*banana*) typicality. Visually inspecting just the explicitly marked *yellow banana*, *brown banana*, and *blue banana* cases suggests a large typicality effect in the overinformative conditions as well as a smaller typicality effect in the informative conditions, such that color-only and color-and-size utterances are less likely to be produced with increasing typicality of the object. [ek: @judith: but this is only true for the overinformative conditions]

The following questions are of interest. First, do we replicate the previously documented typicality effect on redundant color mention beyond the one-item visual inspection? Second, does typicality affect color mention even when color is informative (i.e., technically necessary for establishing unique reference)? Third, are speakers sensitive to the presence of color competitors in their use of color or are typicality effects immune to the nature of the distractor items?

To address these questions we conducted a mixed effects logistic regression predicting color use from fixed effects of typicality, informativeness, and color competitor presence. We used the typicality norms obtained in the *type/object* typicality elicitation study [ek: maybe “*type-only/object typicality elicitation study*” since you called it like that above] reported above (see Figure 14b) as the continuous typicality predictor. Informativeness was coded as a binary variable (color informative vs. color overinformative) as was color competitor presence (absent vs. present). All predictors

<sup>16</sup>See also Appendix E for an independent comparison of our question and the Westerbeek question as applied to typicality norms for the items in Exp. 1. In general, the TYPE-object values are highly correlated with the Westerbeek question values.

were centered before entering the analysis. The model included by-speaker and by-item random intercepts, which was the maximal random effects structure that allowed the model to converge.

There was a main effect of typicality, such that the more typical an object was for the type-only utterance, the lower the log odds of color mention ( $\beta = -4.20$ ,  $SE = 0.45$ ,  $p < .0001$ ), replicating previously documented typicality effects. Stepwise model comparison revealed that including interaction terms was not justified by the data, suggesting that speakers produce more typical colors less often even when the color is in principle necessary for establishing reference (i.e., in the informative conditions). This is notable: speakers are likely to call a yellow banana simply a *banana* even when other bananas are present, presumably because they can rely on listeners drawing the inference that they must have meant the most typical banana. [ek: judith: can we say it like this? In the informative-cc condition, we have 6 type-only utterances in the typical bin (typicality  $i = 0.784$  and zero in the atypical bin (out of 481 utterances in that condition in total); but is that enough for "they are likely"?)]

There was also a main effect of informativeness, such that color mention was less likely when it was overinformative than when it was informative ( $\beta = -5.57$ ,  $SE = 0.33$ ,  $p < .0001$ ). Finally, there was a main effect of color competitor presence, such that color mention was more likely when a color competitor was absent ( $\beta = 0.71$ ,  $SE = 0.16$ ,  $p < .0001$ ). This suggests that speakers are indeed sensitive to the contextual utility of color – color typicality alone does not capture the full set of facts about color mention, as we already saw in Section 3.

## 4.2 Model evaluation

We evaluated the continuous semantics RSA model on the obtained production data from Exp. 2. While the architecture of the model remained the same as that of the model presented in Section 2.2, we briefly review the minor necessary changes, some of which we already mentioned at the beginning of this section. These changes concerned the lexicon and the cost function. We elaborate on each in turn.

### 4.2.1 Lexicon

Whereas for the purpose of evaluating the model in Section 3 we only considered the utterance alternatives *color*, *size*, and *color-size*, we included in this lexicon each color adjective, type noun, and combination of the two. This substantially increased the size of the lexicon to 37 unique utterances. For each combination of utterance  $u$  and object  $o$  that occurred in the experiment, we included a separate semantic value  $x_{u,o}$ , elicited in the norming experiments described in Section 4.1.3 (rather than inferred as done for Exp. 1).<sup>17</sup> For any given context, we assumed the utterance alternatives that correspond to the individually present features and their combinations. For example, for the context in Figure 13d, the set of utterance alternatives was *yellow*, *green*, *pear*, *banana*, *avocado*, *yellow pear*, *green pear*, *yellow banana*, *green banana*, *yellow avocado*, and *green avocado*. [jd: elisa, is this true? [ek: I'm not sure how it is handled in the unified model and when I want to look into it, I get an error. Robert can probably say more to this.]]

---

<sup>17</sup>Ideally, one would like to derive the semantic values for the modified utterances from the semantic values of the individual lexical items. Unfortunately, the compositionality of continuous values is a recognized problem (Kamp & Partee, 1995) and not one we aim to solve here. For our purposes it was therefore sufficient that we had access to the semantic values of the modified utterances elicited experimentally.

### 4.2.2 Cost function

[ek: judith: you already started the former paragraph the same way] Whereas for the purpose of evaluating the model in Section 3 we inferred two constant costs (one for color and one for size), we included here a more complex cost function. In particular, we defined utterance cost  $c(u)$  as follows:

$$c(u) = \beta_F \cdot p(u) + (1 - \beta_F) \cdot l(u) \quad (7)$$

Here,  $p(u)$  is utterance frequency as estimated from the Google Books corpus (years 1950 to 2008) [jd: elisa, was anything done to the values, eg log transform? [ek: the used frequencies are logFrequencies] was it frequencies or probabilities? [ek: I think those are frequencies, but since I'm not 100% sure what you mean, I could be wrong]];  $l(u)$  is the mean empirical length of the utterance in characters in the production data (e.g., sometimes *yellow* was abbreviated as *yel*, leading to an  $l(u)$  smaller than 6); and  $\beta_F$  is a weight that interpolates between length and frequency (when 1, cost is only a function of frequency; when 0, cost is only a function of length). Both  $p(u)$  and  $l(u)$  were normalized to fall into the interval  $[0, 1]$ . The cost function thus prefers short and frequent utterances (e.g., *blue*) over long and infrequent ones (*turquoise-ish bananaesque thing*).

### 4.2.3 Evaluation

[jd: XXX – thus far missing entirely. elisa, robert, what is the status of the bda? see the list below for what needs to go here]

- general procedure (same as for Exp. 1, exclude "other" utterances. same inference algorithm?)
- free parameters:  $\alpha, \beta_c, \beta_F$  – report priors over params [jd: any others? any sort of weight on typicality?]
- include plots: scatterplot of posterior predictive vs empirical proportions, and posterior over parameters
- report correlations with empirical data at different levels: at individual item (u,o) level collapsing across distractors; maybe collapsing across 3 different typicality bins, inspired by elisa's BSc thesis Fig 9? [jd: which others?]
- discuss posteriors over params – is alpha similar? is cost weight still low but with most weight on length?
- potential comparison with different models:
  - should we run a version that has 2 inferred but constant cost terms, one for color and one for type, to justify the more complex cost function?
  - should we run a version that has 2 inferred but constant semantic values, one for color and one for type, to compare with the empirically elicited values? this should do worse because you can't get the typicality effects without varying semantic values

### 4.3 Discussion

In this section we demonstrated that the continuous semantics RSA model predicts color typicality effects in the production of referring expressions. The model employed here did not differ in its architecture from that employed in Section 3, but only in that a) semantic values were assumed to operate at the individual utterance/object level (instead of at the utterance type/object level); b) semantic values were empirically elicited via typicality norming studies (instead of inferred from the data); and c) an utterance's cost was assumed to be a function of its mean empirical length and its corpus frequency as estimated from a large corpus (instead of having a constant utterance type level value).<sup>18</sup>

This suggests that the dynamics at work in the choice of color vs. size and in the choice of color as a function of the object's color typicality are very similar: speakers choose utterances by considering the fine-grained differences in information about the intended referent communicated by the ultimately chosen utterance compared to its competitor utterances. For noisier utterances (e.g., *banana* as applied to a blue banana), including the ‘overinformative’ color modifier is useful because it provides information. For less noisy utterances (e.g., *banana* as applied to a yellow banana), including the color modifier is useless because the unmodified utterance is already highly informative. These dynamics even lead to the color modifier being left out altogether even when there is another object of the same type present that is very atypical, simply because the literal listener asymmetry in probability of choosing the intended referent over the competitor object is big enough.

In the next section, we move beyond the choice of modifier and ask whether continuous semantics RSA provides a good account of content selection in referring expressions more generally. To answer this question we turn to simple nominal referring expressions.

## 5 Unmodified referring expressions: nominal taxonomic level

In this section we investigate whether continuous semantics RSA can account for referring expression production beyond the choice of modifier. To do so, we begin by presenting a second production experiment. This experiment investigates speakers' choice of level of reference in nominal referring expression (*dalmation*, *dog*, or *animal*). As discussed in Section 5, multiple factors have been shown to play a role in the choice of nominal referring expression, including an expression's contextual informativeness, its cognitive cost (short and frequent terms are preferred over long and infrequent ones) cite cite, and its typicality (an utterance is more likely to be used if the object is a good example of it) is that true? cite. yes, caroline put ref in cogsci talk. We then evaluate continuous semantics RSA on the nominal choice dataset by conducting the same type of Bayesian data analysis as reported in the previous section.

### 5.1 Experiment 3: taxonomic level of reference in nominal referring expressions

Exp 2 employed the same procedure as Exp. 1, but each display consisted of three objects.<sup>19</sup> We manipulated the contextual informativeness of each level of reference – subordinate (*dalmatian*), basic (*dog*), and superordinate (*animal*) – by manipulating the distractor items.

<sup>18</sup>See Table 8 for a more extensive overview of the ways in which the models reported across sections differed.

<sup>19</sup>Exp. 2 constitutes a replication of Graf et al. (2016).

### 5.1.1 Method

**Participants** We recruited 58 pairs of participants (116 participants total, the same participants as in Exp. 1) over Amazon’s Mechanical Turk who were each paid \$1.75 for their participation.

**Procedure and materials** The procedure was identical to that of Exp. 1. Participants proceeded through 72 trials. Of these, half were critical trials of interest and half were filler trials (the critical trials from Exp. 1). On critical trials, we varied the level of reference that was sufficient to mention for uniquely establishing reference.

Stimuli were selected from nine distinct domains, each corresponding to distinct basic level categories such as *dog*. For each domain, we selected four subcategories to form our target set (e.g. *dalmatian*, *pug*, *German Shepherd* and *husky*). See Table 11 in Appendix G for a full list of domains and their associated target items. Each domain also contained an additional item which belonged to the same basic level category as the target (e.g., *greyhound*) and items which belonged to the same supercategory but not the same basic level (e.g., *elephant* or *squirrel*). The latter items were used as distractors.

Each trial consisted of a display of three images, one of which was designated as the target object. Each pair of participants saw each target exactly once, for a total of 36 trials per pair. These target items were randomly assigned distractor items which were selected from four different context conditions, corresponding to different communicative pressures (see Figure 16). We refer to these conditions with pairs of numerals specifying which levels of the taxonomy are present in the distractors: (a) item12 contexts contain one distractor of the same basic level and one distractor of the same superlevel (e.g., target: *dalmatian*, distractor 1: *greyhound* (also a dog), distractor 2: *squirrel* (also an animal)); (b) item22 contexts contain two distractors of the same superlevel but different basic level as the target (e.g., target: *husky*, distractors: *hamster* and *elephant*); (c) item23 contexts contain one distractor of the same superlevel and one unrelated item (e.g., target: *pug*, distractor 1: *cow*, distractor 2: *table*); and (d) item33 contexts contain two unrelated items (e.g., target: *German Shepherd*, distractors: *shirt* and *cookie*).

This context manipulation served as a manipulation of utterance informativeness: any target could be referred to at the sub (*dalmatian*), basic (*dog*) or super (*animal*) level. However, the level of reference necessary for uniquely referring differed across contexts: in item12 contexts, the sub level was necessary. In item22 and item23 contexts, the basic level was necessary (though the sub level was also possible). In item33 contexts all three utterances were possible.

### 5.1.2 Data pre-processing and exclusion

We collected 2187 referring expressions. To determine the level of reference for each trial, we followed the following procedure. First, 41 trials on which the listener selected the wrong referent were excluded, leading to the elimination of 1.9% of trials. Then, speakers’ and listeners’ messages were parsed automatically; the referring expression used by the speaker was extracted for each trial and checked for whether it contained the current target’s correct sub, basic or super level term using a simple grep search. In this way, 72.1% of trials were labelled as mentioning a pre-coded level of reference. In the next step, remaining utterances were checked manually to determine whether they contained a correct level of reference term which was not detected by the grep search due to typos or grammatical modification of the expression. In this way, meaning-equivalent alternatives such as *doggie* for *dog*, or reduced forms such as *gummi*, *gummies* and *bears* for *gummy bears* were

counted as containing the corresponding level of reference term. This covered another 15.1% of trials. A total of 12.8% of correct trials were excluded because the utterance consisted only of an attribute of the superclass (*the living thing* for *animal*), of the basic level (*can fly* for *bird*), of the subcategory (*barks* for *dog*) or of the particular instance (*the thing facing left*) rather than a category noun. These kinds of attributes were also mentioned in addition to the noun on trials which were included in the analysis for 8.9% of sub level terms, 19.1% of basic level terms, and 66.7% of super level terms. On 1.2% of trials two different levels of reference were mentioned; in this case the more specific level of reference was counted as being mentioned in this trial. After all exclusion and pre-processing, 1870 cases classified as one of *sub*, *basic*, or *super* entered into the analysis.

### 5.1.3 Results and discussion

Proportions of sub, basic, and super level utterances are shown in the top row of Figure 17. Overall, super level mentions are highly dispreferred (< 2%), so we focus in this section only on predictors of sub over basic level mentions. The clearest pattern of note is that sub level mentions are only preferred in the most constrained context that necessitates the sub level mention for unique reference (item12, e.g. target: dalmatian, distractor: greyhound). Nevertheless, even in these contexts there is a non-negligible proportion of basic level mentions (28%). In the remaining contexts, where the sub and basic level are equally informative, there is a clear preference for the basic level.

What explains these preferences? In order to test for effects of informativeness, length, frequency, and typicality on nominal choice we conducted a mixed effects logistic regression predicting sub over basic level mention from centered predictors for the factors of interest and the maximal random effects structure that allowed the model to converge (random by-speaker and by-target intercepts).

