

Communicative pressures at the semantics-pragmatics interface:

Learning biases may prevent the lexicalization of pragmatic inferences

(– draft October 1, 2016—)

Abstract

Certain classes of lexical meanings enable for pragmatic enrichments in a notably productive fashion. This raises the challenge to justify their regular selection over alternatives that codify semantically what is conveyed pragmatically. To address this issue, we propose a general model that integrates iterated Bayesian learning in the replicator-mutator dynamics. This model generates predictions about the effects of linguistic pressures on selection and transmission in populations of probabilistic language users with varied degrees of pragmatic sophistication. We apply this model in a case study on the (lack of) lexicalization of scalar implicatures. The results show that simpler semantic representations are selected for when languages are pressured towards learnability and compression, provided that pragmatic reasoning can compensate for the disadvantage in expressivity that users of such languages otherwise incur. We argue this result to shed light on the lack of lexicalization of scalar inferences, as well as the semantics-pragmatics distinction more generally.

1 Simplicity, expressivity, and learnability at the semantics-pragmatics interface

In linguistic theorizing it is common to draw a distinction between semantics and pragmatics. Broadly speaking, the former concerns the truth-conditional content of expressions whereas the latter concerns information beyond literal meanings and their composition. An important consequence of this distinction is that the information conveyed by an utterance is seldom, if ever, solely determined by semantics, but rather in tandem with pragmatics. Much research at the semantics-pragmatics interface has been aimed at characterizing expressions in terms of either domain or their interplay. However, an issue that has received little attention is the justification of semantic structure in light of pragmatics. The present investigation seeks to fill this gap by analyzing the effects linguistic pressures have on the selection and pervasiveness of particular lexical meanings under consideration of their possible pragmatic enrichments.

The emergence and change of linguistic structure is influenced by many intertwined factors. These range from biological and socio-ecological to cultural (Steels 2011, Tamariz and Kirby 2016). Social and ecological pressures determine communicative needs, while biology determines the architecture that enables and constrains the means by which they can be fulfilled. Our focus lies on the latter, cultural, factor wherein processes of linguistic change are understood as shaped by language use and its transmission. That is, as a result of a process of cultural evolution.

The idea that linguistic change is influenced by communicative pressures has played a pivotal role in synchronic and diachronic analyses at latest since Zipf's (1949) rationalization of the

approximability of word frequency rankings by a power law distribution as competing hearer and speaker preferences (e.g. Martinet 1962, Horn 1984, Jäger and van Rooij 2007, Jäger 2007, Piantadosi 2014, Kirby et al. 2015). In recent years this line of research has led to a surge in models that analyze communicative pressures within and across generations (see Steels 2015 and Tamariz and Kirby 2016 for recent overviews). Our starting point is given by the overarching argument that has crystalized from the accumulated mathematical, experimental and cross-linguistic evidence in this literature: Natural languages need to be well-adapted to communicative needs within a linguistic community, but also need to be learnable to survive their faithful transmission across generations. More succinctly; natural languages are pressured for expressivity and learnability.

The opposition of expressivity and learnability becomes particularly clear when considering their consequences in the extreme (cf. Kemp and Regier 2012, Kirby et al. 2015). On one end, a language with a single form is easy to learn but lacking in expressivity. On the other, a language that associates a distinct form with all possible meanings its users may want to convey is maximally expressive but challenging to acquire. The most prominent problem that arises from this tension is that of acquiring a language to express a potentially infinite set of meanings through finite means (Kirby 2002). However, this so-called transmission bottleneck is not the only challenge learners confront. A more pressing problem for our purpose is that of selecting particular hypotheses out of a (potentially infinite) space of alternatives compatible with the data learners are exposed to. At the semantics-pragmatics interface this concerns the selection between functionally similar, if not identical, lexical meanings. We assume an integral part of the answer to be that learners are a priori biased towards simpler, more compressed, lexical representations. That is, rational learners should prefer simpler over more complex explanations of the data they witness (Feldman 2000, Chater and Vitányi 2003, Piantadosi et al. 2012a, Kirby et al. 2015, Piantadosi et al. under review).

In the following, we model these pressures using the replicator-mutator dynamics by combining functional pressure on successful communication with effects of learning biases on (iterated) Bayesian learning (Griffiths and Kalish 2007). The semantics-pragmatics distinction and its effect on production and comprehension are made precise through probabilistic models of rational language use in populations and multiple lexica (Frank and Goodman 2012, Franke and Jäger 2014, Bergen et al. 2016). The remainder of this section introduces the individual components of the model as well as the assumptions underlying them. These are: the representation of languages and their use (§1.1), pressures towards expressivity (§1.2) and learnability (§1.3), regulated by the replicator and mutator dynamics, respectively, as well as a bias towards simpler semantic representations, codified as the learners’ prior. After laying out the model, we analyze its predictions in a case-study on the lack of lexicalization of scalar implicatures in §2.