*Frequency* was coded as the difference between the sub and the basic level's log frequency, as extracted from the Google Books Ngram English corpus ranging from 1960 to 2008.

*Length* was coded as the ratio of the sub to the basic level's length. We used the mean empirical lengths in characters of the utterances participants produced. For example, the minivan, when referred to at the subcategory level, was sometimes called "minivan" and sometimes "van" leading to a mean empirical length of 5.71. This is the value that was used, rather than 7, the length of "minivan". That is, a higher frequency difference indicates a *lower* cost for the sub level term compared to the basic level, while a higher length ratio reflects a *higher* cost for the sub level term compared to the basic level.<sup>20</sup>

*Typicality* was coded as the ratio of the target's sub to basic level label typicality.<sup>21</sup> That is, the higher the ratio, the more typical the object was for the sub level label compared to the basic level; or in other words, a higher ratio indicates that the object was relatively atypical for the basic label compared to the sub label. For instance, the panda was relatively atypical for its basic level "bear" (mean rating 0.75) compared to the sub level term "panda bear" (mean rating 0.98), which resulted in a relatively *high* typicality ratio.

*Informativeness* condition was coded as a three-level factor: *sub necessary*, *basic sufficient*, and *super sufficient*, where item22 and item23 were collapsed into *basic sufficient*. Condition was Helmert-coded: two contrasts over the three condition levels were included in the model, comparing

---

<sup>20</sup>We replicate the well-documented negative correlation between length and log frequency ( $r = -.49$  in our dataset).

<sup>21</sup>Typicalities were elicited in a separate norming study that was identical in procedure to that of Exp. 1a. See Appendix F for details about the study.

each level against the mean of the remaining levels (in order: *sub necessary*, *basic sufficient*, *super sufficient*). This allowed us to determine whether the probabilities of type mention for neighboring conditions were significantly different from each other, as suggested by Figure 17.

The log odds of mentioning the sub level term were greater in the *sub necessary* condition than in either of the other two conditions ( $\beta = 2.05$ ,  $SE = .17$ ,  $p < .0001$ ), and greater in the *basic sufficient* condition than in the *super sufficient* condition ( $\beta = .54$ ,  $SE = .15$ ,  $p < .001$ ), suggesting that the contextual informativeness of the sub level mention has a gradient effect on utterance choice.<sup>22</sup> There was also a main effect of typicality, such that the sub level term was preferred for objects that were more typical for the sub level compared to the basic level description ( $\beta = 4.84$ ,  $SE = 1.32$ ,  $p < .001$ , see Figure 18). In addition, there was a main effect of length, such that as the length of the sub level term increased compared to the basic level term (“chihuahua”/“dog” vs. “pug”/“dog”), the sub level term was dispreferred (“chihuahua” is dispreferred compared to “pug”,  $\beta = -.95$ ,  $SE = .27$ ,  $p < .001$ , see Figure 18). The main effect of frequency did not reach significance ( $\beta = .07$ ,  $SE = .10$ ,  $p < .51$ ).

Unsurprisingly, there was also significant by-participant and by-domain variation in sub level term mention. For instance, mentioning the sub over the basic level term was preferred more in some domains (e.g. in the “candy” domain) than in others. Likewise, some domains had a greater preference for basic level terms (e.g. the “shirt” domain). Using the super term also ranged from hardly being observable (e.g. the “flower” domain) to being used more frequently (e.g. in the “table” and “car” domain).

We thus replicate the well-documented preference to refer to objects at the basic level, which is partly modulated by contextual informativeness and partly a result of the basic level term’s cognitive cost and typicality compared to its sub level competitor.

Perhaps surprisingly given the previous literature, we did not observe an effect of frequency on sub level term mention. This may have a number of reasons. For instance, the modality of the experiment may have mattered here: the current study was a written production study, while most studies that have identified frequency as a factor governing production choices are spoken production studies (cite cite). It may be that the cognitive cost of typing longer words may be disproportionately higher than that of producing longer words in speech, thus obscuring a potential effect of frequency.

## 5.2 Model evaluation

Here we show that continuous semantics RSA as presented in Section 2.2 can be straightforwardly extended to modeling the choice of taxonomic level of reference. We include three modifications, while leaving the general framework as is. The first modification concerns the utterance alternatives. The second concerns the elicited typicality values and the resulting fidelity values. The third concerns the cost function. We briefly elaborate on each in turn.

**Utterance alternatives.** Whereas the modifier choice model treats all individual features and feature combinations represented in the display as utterance alternatives, the nominal choice model considers only the three different levels of reference to the target as alternatives, e.g., *dalmatian*, *dog*, *animal*. That is, assuming a German Shepherd as a distractor, *German Shepherd* is not considered

---

<sup>22</sup>Importantly, model comparison between the reported model and one that subsumes basic and super under the same factor level revealed that the three-level condition variable is justified ( $\chi^2(1) = 12.82$ ,  $p < .0004$ ), suggesting that participants don’t simply revert to the basic level unless contextually forced not to.

an alternative. This has consequences for the assumed fidelity values, which we turn to next. [jd: we should probably discuss this in the GD? ie, if we also assumed distractor labels as alternatives, we would have to do the rescaling – would results be different? or the other way round: if we assume in modifier choice only the target’s features are available as alternatives, would results be different?]

**Fidelity values.** Just as we did for capturing color typicality effects in Section ??, we elicited empirical typicality values for object-utterance combinations.<sup>23</sup> For each display, we know the typicality of each object in the display as an instance of the three potential target utterances (capturing, for instance, that the word “dog” describes a dalmatian better than a grizzly bear, but it also describes a grizzly bear better than a tennis ball). This allows us to use the typicality values as fidelity values directly, without rescaling as was necessary in the modifier choice model.

**Cost function.** Recall the pragmatic speaker’s utility function from Section 2.2, where the weighted informativeness term  $\alpha \ln P_{L_0}(o|u)$  traded off against the weighted utterance cost  $\beta_c c(u)$ . In the modifier model we assumed a constant cost for each added modifier. Because all utterance alternatives in the nominal choice model have word length 1, we update the cost function to be composed of each utterance’s length  $\hat{c}_l$  and frequency  $\hat{c}_f$  (as described in the previous section), weighted by free parameters  $\beta_f$  and  $\beta_l$ :

$$P_{S_1}(u|o) \propto e^{\alpha \ln P_{L_0}(o|u) + \beta_f \hat{c}_f + \beta_l \hat{c}_l} \quad (8)$$

To understand the qualitative behavior of the model, we briefly delve into two aspects of the model: first, the effect of typicality on the literal listener (and, in consequence via the pressure to be informative) the speaker. And second, the effect of cost (utterance length and frequency) on the speaker.

### 5.2.1 Typicality effects

**Literal listener behavior.** The literal listener’s probability of choosing the target under different typicalities for the observed utterance are shown in Figure 19. In general: as the target’s typicality as an instance of the utterance increases and the distractors’ typicality decreases, the probability of the literal listener choosing the target increases. Subordinate level terms tend to fall in the upper right quadrant of this graph. Basic level terms in the *sub necessary* conditions tend to fall in the lower right quadrant, while basic level terms in the *basic sufficient* conditions tend to fall in the upper right quadrant as well.

**Pragmatic speaker behavior.** To understand the effect of typicality on the speaker’s behavior it is useful to think about the problem of deciding which taxonomic level to refer at in terms of typicality gain, as we did in Section 4 for the choice between modified and unmodified expression. There, we found that relatively large target (compared to distractor) typicality gains in going from unmodified to modified expressions compared resulted in greater probability of overmodification. Here we observe the same effect in going from a higher (less specific) to a lower (more specific) taxonomic level. This can be seen in Figure 20, which shows the probability of each utterance (sub,

---

<sup>23</sup>See Appendix F for details of typicality elicitation experiment.

Table 7: Overview of simulated distractor typicality (fidelity) values for sub, basic, and super level utterances in simulated conditions. In contrast to the actual experimental conditions, we assume equal typicality values for both distractors.

|           |       | Condition     |                  |                  |
|-----------|-------|---------------|------------------|------------------|
|           |       | sub necessary | basic sufficient | super sufficient |
| Utterance | sub   | 0             | 0                | 0                |
|           | basic | .8            | .1               | 0                |
|           | super | .8            | .8               | 0                |

basic, or super) as a function of absolute target typicality as well as target typicality gain. Target typicality gain is the difference between the target’s sub level typicality and the target’s basic level typicality. Probabilities are shown for contexts with three items, always assuming  $\alpha = 7$ , but manipulating distractor typicality to simulate conditions analogous to our experimental conditions *sub necessary*, *basic sufficient*, and *super sufficient*. Simulated distractor typicalities for sub, basic, and super level reference are shown in Table 7.

The blue areas in the graph indicate highest-probability regions. For example, as expected in the *sub necessary* condition, the sub level term is the most likely one. However, in certain cases the basic level term also receives non-zero probability, notably when the target is a better instance of the basic than the sub level term, or (not pictured) when the typicality of the distractor as an instance of the basic level term is very low (e.g., the typicality of the koala bear as an instance of “bear” was only 0.50). Indeed, the grizzly (with high typicality for basic level “bear”, .97) is referred to as “bear” rather than “grizzly bear” in 85% of *sub necessary* conditions when the koala is the distractor.

In the *basic sufficient* conditions, sub level reference is nevertheless strongly predicted when target sub typicality gain is positive (i.e., when the target is a much better instance of the sub than of the basic level term). An example of such a case is the panda bear, who received a sub level typicality of .98 and a basic level typicality of only .75. Indeed, even when basic level reference was sufficient, the panda was referred to as the “panda” 81% of the time.

These patterns mirror the typicality effects obtained via the mixed effects regression.

### 5.2.2 Cost effects

The additional effect of cost on nominal choice is straightforward: the costlier an utterance (relative to its alternatives), the less likely it is to be used. This pattern, too, is one observed in the mixed effects regression. For instance, the (short, less costly) pug is almost three times as likely as the (long, more costly) German Shepherd to be referred to by its subordinate level term in the *basic sufficient* and *super sufficient* conditions, where subordinate level reference is unnecessary.

In Section 5.2 we showed that continuous semantics RSA captures the right kinds of qualitative effects as observed in the mixed effects regression. In the next section we evaluate how well the model captures nominal choice preferences quantitatively.

### 5.3 Model evaluation: nominal choice

In order to evaluate continuous semantics RSA for nominal choice, we repeated the same Bayesian data analysis as reported in Section 3.2 and Section ?? to generate model predictions and infer likely parameter values. We did so by conditioning on the observed production data (coded into *sub*, *basic*, and *super* level mentions as described above) and integrating over the three free parameters  $\alpha \sim \mathcal{U}(0, 20)$ ,  $\beta_f \sim \mathcal{U}(0, 5)$ ,  $\beta_l \sim \mathcal{U}(0, 5)$ .

Point-wise maximum a posteriori (MAP) estimates of the model’s posterior predictives for each combination of utterance and informativeness condition (collapsing across different items) are compared to empirical data in Figure 17. The model clearly captures the preference towards sub level mentions in the *sub necessary* conditions and the basic level preference in all other conditions. It also captures the further decrease in sub level mentions in the *super sufficient condition*. However, it does overpredict super level mentions, though not as badly as models that either assume a deterministic semantics or that ignore utterance cost.<sup>24</sup> At this level, the model achieves a correlation of  $r = .94$ . Computing correlations additionally on the by-target level yields a correlation of  $r = .84$  (see also the scatterplot in Figure 21).

Parameter posteriors are shown in Figure 22. Both informativeness and length receive significant weight. In contrast, the effect of frequency appears to be much weaker with a MAP of .1 and the HDIs overlapping with 0. This mirrors the null effect of frequency found in the regression analysis. However, a large number of cases also received a non-zero frequency weight.

In order to ascertain whether typicality as incorporated in the continuous semantics was indeed contributing to the explanatory power of the model, we ran an additional Bayesian data analysis with an added typicality weight parameter  $\beta_t \in [0, 1]$ . This parameter interpolated between empirical typicality values (when  $\beta_t = 1$ ) and deterministic (i.e., 0 or 1) a priori values based on the true taxonomy (when  $\beta_t = 0$ ). We found a MAP estimate for  $\beta_t$  of .95, HDI = [0.82,.99], strongly indicating that it is useful to incorporate empirical typicality values and thus providing further support for the value of non-deterministic truth functions in modeling referring expressions.

## 6 General Discussion

[jd: make sure to include adele goldberg’s typicality effects in kids’ generalizations to the basic level in discussion (srcd poster and submitted manuscript with lauren emberson)]

[jd: These effects can also be thought of in terms of Levinsonian stereotype inferences (Levinson, 2000): the speaker reasons that just saying *banana* will lead a listener to believe they are referring to a stereotypical instance of banana (i.e., a yellow one). In order to prevent this (incorrect) stereotype inference they add a modifier to refer to the blue one.]

### 6.1 Summary

How do speakers choose a referring expression? Here we have shown that they do so by trading off various factors: the contextual informativeness of the referring expression on the one hand, and the cognitive cost of the expression on the other. Importantly, computing contextual informativeness with respect to a *non-deterministic* underlying semantics was crucial for capturing various aspects

<sup>24</sup>The reader is referred to Appendix H for a comparison of the models containing a) only informativeness with deterministic semantics; b) only informativeness with continuous semantics; c) informativeness with deterministic semantics and cost; d) informativeness with continuous semantics and cost (the current model).