1.1 Lexica and linguistic behavior

Lexica codify the truth-conditions of a language’s expressions, i.e., its semantics. A convenient way to represent them is by $(|S|, |M|)$ -Boolean matrices, where S is a set of states of affairs – or meanings – to convey and M a set of messages of the language (Franke and Jäger 2014). For example, the following two fragments determine the truth-conditions of two messages, m_1 and m_2 , for two states, s_1 and s_2 .

$$L_a = \begin{matrix} & \begin{matrix} m_1 & m_2 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \end{matrix} \qquad L_b = \begin{matrix} & \begin{matrix} m_1 & m_2 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \end{matrix}$$

According to lexicon L_a the message m_1 is true of state s_1 as well as of s_2 . In contrast, m_1

is only true of s_1 in L_b . Otherwise, the two languages are truth-conditionally equivalent.

To make the distinction between semantics and pragmatics precise, we distinguish between two kinds of linguistic behavior. *Literal interlocutors* produce and interpret messages literally. That is, their linguistic choices are guided only by their lexica. In contrast, *pragmatic interlocutors* engage in mutual reasoning to inform their choices. For instance, a rational speaker of L_a who reasons about her addressee may use m_1 to exclusively signal state s_1 given that s_2 can unambiguously be conveyed through m_2 . Accordingly, rational hearers that expect their interlocutor to reason along these lines will interpret ambiguous m_1 as s_1 . An important consequence of this pragmatic process is that it renders L_a indistinguishable from L_b in terms of expressivity.

Following models of rational language use such as Rational Speech Act models (Frank and Goodman 2012) and their game-theoretic counterparts (Benz et al. 2005a, Franke 2009, Franke and Jäger 2014), this kind of signaling behavior is captured by a reasoning hierarchy. The hierarchy’s bottom, level 0, corresponds to literal language use. Pragmatic language users of level $n + 1$ behave rationally according to (expected) level n behavior of their interlocutors. (1) and (3) specify the behavior of literal and pragmatic hearers of a language L . Mutatis mutandis for the literal and pragmatic speakers in (2) and (4).

$$H_0(s|m; L) \propto pr(s)L_{sm} \quad (1)$$

$$S_0(m|s; L) \propto \exp(\lambda L_{sm}) \quad (2)$$

$$H_{n+1}(s|m; L) \propto pr(s)S_n(m|s; L) \quad (3)$$

$$S_{n+1}(m|s; L) \propto \exp(\lambda H_n(s|m; L)^\alpha) \quad (4)$$

According to (1) a literal hearer’s interpretation of a message m as a state s depends on her lexicon and her prior over states, $pr \in \Delta(S)$. The literal speaker’s behavior as given in (2) is regulated by a soft-max parameter λ , $\lambda \geq 1$ (Luce 1959, Sutton and Barto 1998). As λ increases, choices made in production are more rational. That is, higher values lead to behavior that is increasingly in line with expected utility maximization.

Pragmatic behavior is similar to its literal counterpart. As in our informal sketch, their difference lies in that level $n + 1$ speakers/hearers reason about level n hearer/speaker behavior instead of solely relying on their lexicon. That is, they reason about how a rational level n interlocutor would use or interpret a message, and behave according to these expectations. Pragmatic production is further regulated by a parameter α which controls the tension between semantics and pragmatics, $\alpha \in (0, 1]$. Lower values lead to more literal production whereas higher values lead to stronger pragmatic behavior.

The combination of a lexicon with its use, i.e., a level in the reasoning hierarchy, is called a type. These are the basic units on which our population dynamics operate.

1.2 Replication & expressivity

Expressivity has received particular attention from investigations using evolutionary game theory (Nowak and Krakauer 1999, Nowak et al. 2000; 2002). Under this view, a type’s ability to convey and interpret information successfully confers it a higher fitness, this measure being relative to the success of other types in the population. In the simplest models fitness directly translates into the proportion of types present in the subsequent population. As a consequence, fitter types flourish while unfit ones are driven out. This association of communicative success within a population with changes in the proportion of types present in it creates a feedback loop that pressures the population towards greater expressivity. The replicator equation gives us the means to make these dynamics precise.