Table 8: Overview of the models used for the three different production datasets Color/size (Exp. 1), Color typicality (Exp. 2), and Nominal choice (Exp. 3).

|                                  | Color/size  | Color typicality  | Nominal choice   |
|----------------------------------|---|---|--|
| Lexicon $\mathcal{L}(u, o)$ size | 2 (color, size)   | 1 for each $u, o$ combination   | 1 for each $u, o$ combination                                  |
| Semantic values were...          | inferred  | elicited experimentally   | elicited experimentally  |
| Compositionality                 | $\mathcal{L}(u_{\text{size}}, o) \times \mathcal{L}(u_{\text{color}}, o)$ | $\mathcal{L}(u_{\text{color}}, o) \times \mathcal{L}(u_{\text{type}}, o)$ | NA   |
| Cost function $c(u)$             | 2 constant values (1 each for color and size)                             | $\beta_F p(u) + (1 - \beta_F)l(u)$  | $\beta_F p(u) + (1 - \beta_F)l(u)$                             |
| Costs were...                    | inferred  | estimated: $l(u)$ from production data and $p(u)$ from corpora            | estimated: $l(u)$ from production data and $p(u)$ from corpora |
| Set of alternatives              | all contextually available feature combinations (color, size)             | all contextually available feature combinations (type, color)             | 3 target alternatives (level of reference)                     |
| Grammatical penalty              | none  | [jd: negative type cost param?]   | none   |
| [jd: Typicality weight]          | none  | [jd: yes?]  | [jd: yes?]   |

of speakers’ referring behavior. First, the continuous semantics allowed us to capture the basic well-documented asymmetry for speakers to be more likely to redundantly use color adjectives rather than size adjectives. In addition, it predicted an interaction between sufficient dimension and scene variation on the probability of redundancy, which was very clearly borne out in the data: increased scene variation resulted in a much greater increase in redundant color than in redundant size adjective use. Finally, the non-determinism in the semantics gave rise to well-documented effects of typicality in both modifier choice and noun choice. A modifier was more likely to be mentioned redundantly when the object was a substantially less good instance of the unmodified than of the modified expression. Analogously, a noun at a taxonomically lower level than necessary for establishing reference was more likely to be mentioned when the object was a substantially less good instance of the higher than of the lower level.

We have thus shown that with one key innovation – a continuous semantics – one can retain the assumption that speakers rationally trade off informativeness and cost of utterances in language production. Rather than being wastefully overinformative, adding redundant modifiers or referring at a lower taxonomic level than strictly necessary *is* in fact informative when the *prima facie* sufficiently informative expression is substantially noisier than its redundant/overly specific counterpart. This innovation thus not only provides a unified explanation for a number of key patterns within the overinformative referring expression literature that have thus far eluded a unified explanation; it also extends to the domain of nominal choice.

In the following we discuss a number of intriguing questions this work raises and avenues for future research that it suggests.

## 6.2 ‘Overinformativeness’

This work challenges the traditional notion of overinformativeness in the linguistic and psychological literature (Engelhardt, Bailey, & Ferreira, 2006b; ?, ?). The reason that redundant referring expressions became interesting for psycholinguists to study is because they seem to constitute a clear violation of rational theories of language production. For example, Grice’s Quantity-2 maxim, which asks of speakers to “not make [their] contribution more informative than is required” (Grice, 1975), appears violated by any redundant referring expression – if size is the only feature that distinguishes the target object from the rest, the mention of color seems more informative than required.

This conception of (over-)informativeness assumes that all modifiers are born equal – i.e., that there are no a priori differences in the utility of mentioning different properties of an object. Under this conception of modifiers, there are hard lines between modifiers that are and aren’t informative in a context. However, what we have shown here is that under a continuous semantics, a modifier that would be regarded as overinformative under the traditional conception may nevertheless add some information about the referent. In particular, the more visual variation there is in the scene and the less noisy the redundant modifier is compared to the modifier that selects the dimension that uniquely singles out the target, the more information it adds about the referent, and the more likely it therefore is to be mentioned. This work thus challenges the traditional notion of utterance overinformativeness by providing an alternative that nicely captures the quantitative variation observed in speakers’ production in a principled way while still assuming that speakers are aiming to be informative.

What, then, would count as an overinformative utterance under continuous semantics RSA? RSA shifts the bar for overinformativeness and turns it into a graded notion: the less expected the use of a redundant modifier is (given knowledge of, e.g., utterance noise, cost, scene variation, and typicality), the more the use of that modifier will be considered overinformative.

## 6.3 Comprehension

While the account proposed in this paper is not directly concerned with predicting listeners’ behavior in interpreting referring expressions, it can be extended to do so relatively straightforwardly. RSA models typically assume that listeners, in interpreting utterances, are doing so by reasoning about their model of the speaker. In this paper we have provided precisely such a model of the speaker. In what way should the predicted speaker probabilities enter into comprehension? Here we can make a direct connection to surprisal theory in sentence processing (?, ?), where it has been shown that the effort involved in processing a sentence is a function of how surprising that sentence is under the listener’s language model. While in these studies surprisal is usually estimated from syntactically parsed corpora, here we are providing a speaker model from which we can derive estimates of *pragmatic surprisal*. Generally, the more likely a redundant utterance is, the easier it should be to process in context. We have shown that redundant expressions are more likely than minimal expressions when the distinguishing dimension is relatively noisy and scene variation is relatively high. In situations like these, one would thus expect the redundant expression to be easier to process than in cases where the redundant expression is relatively less likely.

Is there evidence that listeners do behave in accordance with this prediction? While we have not run processing studies ourselves, we can look into the literature. Indeed, there is evidence that in situations where the redundant modifier does provide some information about the referent,

Table 9: Fidelity across models alongside the effects from Table 2 that each model captures.

| Exp. | Model                            | Fidelity level   | How obtained                                  | Effect(s)   |
|------|----------------------------------|------------------|---|---|
| 1    | basic non-deterministic          | modifier type    | inferred                                      | color/size asymmetry & scene variation                    |
| 1    | typicality (modified/unmodified) | utterance-object | elicited (nominal, color) and inferred (size) | color/size asymmetry, scene variation, & color typicality |
| 2    | typicality (level of reference)  | utterance-object | elicited                                      | basic level preference & subordinate level mention        |

listeners are faster to respond and select the intended referent when they observe a redundant referring expression than when they observe a minimal one (Arts et al., 2011; Paraboni et al., 2007). However, there is also evidence that redundancy sometimes incurs a processing cost: both Engelhardt, Demiral, and Ferreira (2011) and Davies and Katsos (2013) (Exp. 2) found that listeners were slower to identify the target referent in response to redundant compared to minimal utterances. It is useful to examine the stimuli they used. In the Engelhardt et al study, there was only one distractor that varied in type, i.e., type was sufficient for establishing reference. This distractor varied either in size or in color. Thus, scene variation was very low and overinformative expressions therefore likely surprising. Interestingly, the incurred cost was greater for redundant size than for redundant color modifiers, in line with the RSA predictions that color should be generally more likely to be used redundantly than size. In the Davies et al study, the ‘overinformative’ conditions contained displays of four objects which differed in type. Stimuli were selected via a production pre-test: only those objects that in isolation were not referred to with a modifier were selected for the study. That is, stimuli were selected precisely on the basis that redundant modifier use would be unlikely.

While the online processing of redundant referring expressions is yet to be systematically explored under the continuous semantics RSA account, this cursory overview of the patterns reported in the existing literature suggests that pragmatic surprisal (i.e., negative log-transformed speaker probabilities) may be a plausible linking function from model predictions to processing times.

#### 6.4 Fidelity

The model crucially relies on a continuous semantics to capture the effects we have reported in this paper. But what is the nature of this non-determinism? What does it represent? For the purpose of Exp. 1 (modifier choice), fidelity initially applied at the modifier *type* level. The semantics of modifiers was underlyingly truth-conditional and the fidelity term captured the probability that a modifier’s truth conditions would accidentally be inverted. This model included only two fidelity terms, one for size and one for color. We then extended the notion of fidelity to apply at the level of utterance-object combinations (e.g., *golf ball* vs. *pink golf ball* as applied to a pink golf ball) to account for color typicality effects. In this instantiation of the model, fidelity differed for every utterance-object combination and captured how good of an instance of an utterance an object was. Similarly, in Exp. 2 (nominal choice) fidelity differed for every utterance-object combination (e.g., *dog* vs. *dalmatian* as applied to a dalmatian). This is summarized in Table 9.

What we have said nothing about thus far is where these numbers come from; in particular,

which aspects of our experience – linguistic, perceptual, conceptual, communicative – they represent. We will offer some speculative remarks and directions for future research here.

First, it is possible that the numbers represent the difficulty associated with verifying whether the property denoted by the utterance holds of the object. This difficulty may be perceptual – for example, it may be relatively easier to visually determine of an object whether it is red than whether it is big. Similarly, at the object-utterance level, it may be easier to determine of a yellow banana than of a blue banana whether it exhibits banana-hood, in consequence yielding a lower typicality value for a blue banana than for a yellow banana as an instance of *banana*. It may also be conceptual – for example, it may be easier to determine whether a box belongs to John than whether **XXX**.

Another possibility is that the numbers represent aspects of agents' prior beliefs (world knowledge) about the correlations between features of objects. For example, conditioning on bananahood holding of objects and asking for the relative probabilities of various colors obtaining in that set will yield a high number for yellow and a low one for blue.<sup>25</sup>

Another hypothesis is that the numbers capture the past probability of communicative success in using a particular utterance (e.g., *banana*) to refer to an object with a particular set of features (e.g., blue bananas vs. yellow bananas). However, this probability is likely itself not independent of the first two possibilities discussed.

Finally, it is also possible that the numbers are simply an irreducible part of the lexical entry of each utterance-object pair. This seems unlikely, given that this would require a separate number for each utterance and object token. It also suggests that the numbers should not be updated in response to further exposure of objects. For example, if the numbers were a fixed component of the lexical entry *banana*, then even being exposed to a large number of blue bananas should not change the value. This seems unlikely but deserves to be investigated further.

## 6.5 Audience design

One question which has plagued the literature on language production is that of whether, and to what degree, speakers actually tailor their utterances to their audience (Clark & Murphy, 1982; Horton & Keysar, 1996; Brown-Schmidt & Heller, 2014). This is also known as the question of *audience design*. With regards to redundant referring expressions, the question is whether speakers produce redundant expressions because they can't help it (i.e., due to internal production pressures) or specifically because it is helpful for their interlocutor (i.e., due to considerations of audience design).

continuous semantics RSA seems to make a claim about this issue: the non-determinism is located in the literal listener component, with respect to which speakers are trying to be informative. That is, it would seem that speakers produce referring expressions that are tailored to their listeners. However, this is misleading. The ontological status of the literal listener is as a “dummy component” that allows the pragmatic recursion to get off the ground. Actual listeners are, in line with previous work, more likely fall into the class of pragmatic  $L_1$  listeners; listeners who reason about the speaker's intended meaning via Bayesian inference (M. C. Frank & Goodman, 2012; Goodman & Stuhlmüller, 2013).<sup>26</sup>

---

<sup>25</sup>Though these probabilities cannot directly match up with the elicited typicality values, given that probabilities will have to sum up to 1, while typicality values were not normalized.

<sup>26</sup>But see Franke and Degen (2016) for an evaluation of the distribution of listener and speaker types in Quantity inferences.

Because RSA is a computational-level theory (Marr, 1982) of language use, it does not claim that speakers *actually*, *consciously* consult an internal model of a listener every time they choose an utterance, just that the distribution of utterances they use reflect informativity with respect to such a model. It is possible that this distribution is cached or computed using some other algorithm that doesn't explicitly involve a listener component.

Thus, the RSA model as formulated here remains agnostic about whether speakers' (over-)informativeness should be considered geared towards listeners' needs or simply a production-internal process.

## 6.6 Other factors affecting redundancy

continuous semantics RSA as presented in this paper straightforwardly accounts for effects of typicality, cost, and scene variation on redundancy in referring expressions. However, other factors have been identified as contributing to redundancy. For example, Rubio-Fernandez (2016) has shown that colors are mentioned more often redundantly for clothes than for geometrical shapes. Her explanation: knowing an object's color is generally more useful for clothing than it is for shapes. While she doesn't provide a detailed explanation for why this is the case, it is plausible that agents' knowledge of *goals* may be relevant here. For example, knowing the color of clothing is relevant to the goal of deciding what to wear or buy. In contrast, knowing the color of geometrical shapes is rarely relevant to any everyday goal agents might have. While the RSA model as implemented here does not accommodate an agent's goals, it can be extended to do so via projection functions, as has been done for capturing figurative language use (e.g., Kao, Wu, Bergen, & Goodman, 2014) or question-answer behavior (Hawkins, Stuhlmüller, Degen, & Goodman, 2015). This should be explored further in future research.

[jd: a note on incrementality? eg, pechmann says incrementality is to blame for redundancy: we retrieve words when we can, and colors are easier to retrieve, so we throw them out there regardless of whether or not they're redundant. The problem with this is that this makes a prediction about the order of adjectives; in particular, the preferred order should be reversed. Pechmann does find some instances of this, but not very many. But there are other ways incrementality could play a role. For example, throwing out the color word may help when the noun is hard to retrieve. This predicts that in languages with post-nominal adjectives, where you can't use this as a delay strategy for holding off on planning the noun, there should be less color redundancy; indeed, Rubio-Fernandez 2016 shows this for Spanish. The dynamic nature of language processing plays a role in other ways, too: it allows us to update our beliefs about individual speakers' use of modifiers and generate better expectations about upcoming input. For example, Pogue et al 2016 have shown that listeners, after being exposed to consistently overinformative speakers, stop drawing early contrastive inferences based on modifier use.]

## 6.7 Extensions to other language production phenomena

In this paper, we have focused on providing an account of content selection (Gatt et al., 2013) in modified referring expressions on the one hand (i.e., when to mention an object's size or color) and in nominal referring expressions on the other (i.e., at which taxonomic level to refer to an object). Future work should investigate whether these models can be merged to jointly account for the choice of content expressed in modifiers and in nouns. Further, in order to scale up to more naturalistic conversational domains it will be necessary to consider richer language models.

Recall that we treated different color names (e.g., *pink* and *purple*) as simply a color mention. Similarly, we treated different nouns that clearly referred at the same level (e.g., *grizzly* and *grizzly bear*) as simple sub level mentions. For the purpose of predicting not only content selection but also utterance choice, a richer inventory of utterance alternatives will need to be explored. An interesting question is how this approach can be extended to other referring expressions mentioned in the Introduction, e.g., names, pronouns, or referring expressions with post-nominal modification.

However, future research should also investigate the very intriguing potential for this approach to be extended to any language production phenomenon that involves content selection. For example, there is a large literature on optional instrument mentions. P. Brown and Dell (1987) showed that atypical instruments are more likely to be mentioned than typical ones – if a stabbing occurred with an icepick, speakers prefer “The man was stabbed with an ice pick” rather than “The man was stabbed”. If instead a stabbing occurred with a knife, “The man was stabbed” is preferred over “The man was stabbed with a knife”). This is very much parallel to the case of atypical color mention. While P. Brown and Dell (1987)’s account of the effect is that speakers do or don’t mention instruments for speaker-internal ego-centric reasons, later evidence suggests an explanation that is rather more driven by audience design considerations. Lockridge and Brennan (2002) replicated the original finding in a story retelling scenario while also manipulating whether or not addressees saw pictures of the actions. Without pictures, speakers produced even more mentions of atypical objects (presumably to prevent addressees from forming a faulty mental model of the situation), suggesting that the typicality effect is in fact an audience design effect.

More generally, the approach should extend to any content selection phenomenon that affords a choice between a more or less specific chunk of linguistic signal. Whenever the chunk adds sufficient information, it should be included. This is related to surprisal theories of production like Uniform Information Density (UID, Jaeger, 2006; Levy & Jaeger, 2007; A. Frank & Jaeger, 2008; Jaeger, 2010), where it has been found that speakers are more likely to omit linguistic signal if the underlying meaning or syntactic structure is highly predictable. Importantly, UID diverges from ours in that ours is (thus far) an account of *content selection*, while UID is an account of the choice between meaning-equivalent alternative *utterances*.

## 6.8 Conclusion

In conclusion, we have provided an account of redundant referring expressions that challenges the traditional notion of overinformativeness, unifies multiple language production literatures, and has the potential for many further extensions. For the time being, we take this work to suggest that, rather than being wastefully overinformative, speakers are rationally redundant.

[jd: What else needs to be included in GD?]

## A Effects of semantic value on utterance probabilities

Here we visualize the effect of different adjective types’ semantic value on the probability of producing the insufficient color-only utterance (*blue pin*), the sufficient size-only utterance (*small pin*), or the redundant color-and-size utterance (*small blue pin*) to refer to the target in context Figure 1a under varying  $\alpha$  values, in Figure 23. This constitutes a generalization of Figure 4, which is duplicated in row 6 ( $\alpha = 30$ ).

## B Validation of interactive web-based written production paradigm

make sure to discuss why overall we have lower overspecification rates – probably because of color typicality!! we had pretty typical colors in our stimuli

## C Pre-experiment quiz

Before continuing to the main experiment, each participant was required to correctly respond “True” or “False” to the following statements. Correct answers are given in parentheses after the statement.

- The speaker can click on an object. (False)
- The listener wants to click on the object that the speaker is telling them about. (True)
- The target is the object which has the red circle around it. (False)
- Only the speaker can send messages. (False)
- There are a total of 72 rounds. (True)
- The locations of the three objects are the same for the speaker and the listener. (False)

## D Exp. 1 items

The following table lists all 36 object types from Exp. 1 and the colors they appeared in:

| Object        | Colors         | Object      | Colors         |
|---------------|----------------|-------------|----------------|
| avocado       | black, green   | balloon     | pink, yellow   |
| belt          | black, brown   | bike        | purple, red    |
| billiard ball | orange, purple | binder      | blue, green    |
| book          | black, blue    | bracelet    | green, purple  |
| bucket        | pink, red      | butterfly   | blue, purple   |
| candle        | blue, red      | cap         | blue, orange   |
| chair         | green, red     | coat hanger | orange, purple |
| comb          | black, blue    | cushion     | blue, orange   |
| flower        | purple, red    | frame       | green, pink    |
| golf ball     | blue, pink     | guitar      | blue, green    |
| hair dryer    | pink, purple   | jacket      | brown, green   |
| napkin        | orange, yellow | ornament    | blue, purple   |
| pepper        | green, red     | phone       | pink, white    |
| rock          | green, purple  | rug         | blue, purple   |
| shoe          | white, yellow  | stapler     | purple, red    |
| thumb tack    | blue, red      | tea cup     | pink, white    |
| toothbrush    | blue, red      | turtle      | black, brown   |
| wedding cake  | pink, white    | yarn        | purple, red    |

## E Typicality effects in Exp. 1

To assess whether we replicate the color typicality effects previously reported in the literature (Sedivy, 2003; Westerbeek et al., 2015; Rubio-Fernandez, 2016), we elicited color typicality norms for each of the items in Exp. 1 and then included typicality as an additional predictor of redundant adjective use in the regression analysis reported in Section 3.1.3.

### E.1 Methods

#### E.1.1 Participants

We recruited 60 participants over Amazon’s Mechanical Turk who were each paid \$0.25 for their participation.

#### E.1.2 Procedure and materials

On each trial, participants saw one of the big versions of the items used in Exp. 1 and were asked to answer the question “How typical is this for an  $X$ ?” on a continuous slider with endpoints labeled “very atypical” to “very typical.”  $X$  was a referring expression consisting of either only the correct noun (e.g., *stapler*) or the noun modified by the correct color (e.g., *red stapler*). Figure 24 shows an example of a modified trial.

Each participant saw each of the 36 objects once. An object was randomly displayed in one of the two colors it occurred with in Exp. 1 and was randomly displayed with either the correct modified utterance or the correct unmodified utterance, in order to obtain roughly equal numbers of object-utterance combinations.

Importantly, we only elicited typicality norms for unmodified utterances and utterances with color modifiers, but not utterances with size modifiers. This was because it is impossible to obtain size typicality norms for objects presented in isolation, due to the inherently relational nature of size adjectives. Consequently, we only test for the effect of typicality on *size-sufficient* trials, i.e. when color is redundant.

### E.2 Results and discussion

We coded the slider endpoints as 0 (“very atypical”) and 1 (“very typical”), essentially treating each response as a typicality value between 0 and 1. For each combination of object, color, and utterance (modified/unmodified), we computed that item’s mean. Mean typicalities were generally lower for unmodified than for modified utterances: mean typicality for unmodified utterances was .67 ( $sd=.17$ ,  $mode=.76$ ) and for modified utterances .75 ( $sd=.12$ ,  $mode=.81$ ). This can also be seen on the left in Figure 25. Note that, as expected given how the stimuli were constructed, typicality was generally skewed towards the high end, even for unmodified utterances. This means that there was not much variation in the difference in typicality between modified and unmodified utterances. We will refer to this difference as *typicality gain*, reflecting the overall gain in typicality via color modification over the unmodified baseline. As can be seen on the right in Figure 25, in most cases typicality gain was close to zero.

This makes the typicality analysis difficult: if typicality gain is close to zero for most cases (and, taking into account confidence intervals, effectively zero), it is hard to evaluate the effect of typicality on redundant adjective use. In order to maximize power, we therefore conducted the

Table 10: Model coefficients, standard errors, and p-values. Significant p-values are bolded.

|                                       | Coef $\beta$ | SE( $\beta$ ) | <i>p</i>         |
|---------------------------------------|--------------|---------------|------------------|
| Intercept                             | -1.85        | 0.34          | <b>&lt;.0001</b> |
| Scene variation                       | 4.29         | 1.16          | <b>&lt;.001</b>  |
| Sufficient property                   | 2.72         | 0.60          | <b>&lt;.0001</b> |
| Scene variation : Sufficient property | 0.88         | 2.12          | <0.68            |
| Sufficient property : Typicality gain | 9.43         | 2.68          | <b>&lt;.001</b>  |

analysis only on those items for which for at least one color the confidence intervals for the modified and unmodified utterances did not overlap. There were only four such cases: *(pink) golfball*, *(pink) wedding cake*, *(green) chair*, and *(red) stapler*, for a total of 231 data points.

Predictions differ for size-sufficient and color-sufficient trials. Given the typicality effects reported in the literature and the predictions of continuous semantics RSA, we expect greater redundant color use on size-sufficient trials with *increasing* typicality gain. The predictions for redundant size use on color-sufficient trials are unclear from the previous literature. continuous semantics RSA, however, predicts greater redundant size use with *decreasing* typicality gain: small color typicality gains reflect the relatively low out-of-context utility of color. In these cases, it may be useful to redundantly use a size modifier even if that modifier is noisy. If borne out, these predictions should surface in an interaction between sufficient property and typicality gain. Visual inspection of the empirical proportions of redundant adjective use in Figure 26 suggests that this pattern is indeed borne out.

In order to investigate the effect of typicality gain on redundant adjective use, we conducted a mixed effects logistic regression analysis predicting redundant over minimal adjective use from fixed effects of scene variation, sufficient dimension, the interaction of scene variation and sufficient property, and the interaction of typicality gain and sufficient property. This is the same model as reported in Section 3.1.3, with the only difference that the interaction between sufficient property and typicality gain was added. All predictors were centered before entering the analysis. The model contained the maximal random effects structure that allowed it to converge: by-participant and by-item (where item was a color-object combination) random intercepts.

The model summary is shown in Table 10. We replicate the effects of sufficient property and scene variation observed earlier on this smaller dataset. Crucially, we observe a significant interaction between sufficient property and typicality gain.<sup>27</sup> Simple effects analysis reveals that this interaction is due to a positive effect of typicality gain on redundant adjective use in the size-sufficient condition ( $\beta = 4.47$ ,  $SE = 1.65$ ,  $p < .007$ ) but a negative effect of typicality gain on redundant adjective use in the color-sufficient condition ( $\beta = -5.77$ ,  $SE = 2.49$ ,  $p < .03$ ).

An important point is of note: the typicality elicitation procedure we employed here is somewhat different from that employed by Westerbeek et al. (2015), who asked their participants “How typical is this color for this object?” We did this for conceptual reasons: the values that go into the

<sup>27</sup>Conducting the same analysis on the entire dataset (i.e., using all of the noisy typicality estimates, replicated the scene variation and sufficient property effects. The interaction of typicality gain and sufficient property went in the same direction numerically, but failed to reach significance ( $\beta = 1.52$ ,  $SE = 1.45$ ,  $p < .29$ ).

semantics of the RSA model are most easily conceptualized as the typicality of an object as an instance of an utterance. While the typicality of a feature for an object type no doubt plays into how good of an instance of the utterance the object is, deriving our typicalities from the statistical properties of the subjective distributions of features over objects is beyond the scope of this paper. However, in a separate experiment we did ask participants the Westerbeek question. The correlation between mean typicality ratings from the Westerbeek version and the unmodified “How typical is this for  $X$ ” version was .75. The correlation between the Westerbeek version and the modified version was .64. The correlation between the Westerbeek version and typicality gain was -.52.

For comparison, including typicality means obtained via the Westerbeek question as a predictor instead of typicality gain on the four high-powered items replicated the significant interaction between typicality and sufficient property ( $\beta = -6.77$ ,  $SE = 1.88$ ,  $p < .0003$ ). Simple effects analysis revealed that the interaction is again due to a difference in slope in the two sufficient property conditions: in the size-sufficient condition, color is less likely to be mentioned with increasing color typicality ( $\beta = -3.66$ ,  $SE = 1.18$ ,  $p < .002$ ), whereas in the color-sufficient condition, size is more likely to be mentioned with increasing color typicality ( $\beta = 3.09$ ,  $SE = 1.45$ ,  $p < .04$ ).<sup>28</sup>

We thus overall find moderate evidence for typicality effects in our dataset. Typicality effects are strong for those items that clearly display typicality differences between the modified and unmodified utterance, but much weaker for the remaining items. That the evidence for typicality effects is relatively scarce is no surprise: the stimuli were specifically designed to minimize effects of typicality. However, the fact that both ways of quantifying typicality predicted redundant adjective use in the expected direction suggests that with more power or with stimuli that exhibit greater typicality variation, these effects may show up more clearly.

## F Experiment 3a: typicality norms for Experiment 3

Analogous to the color typicality norms elicited for utterances in Exp. 1, we elicited typicality norms for utterances in Exp. 3. The elicited typicalities were used in the Bayesian Data Analysis reported in Section 5.3.

### F.0.1 Methods

**Participants** We recruited 240 participants over Amazon’s Mechanical Turk who were each paid \$0.50 for their participation.

**Procedure and materials** On each trial, participants saw one of the images used in Exp. 2 and were asked to answer the question “How typical is this for an  $X$ ?” on a continuous slider with endpoints labeled “very atypical” to “very typical.”  $X$  was a nominal referring expression. In contrast to Exp. 1a, where we only elicited typicality norms for utterance-object pairs where the object was in the extension of the utterance under a deterministic semantics (e.g., here *dalmatian*, *dog*, or *animal* for a dalmatian), in this norming study we also elicited norms for utterance-object pairs where that was not clearly the case (e.g., *a bear* for a bison, *a car* for an ambulance, or *a snack* for a lobster). However, we did not test all utterance-object combinations, which would have led to an explosion of conditions. Instead, we tested each target object with its three utterances

---

<sup>28</sup>Again, conducting this analysis on the entire dataset yielded only a marginal interaction of sufficient property and color typicality in the right direction ( $\beta = -1.10$ ,  $SE = .64$ ,  $p < .09$ ).

(e.g., the dalmatian was paired with *dalmatian*, *dog*, and *animal*; the pug was paired with *pug*, *dog*, and *animal*, etc.). That yielded a total of 108 combinations – four targets in nine domains with three utterances each. We further tested each distractor item that shared the target’s superclass category (*dist-samesuper*, e.g., cows share the superclass category *animal* with dogs) on both the basic level and the super level term (e.g., *dog* for cow and *animal* for cow), for a total of 469 combinations. Finally, we also tested each distractor of a different super category than the target on the target’s super level term (*dist-diffsuper*, e.g., *animal* for socks). This yielded another 168 combinations. Overall, we obtained typicality norms for 745 object-utterance combinations. All other object-utterance combinations were assumed to have typicality 0.

Each participant rated 45 items: 7 targets, 10 dist-diffsuper, and 28 dist-samesuper cases. These were randomly sampled from the overall pool of items in each category.

#### F.0.2 Results and discussion

Each combination was rated at least 5 times and at most 27 times. We coded the slider endpoints as 0 (“very atypical”) and 1 (“very typical”). In order to evaluate the model, we used each object-utterance combination’s typicality mean as input.

Typicality ratings by item type (target, dist-samesuper, dist-diffsuper) and utterance type (sub, basic, super) are visualized in Figure 27. As expected, typicality was close to 0 for dist-diffsuper cases and for sub/basic terms used with dist-samesuper cases. However, even for these cases, there was some variation.

For targets, typicality of the object for the utterance decreased with increasing reference level, mirroring the typicality ratings obtained for Exp. 1 – a particular object is a better instance of the more specific term than of the more general term for that object.

## G Experiment 3 items

The following table lists all items used in Exp. 3 and the mean empirical utterance lengths that participants produced to refer to them:

## H Nominal choice model comparison

[jd: This isn’t model comparison in the technical sense, just a side-by-side look at the different models. Leave it in or throw out?]

Here we report correlations, MAP estimates of posterior predictives collapsed across targets and items, and scatterplots of posterior predictive MAP estimates on the by-target level for the model containing a) only informativeness with deterministic semantics; b) informativeness with deterministic semantics and cost; c) only informativeness with continuous semantics; d) informativeness with continuous semantics and cost (the model reported in the main text). Table 12 shows correlations. Figure 28 shows the collapsed patterns for utterance choice. Figure 29 shows the scatterplots.

Table 11: List of domains and associated superordinate category, target stimuli, and mean length (standard deviation) in characters of actually produced subordinate level utterances in Exp. 2.

| Domain | Super     | Targets         | Mean sub length (sd) |
|--------|-----------|-----------------|----------------------|
| bear   | animal    | black bear      | 9.9 (.14)            |
|        |           | polar bear      | 8.8 (.35)            |
|        |           | panda bear      | 5.5 (.2)             |
|        |           | grizzly bear    | 9 (.98)              |
| bird   | animal    | eagle           | 4.9 (.1)             |
|        |           | parrot          | 6.1 (.13)            |
|        |           | pigeon          | 5.9 (.22)            |
|        |           | hummingbird     | 10.1 (.5)            |
| candy  | snack     | MnMs            | 4.4 (.49)            |
|        |           | skittles        | 6.9 (.43)            |
|        |           | gummy bears     | 8.5 (.47)            |
|        |           | jelly beans     | 9.3 (.44)            |
| car    | vehicle   | SUV             | 3 (0)                |
|        |           | minivan         | 5.7 (.27)            |
|        |           | sports car      | 9.8 (.23)            |
|        |           | convertible     | 11.1 (.2)            |
| dog    | animal    | pug             | 3 (.08)              |
|        |           | husky           | 4.7 (.22)            |
|        |           | dalmatian       | 8.8 (.18)            |
|        |           | German Shepherd | 13.1 (.82)           |
| fish   | animal    | catfish         | 6.6 (.4)             |
|        |           | goldfish        | 7.9 (.22)            |
|        |           | swordfish       | 8 (.43)              |
|        |           | clownfish       | 9.1 (.38)            |
| flower | plant     | rose            | 4 (0)                |
|        |           | tulip           | 4.4 (.18)            |
|        |           | daisy           | 5.9 (.55)            |
|        |           | sunflower       | 9 (.11)              |
| shirt  | clothing  | T-shirt         | 6.4 (.48)            |
|        |           | polo shirt      | 6.7 (.79)            |
|        |           | dress shirt     | 11 (0)               |
|        |           | Hawaii shirt    | 12.6 (.46)           |
| table  | furniture | picnic table    | 9.7 (.58)            |
|        |           | dining table    | 12 (0)               |
|        |           | coffee table    | 9.1 (.95)            |
|        |           | bedside table   | 8.3 (.68)            |

Table 12: Correlations ( $r$  and  $R^2$ ) of posterior predictive MAPs of four different models (see main text) with empirical proportions of sub, basic, and super level choices.

|                |           | Model         |               |                   |                   |
|----------------|-----------|---------------|---------------|-------------------|-------------------|
|                |           | deterministic | deterministic | non-deterministic | non-deterministic |
| Semantics Cost |           | no            | yes           | no                | yes               |
| $r$            | collapsed | .85           | .88           | .86               | .94               |
|                | by-target | .63           | .71           | .71               | .84               |
| $R^2$          | collapsed | .72           | .77           | .74               | .89               |
|                | by-target | .40           | .51           | .51               | .70               |

## References

- Arts, A., Maes, A., Noordman, L., & Jansen, C. (2011). Overspecification facilitates object identification. *Journal of Pragmatics*, 43(1), 361–374. Retrieved from <http://dx.doi.org/10.1016/j.pragma.2010.07.013> doi: 10.1016/j.pragma.2010.07.013
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. doi: 10.18637/jss.v067.i01
- Baumann, P., Clark, B., & Kaufmann, S. (2014). Overspecification and the Cost of Pragmatic Reasoning about Referring Expressions. In *Proceedings of the 36th annual conference of the cognitive science society* (pp. 1898–1903). Austin, TX: Cognitive Science Society.
- Belke, E., & Meyer, A. S. (2002). Tracking the time course of multidimensional stimulus discrimination: Analyses of viewing patterns and processing times during same-different decisions. *European Journal of Cognitive Psychology*, 14(2), 237–266. doi: 10.1080/09541440143000050
- Bergen, L., Levy, R., & Goodman, N. (2016). Pragmatic reasoning through semantic inference. *Semantics and Pragmatics*, 9(1984), 1–46. Retrieved from <http://semprag.org/article/view/sp.9.20> doi: 10.3765/sp.9.20
- Bloomfield, L. (1933). *Language*. New York: Holt.
- Brennan, S. E., & Clark, H. H. (1996, nov). Conceptual pacts and lexical choice in conversation. *Journal of experimental psychology. Learning, memory, and cognition*, 22(6), 1482 – 1493. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/8921603>
- Brown, P., & Dell, G. (1987). Adapting Production to Comprehension : Mention of Instruments. *Cognitive Psychology*, 472, 441–472.
- Brown, R. (1958). Words and things.
- Brown-Schmidt, S., & Heller, D. (2014). What language processing can tell us about perspective taking: A reply to Bezuidenhout (2013). *Journal of Pragmatics*, 60, 279–284. Retrieved from <http://dx.doi.org/10.1016/j.pragma.2013.09.003> doi: 10.1016/j.pragma.2013.09.003
- Clark, H. H., & Murphy, G. L. (1982). Audience Design in Meaning and Reference. *Advances in Psychology*, 9(C), 287–299. doi: 10.1016/S0166-4115(09)60059-5
- Cohen, B., & Murphy, G. L. (1984). Models of concepts. *Cognitive science*, 8(1), 27–58.
- Dale, R. (1989). Cooking up referring expressions. *Proceedings of the 27th annual meeting on Association for Computational Linguistics (ACL'89)*, 68–75. Retrieved from <http://portal.acm.org/citation.cfm?doid=981623.981632> doi: 10.3115/981623.981632

- Dale, R., & Reiter, E. (1995). Computational Interpretations of the Gricean Maxims in the Generation of Referring Expressions . *Cognitive Science*, 19, 233 – 263.
- Davies, C., & Katsos, N. (2013). Are speakers and listeners 'only moderately Gricean'? An empirical response to Engelhardt et al. (2006). *Journal of Pragmatics*, 49(1), 78–106. Retrieved from <http://dx.doi.org/10.1016/j.pragma.2013.01.004> doi: 10.1016/j.pragma.2013.01.004
- Degen, J., Franke, M., & Jäger, G. (2013). Cost-based pragmatic inference about referential expressions. In *Cogsci*.
- Degen, J., Franke, M., & Jäger, G. (2013). Cost-Based Pragmatic Inference about Referential Expressions. In *Proceedings of the 35th annual conference of the cognitive science society*.
- Engelhardt, P. E., Bailey, K., & Ferreira, F. (2006a, may). Do speakers and listeners observe the Gricean Maxim of Quantity? *Journal of Memory and Language*, 54(4), 554–573. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0749596X05001518> doi: 10.1016/j.jml.2005.12.009
- Engelhardt, P. E., Bailey, K., & Ferreira, F. (2006b, may). Do speakers and listeners observe the Gricean Maxim of Quantity? *Journal of Memory and Language*, 54(4), 554–573. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0749596X05001518> doi: 10.1016/j.jml.2005.12.009
- Engelhardt, P. E., Demiral, S. B., & Ferreira, F. (2011). Over-specified referring expressions impair comprehension: An ERP study. *Brain and Cognition*, 77(2), 304–314. doi: 10.1016/j.bandc.2011.07.004
- Frank, A., & Jaeger, T. F. (2008). Speaking rationally: Uniform information density as an optimal strategy for language production. In *The 30th annual meeting of the cognitive science society*.
- Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, 336, 998.
- Franke, M., & Degen, J. (2016). Reasoning in Reference Games : Individual- vs . Population-Level Probabilistic Modeling. *PLoS ONE*, 11(5), 1–25. doi: 10.1371/journal.pone.0154854
- Gatt, A., Krahmer, E., van Deemter, K., & van Gompel, R. P. (2014). Models and empirical data for the production of referring expressions. *Language, Cognition and Neuroscience*, 29(8), 899–911.
- Gatt, A., van Gompel, R. P. G., Krahmer, E., & van Deemter, K. (2011). Non-deterministic attribute selection in reference production. In *Proceedings of the workshop on production of referring expressions: Bridging the gap between empirical, computational and psycholinguistic approaches to reference (pre-cogsci11)*. Boston. Retrieved from E:\$\backslash\$backslash\$Disser\$\backslash\$backslash\$Bibliography\$\backslash\$backslash\$gatt2011non.pdf
- Gatt, A., van Gompel, R. P. G., van Deemter, K., & Krahmer, E. (2013). Are we Bayesian referring expression generators? In *Proceedings of the workshop on production of referring expressions: Bridging the gap between computational and cognitive approaches to reference (pre-cogsci'13)*.
- Goodman, N. D., & Frank, M. C. (2016). Pragmatic Language Interpretation as Probabilistic Inference. *Trends in Cognitive Sciences*, 20(11), 818–829. Retrieved from <http://dx.doi.org/10.1016/j.tics.2016.08.005> doi: 10.1016/j.tics.2016.08.005
- Goodman, N. D., & Stuhlmüller, A. (2013, jan). Knowledge and implicature: modeling language understanding as social cognition. *Topics in Cognitive Science*, 5(1), 173–84. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/23335578> doi: 10.1111/tops.12007
- Goodman, N. D., & Stuhlmüller, A. (electronic). *The design and implementation of probabilistic programming languages*. Retrieved 2015/1/16, from <http://dippl.org>

- Graf, C., Degen, J., Hawkins, R. X. D., & Goodman, N. D. (2016). Animal , dog , or dalmatian ? Level of abstraction in nominal referring expressions. In A. Papafragou, D. Grodner, D. Mirman, & J. Trueswell (Eds.), *Proceedings of the 38th annual conference of the cognitive science society* (pp. 2261–2266). Austin, TX: Cognitive Science Society.
- Grice, H. P. (1975). Logic and Conversation. *Syntax and Semantics*, 3, 41–58. Retrieved from <http://books.google.com/books?hl=en&lr={\&}amp;id=hQCz0maGeVYC{\&}amp;oi=fnd{\&}amp;pg=PA121{\&}amp;dq=Logic+and+conversation{\&}amp;ots=j7aijUymwm{\&}amp;sig=iV1rz1eEm4ns6bQ6CevIURXFV04>
- Hawkins, R. X. D. (2015). Conducting real-time multiplayer experiments on the web. *Behavior Research Methods*, 47(4), 966-976.
- Hawkins, R. X. D., Stuhlmüller, A., Degen, J., & Goodman, N. D. (2015). Why do you ask ? Good questions provoke informative answers . In *Proceedings of the 37th annual conference of the cognitive science society*.
- Herrmann, T., & Deutsch, W. (1976). *Psychologie der Objektbenennung*. Huber.
- Hoffmann, J., & Ziessler, C. (1983). Objektidentifikation in künstlichen begriffshierarchien. *Zeitschrift für Psychologie mit Zeitschrift für angewandte Psychologie*.
- Horton, W., & Keysar, B. (1996). When do speakers take into account common ground? *Cognition*, 59, 91–117.
- Huetting, F., & Altmann, G. T. M. (2011). Looking at anything that is green when hearing "frog": how object surface colour and stored object colour knowledge influence language-mediated overt attention. *Quarterly journal of experimental psychology* (2006), 64(1), 122–145. doi: 10.1080/17470218.2010.481474
- Jaeger, T. F. (2006). *Redundancy and Reduction in Spontaneous Speech* (Unpublished doctoral dissertation). Stanford University.
- Jaeger, T. F. (2010). Redundancy and reduction: speakers manage syntactic information density. *Cognitive Psychology*, 61(1), 23–62.
- Johnson, K. E., & Mervis, C. B. (1997). Effects of varying levels of expertise on the basic level of categorization. *Journal of Experimental Psychology: General*, 126(3), 248–277. Retrieved from <http://papers3://publication/uuid/F6C763F3-CAD3-479E-BFED-1EF941293840> doi: 10.1037/0096-3445.126.3.248
- Jolicoeur, P., Gluck, M. A., & Kosslyn, S. M. (1984). Pictures and names: Making the connection. *Cognitive Psychology*, 16(2), 243–275.
- Kamp, H., & Partee, B. (1995, nov). Prototype theory and compositionality. *Cognition*, 57(2), 129–91. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/8556840>
- Kao, J., Wu, J., Bergen, L., & Goodman, N. D. (2014). Nonliteral understanding of number words. *Proceedings of the National Academy of Sciences of the United States of America*, 111(33), 12002–12007.
- Kennedy, C., & McNally, L. (2005). Scale Structure, Degree Modification, and the Semantics of Gradable Predicates. *Language*, 81(2), 345–381. Retrieved from <http://muse.jhu.edu/content/crossref/journals/language/v081/81.2kennedy.pdf> doi: 10.1353/lan.2005.0071
- Koolen, R., Gatt, A., Goudbeek, M., & Krahmer, E. (2011). Factors causing overspecification in definite descriptions. *Journal of Pragmatics*, 43(13), 3231–3250. Retrieved from <http://dx.doi.org/10.1016/j.pragma.2011.06.008> doi: 10.1016/j.pragma.2011.06.008
- Koolen, R., Goudbeek, M., & Krahmer, E. (2013). The effect of scene variation on the redundant use of color in definite reference. *Cognitive Science*, 37(2), 395–411. doi: 10.1111/cogs.12019

- Levinson, S. C. (1983). Pragmatics (cambridge textbooks in linguistics).
- Levinson, S. C. (2000). *Presumptive Meanings - The Theory of Generalized Conversational Implicature*. MIT Press.
- Levy, R., & Jaeger, T. F. (2007). Speakers optimize information density through syntactic reduction. In B. Schröckopf, J. Platt, & T. Hoffman (Eds.), *Advances in neural information processing systems* (Vol. 19, pp. 849–856). Cambridge, MA: MIT Press. Retrieved from <http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Speakers+optimize+information+density+through+syntactic+reduction>
- Lockridge, C. B., & Brennan, S. E. (2002, sep). Addressees' needs influence speakers' early syntactic choices. *Psychonomic bulletin & review*, 9(3), 550–7. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/12412896>
- Luce, R. D. (1959). *Individual Choice Behavior: A Theoretical Analysis*. New York: Wiley.
- Maes, A., Arts, A., & Noordman, L. (2004). Reference Management in Instructive Discourse. *Discourse Processes: A Multidisciplinary Journal*, 37(2), 117–144. Retrieved from [http://proxy.lib.uiowa.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=eric&AN=EJ682763&site=ehost-live\\$backslash\\$http://www.leaonline.com](http://proxy.lib.uiowa.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=eric&AN=EJ682763&site=ehost-live$backslash$http://www.leaonline.com) doi: 10.1207/s15326950dp3702\_3
- Marr, D. (1982). *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. San Francisco: W.H. Freeman.
- Mitchell, M. (2013). Typicality and object reference. *Proceedings of the 35th ...*, 3062–3067. Retrieved from <http://csjarchive.cogsci.rpi.edu/Proceedings/2013/papers/0547/paper0547.pdf>
- Murphy, G. L., & Smith, E. E. (1982). Basic-level superiority in picture categorization. *Journal of verbal learning and verbal behavior*, 21(1), 1–20.
- Nadig, A. S., & Sedivy, J. C. (2002, jul). Evidence of Perspective-Taking Constraints in Children's On-Line Reference Resolution. *Psychological Science*, 13(4), 329–336. Retrieved from <http://pss.sagepub.com/lookup/doi/10.1111/j.0956-7976.2002.00460.x> doi: 10.1111/j.0956-7976.2002.00460.x
- Oldfield, R. C., & Wingfield, A. (1965). Response latencies in naming objects. *Quarterly Journal of Experimental Psychology*, 17(4), 273–281.
- Olson, D. R. (1970). Language and thought: Aspects of a cognitive theory of semantics. *Psychological review*, 77(4), 257.
- Paraboni, I., van Deemter, K., & Masthoff, J. (2007). Generating Referring Expressions: Making Referents Easy to Identify. *Computational Linguistics*, 33(2), 229–254. doi: 10.1162/coli.2007.33.2.229
- Pechmann, T. (1989). Incremental speech production and referential overspecification. *Linguistics*, 27(1), 89–110. doi: 10.1515/ling.1989.27.1.89
- R Core Team. (2017). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Rohde, H., Seyfarth, S., Clark, B., Jäger, G., & Kaufmann, S. (2012). Communicating with Cost-based Implicature: a Game-Theoretic Approach to Ambiguity. In *Proceedings of the 16th workshop on the semantics and pragmatics of dialogue* (pp. 107 – 116).
- Rosch, E. (1973, may). Natural categories. *Cognitive Psychology*, 4(3), 328–350. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/0010028573900170> doi: 10.1016/0010-0285(73)90017-0

- Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic objects in natural categories. *Cognitive Psychology*, 8(3), 382–439. doi: 10.1016/0010-0285(76)90013-X
- Rubio-Fernandez, P. (2016). How redundant are redundant color adjectives? An efficiency-based analysis of color overspecification. *Frontiers in Psychology*, 7(153). doi: 10.3389/fpsyg.2016.00153
- Scontras, G., Degen, J., & Goodman, N. D. (2017). Subjectivity Predicts AdjectiveOrdering Preferences. *Open Mind: Discoveries in Cognitive Science*, 1(1), 53–65. doi: 10.1162/opmi
- Sedivy, J. C. (2003, jan). Pragmatic versus form-based accounts of referential contrast: evidence for effects of informativity expectations. *Journal of psycholinguistic research*, 32(1), 3–23. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/12647560>
- Sproat, R., & Shih, C. (1991). The cross-linguistic distribution of adjective ordering restrictions. In *Interdisciplinary approaches to language* (pp. 565–593). Springer Netherlands.
- Tanaka, J. W., & Taylor, M. (1991a). Object categories and expertise: Is the basic level in the eye of the beholder? *Cognitive psychology*, 23(3), 457–482.
- Tanaka, J. W., & Taylor, M. (1991b). Object categories and expertise: Is the basic-level in the eye of the beholder? *Cognitive Psychology*, 23, 457–482. doi: 10.1016/0010-0285(91)90016-H
- 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(July), 1–12. Retrieved from <http://journal.frontiersin.org/Article/10.3389/fpsyg.2015.00935/abstract> doi: 10.3389/fpsyg.2015.00935
- Zadeh, L. A. (1965). Fuzzy sets. *Information and control*, 8(3), 338–353.

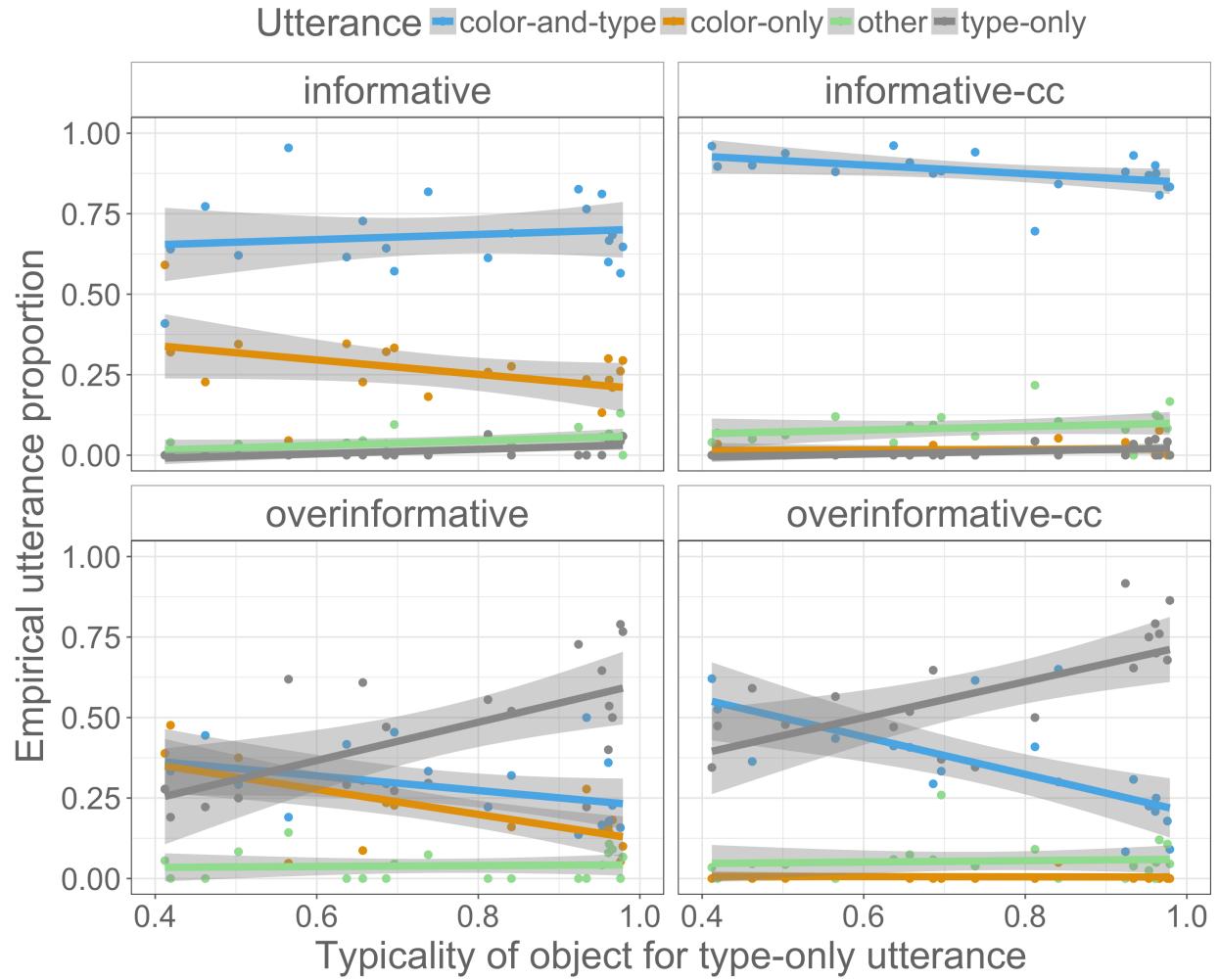


Figure 15: For each target, proportion of color-only (*yellow*), type-only (*banana*), color-and-type (*yellow banana*), and other (*funky carrot*) utterances as a function of mean object typicality for the type-only utterance (e.g., *banana*), across conditions. COLOR *banana* cases are circled in their respective color. [ek: judith, if this is ok, let me know and I will circle the banana cases.]

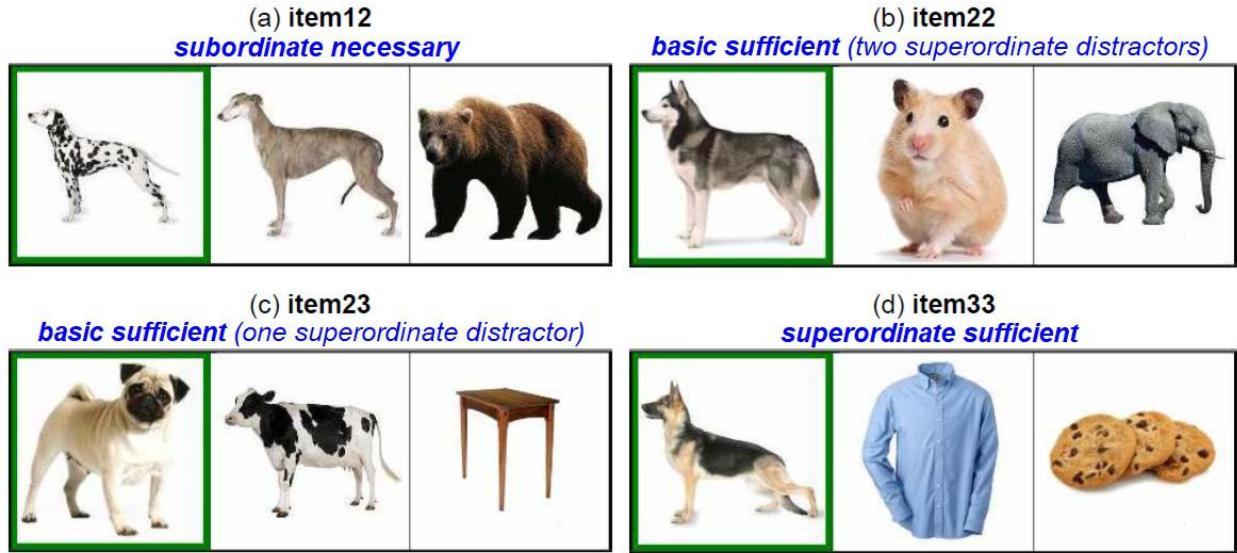


Figure 16: Example contexts in which different levels of reference are necessary for establishing unique reference to the target marked with a green border: (a) subordinate necessary (*dalmatian*); (b, c) basic sufficient (*dog*) and subordinate possible (*husky*, *pug*); (d) superordinate sufficient (*animal*) and basic or subordinate possible (*dog*, *German Shepherd*). [ek: I would make the item descriptions clearer (instead of item12,...) or leave them out][jd: let's get rid of the itemXX notation entirely – caroline, can you make sure that's the case in the basic level section, including this figure?]

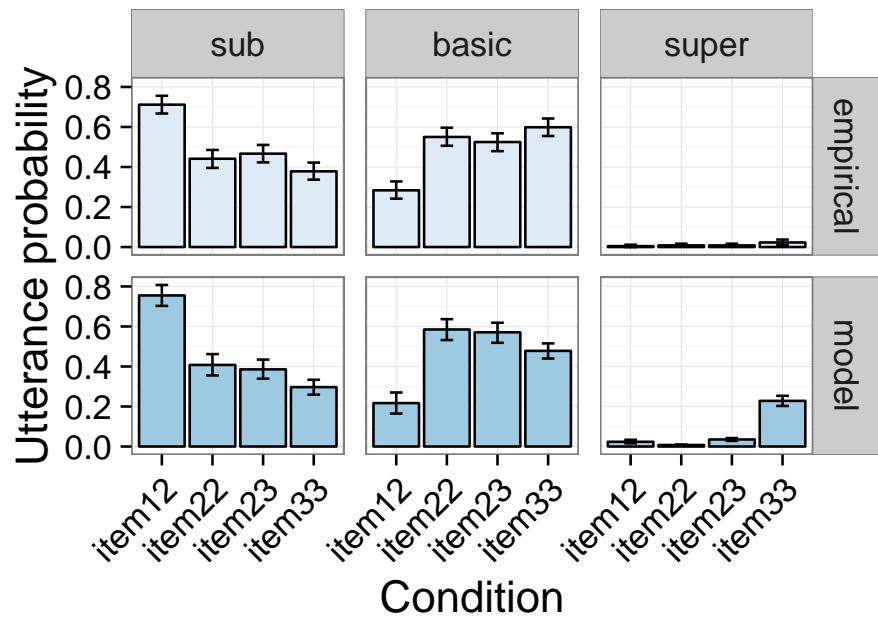


Figure 17: Utterance probabilities across different conditions. Columns indicate utterances, rows indicate data type (empirical proportion, MAP estimates of posterior predictives for full model with cost and continuous semantics).

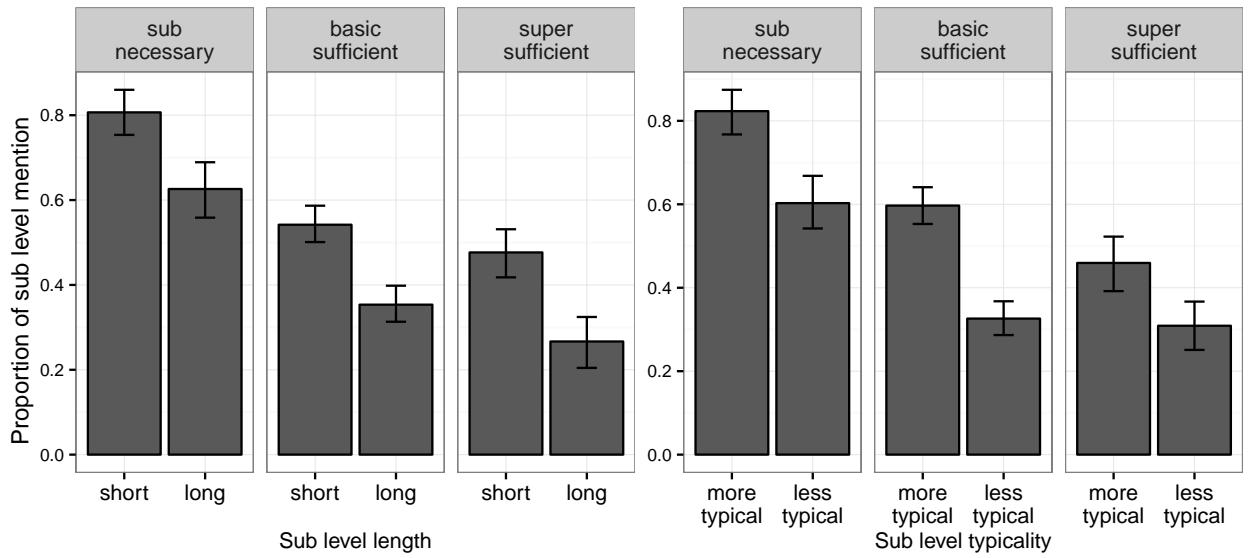


Figure 18: Proportion of sub level (over sub and basic level) terms across conditions. Left: when the sub length is relatively short (.67,1.82] or long [1.82,4.3) compared to the basic level term. Right: when the target object was relatively more [1.06,1.91) or less (.88,1.06] typical for the sub compared to the basic level term. Intervals were generated by splitting data into groups of roughly equal numbers of observations.

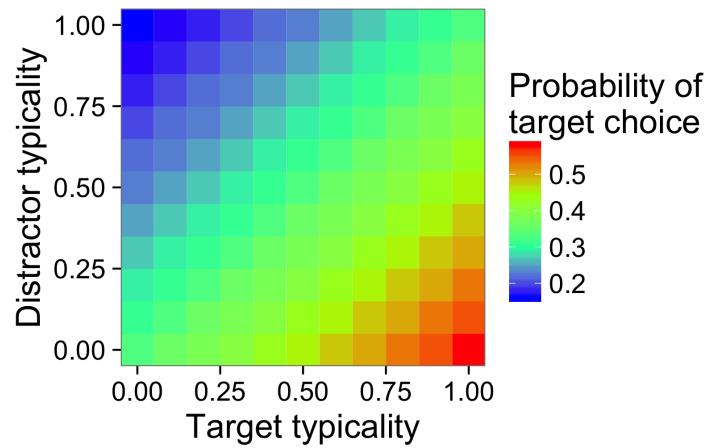


Figure 19: Literal listener probability of choosing the target under different typicalities of the target (x-axis) or the distractors (y-axis for the observed utterance. For simplicity we assume equal typicality of both distractors. The remaining probability mass for each case is thus uniformly distributed over both distractors.  
[jd: say where sub/basic/super level terms tend to fall in heatmap once you've regenerated plot]

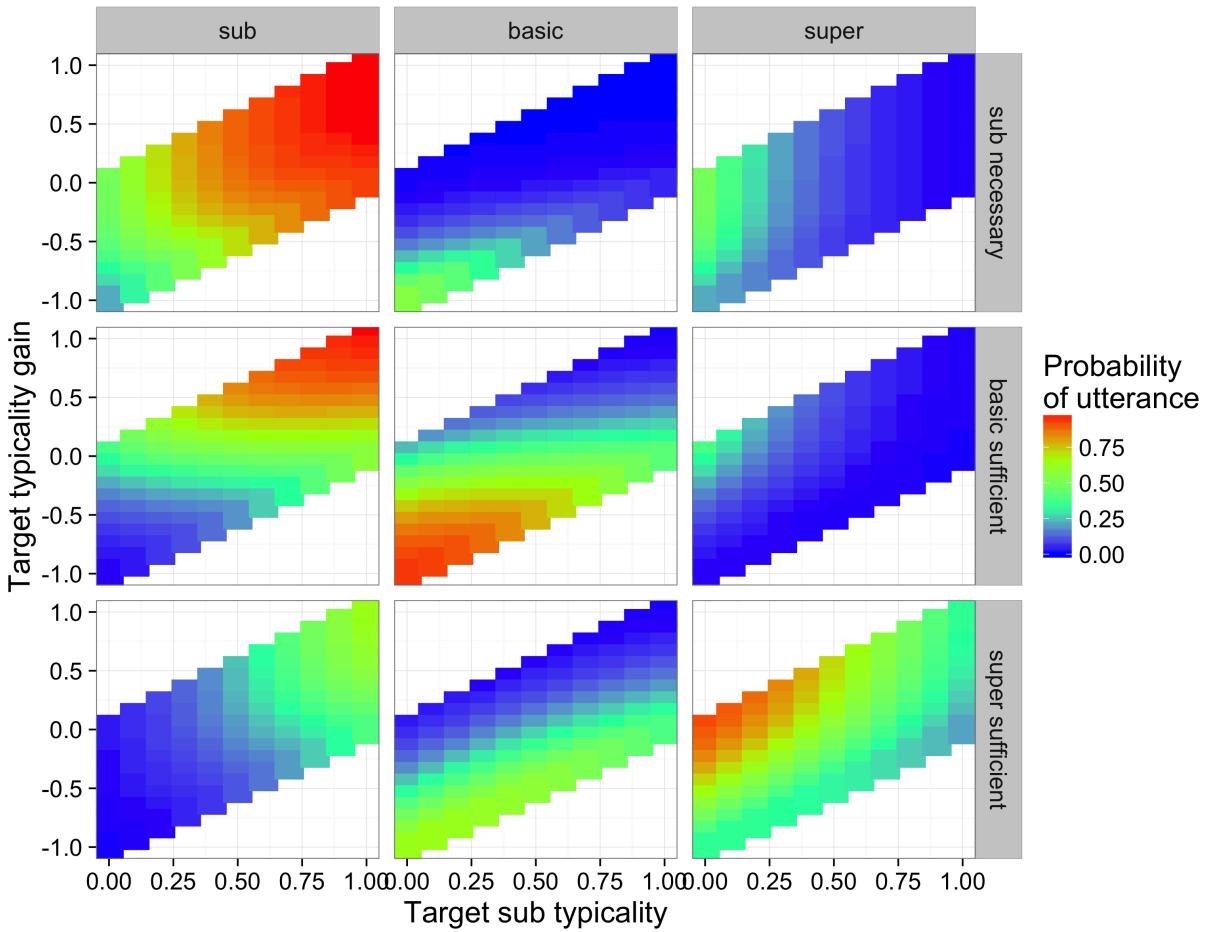


Figure 20: Pragmatic speaker probability of choosing each utterance (sub, basic, super) under varying absolute target sub typicalities (x-axis) and target typicality gains (difference between sub and basic level term typicality, y-axis), assuming equal typicality values for both distractors. Rows indicate different simulated conditions. [jd: go back into script to see how you computed typicality gain here]

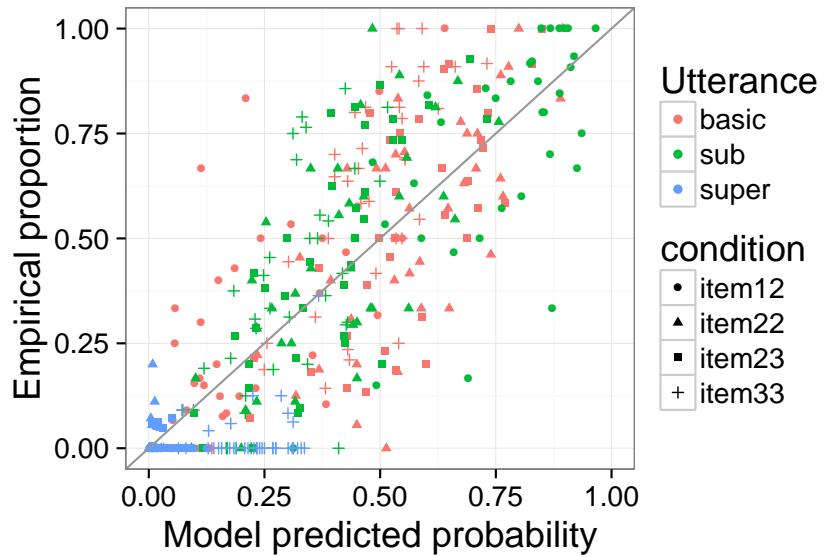


Figure 21: Scatterplot of by-target empirical utterance proportions against model posterior predictive MAP estimates. Gray line indicates perfect correlation line.

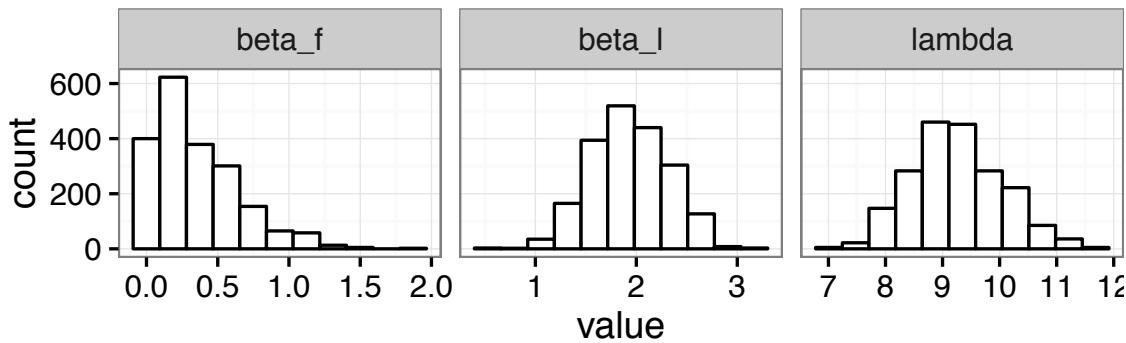


Figure 22: Posterior distribution over model parameters. Maximum a posteriori (MAP)  $\beta_f = 0.10$ , 95% highest density interval (HDI) = [0.002,0.95]; MAP  $\beta_l = 1.85$ , HDI = [1.23,2.65]; MAP  $\alpha = 9.19$ , HDI = [7.72,10.80].

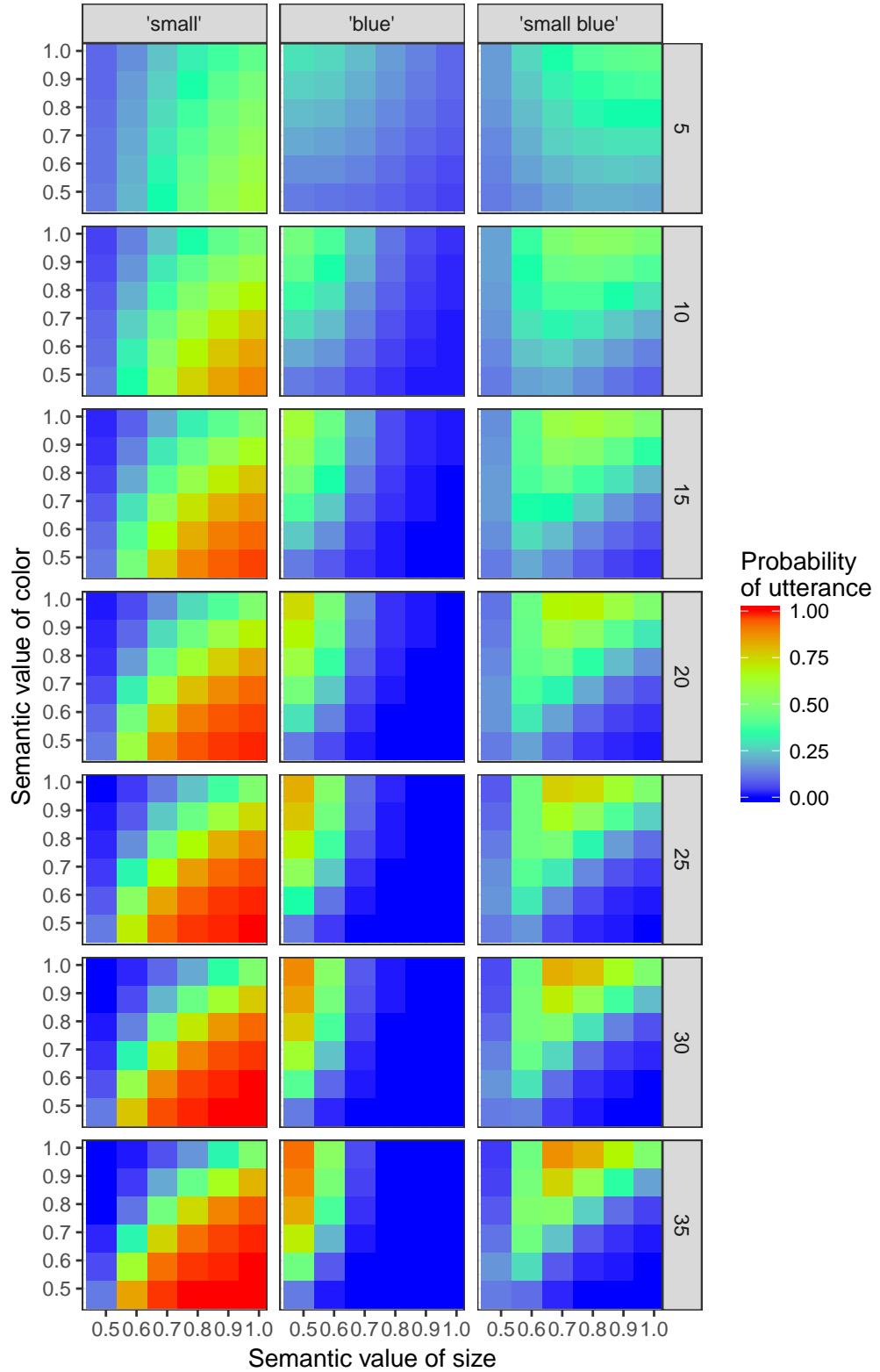


Figure 23: Probability of producing sufficient *small pin*, insufficient *blue pin*, and redundant *small blue pin* in contexts as depicted in Figure 1a, as a function of semantic value of color and size utterances and varying  $\alpha$  row-wise (for  $\beta_c = 0$ ). 63

How typical is this for a red stapler?



Figure 24: A modified example trial from the typicality elicitation experiment.

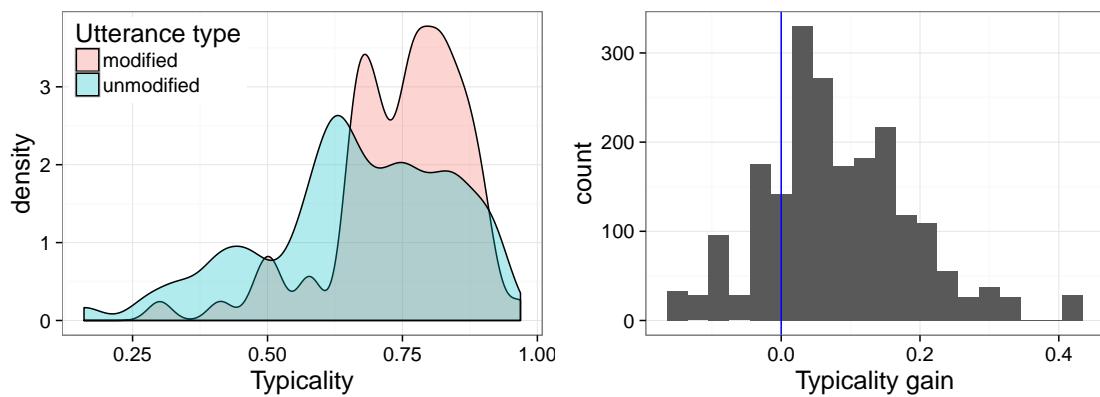


Figure 25: Typicality densities for modified and unmodified utterances (left) and histogram of typicality gains (differences between modified and unmodified typicalities, right).

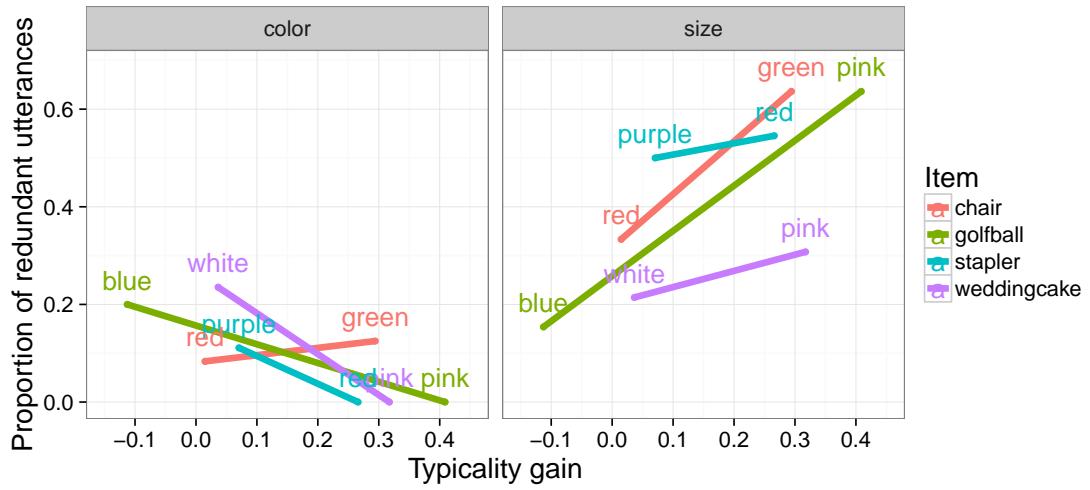


Figure 26: Utterance probability for four items as a function of difference in typicality between modified and unmodified utterance (x-axis) and sufficient dimension (columns).

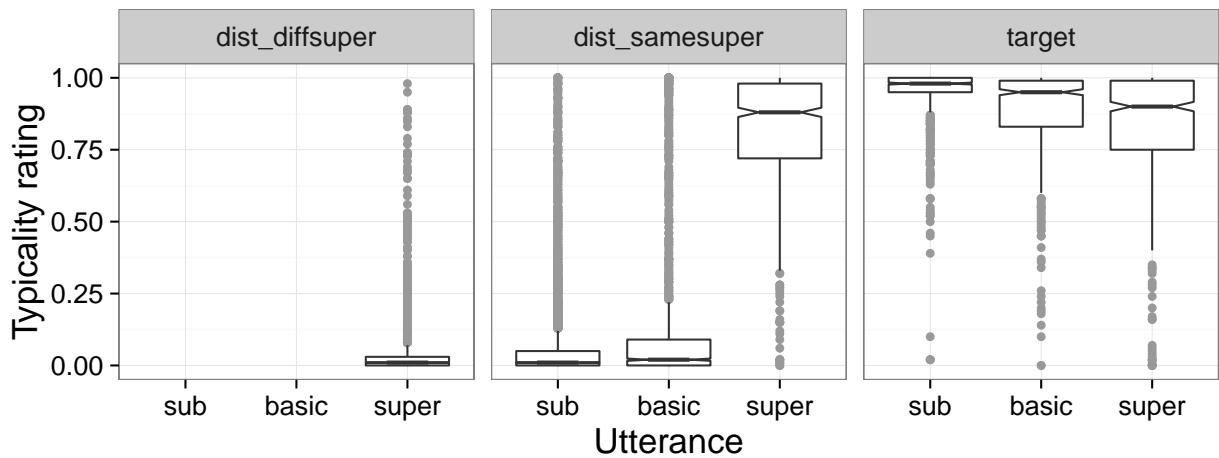


Figure 27: Boxplots of typicity ratings. The lower and upper hinges correspond to the first and third quartiles (the 25th and 75th percentiles). Upper and lower whiskers extend from the respective hinge to the highest and lowest values that are within 1.5 times the inter-quartile range of the hinge. Outliers are indicated as gray dots.

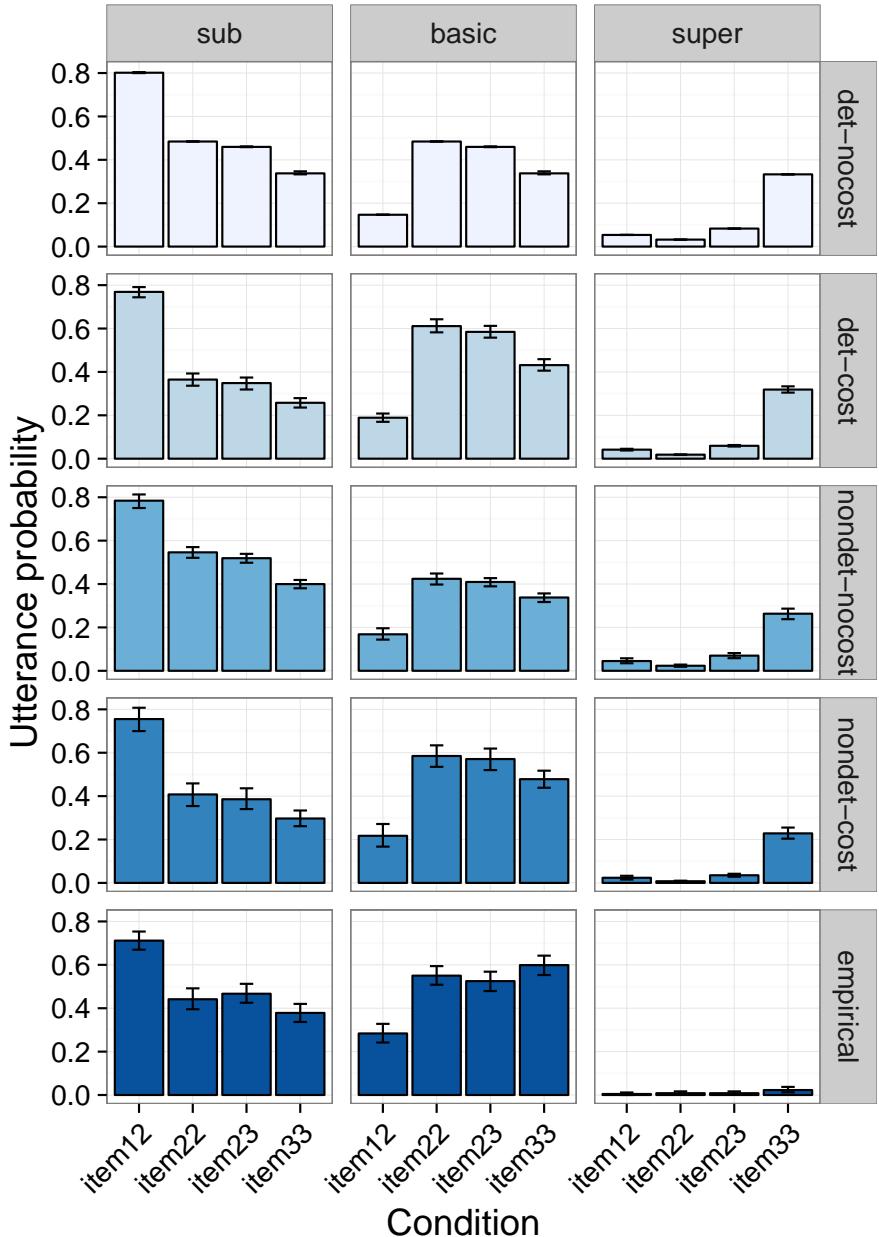


Figure 28: Utterance probabilities across different conditions. Columns indicate utterances, rows indicate data type (empirical proportion, MAP estimates of posterior predictives for the four different models).

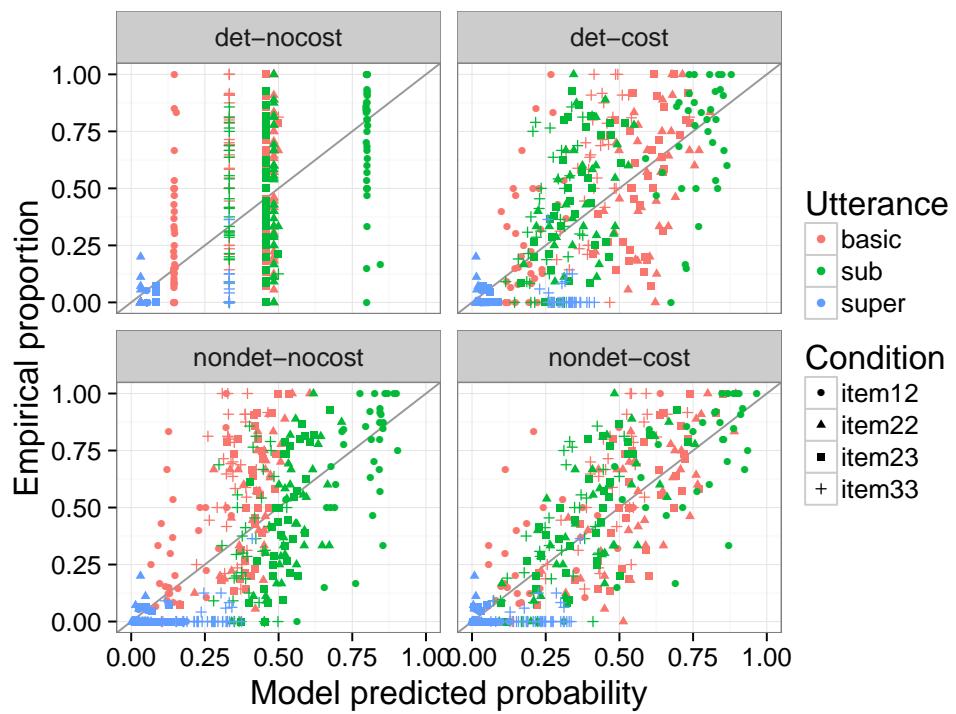


Figure 29: Scatterplot of by-target empirical utterance proportions against model posterior predictive MAP estimates for the four different models. Gray line indicates perfect correlation line.