The proportion of types in a given population is codified in a vector x , where x_i is type i 's proportion. The fitness of a type i , f_i , is given by its expected utility in this population, $f_i = \sum_j x_j \text{EU}(t_i, t_j)$. That is, its fitness is the sum of its expected communicative success in interacting with other types weighted by the latter's population share. The expected utility of i and j is obtained by considering the average communicative success of i conveying information to j and vice versa: $\text{EU}(t_i, t_j) = [U_S(t_i, t_j) + U_R(t_i, t_j)]/2$. $U_S(x, y)$ and $U_R(x, y)$ are $\sum_s P(s) \sum_m S_n(m|s; L) \sum_{s'} R_o(s'|m; L) \delta(s, s')$ and $U_S(y, x)$, respectively, for n and o being the reasoning level of x and y , and $\delta(s, s') = 1$ iff $s = s'$ and 0 otherwise. This quantity is symmetric, reflecting the probability of two types' mutual understanding. The average fitness of the population is given by Φ , $\Phi = \sum_i x_i f_i$. This term serves as a normalizing constant for the (discrete) replicator equation: $\dot{x}_i = \frac{x_i f_i}{\Phi}$.

Under its biological interpretation the replicator equation captures the idea of fitness-relative selection whereby fitter types produce more offspring, leading to their propagation in subsequent generations. Similarly, many aspects of natural language are subject to processes of transmission and change across varied time-spans. Amongst others, the replicator equation can be understood as a model of language acquisition across generations, as e.g. in Nowak et al. 2002, but also as a process of horizontal adaptation (see Benz et al. 2005b:§3.3 for discussion). In the following, we take the latter view in assuming that interlocutors may adapt their types to that which works best within their population.

In their series of papers on language evolution, Nowak and colleagues did not only consider expressivity but also recognized the central role that learnability plays in its transmission. On the one hand, linguistic production can be prone to errors (Nowak and Krakauer 1999). On the other, multiple languages may be compatible with the data learners are exposed to (Nowak et al. 2002). Both sources introduce variation in language acquisition. In keeping the analogy to evolutionary processes this variation can be likened to mutation to the effect that a type's offspring may adopt a different type than that of its parent. Importantly, the resulting generational turnovers should depend on the relative learnability of a type. For this purpose, we turn to a different strand of research in cultural evolution: *iterated learning*.

1.3 Mutation & learning

Iterated learning is a process in which the behavior of one individual serves as learning input for another, who's behavior subsequently serves as input for a new learner, and so on. In language this process can be thought of as a progression through chains of parents and children where the parent produces linguistic data from which the child infers a language. The latter, now a parent, goes on to produce linguistic data for a new generation of naïve learners. Following Griffiths and Kalish (2007) we model iterated learning as a process of Bayesian inference in which learners combine the likelihood of a type producing the learning data with prior inductive biases. They then select a type to adopt from the resulting posterior distribution.

Due to the pressure towards learnability it exerts, iterated learning generally leads to simpler and more regular languages (see Kirby et al. 2014 and Tamariz and Kirby 2016 for recent surveys). Importantly, experimental and mathematical results suggest the outcome of this process to reflect the learners' a priori biases. In a Bayesian setting these biases are codified in the prior $P \in \Delta(\mathcal{T})$. A way to think about this prior is as the amount of data a learner would require in order to adopt a language (cf. Griffiths and Kalish 2007:450). Or, in our case, a combination of a lexicon and a signaling behavior. **INTRODUCE SIMPLICITY BIAS HERE**

In linguistics, a drive for simplicity has been argued to underpin speaker preferences for brevity and ease of articulation, as well as to pressure languages towards lexical ambiguity and grammatical compression (Zipf 1949, Grice 1975, Piantadosi et al. 2012b, Kirby et al. 2015). As

a broader cognitive principle, the use of simplicity as means to select between hypotheses has a long standing tradition. Crucially, Chater and Vitányi (2003) give a number of compelling arguments for simplicity on both mathematical and empirical grounds.

The extent of the prior’s influence has been shown to strongly depend on the learning strategy assumed to underly the inference process. While simulation results suggested that weak biases could be magnified by exposing learners to only small data samples (Brighton 2002), the mathematical characterization provided by Griffiths and Kalish (2007) showed that, instead, iterated learning converged to the prior. That is, the distribution over languages in a population, or the likelihood of learning a language when taking an individual’s perspective, corresponds to the learners’ prior distribution irrespective of the amount of input given to learners. This divergence in predictions can be traced back to differences in the selection of hypotheses from the posterior. On the one extreme, Griffith & Kalish’s prediction of convergence to the prior holds for learners that sample from the posterior. On the other, more deterministic strategies such as the selection of the type with the highest posterior probability, so-called *maximum a posterior estimation* (MAP), increase the prior’s influence (Griffiths and Kalish 2007, Kirby et al. 2007). In the following, we parametrize the posterior, $P(t_i|d)^l$, to obtain a range of learning strategies that live between posterior sampling and MAP, $l \geq 1$. When $l = 1$ learners sample from the posterior. As l increases towards infinity, the learners’ propensity to maximize the posterior increases.

We combine the replicator dynamics with iterated learning by codifying the latter as a transition matrix Q . Just as in standard mutator dynamics, Q_{ij} indicates the probability of the children of a parent of type i adopting type j . However, to make this process depend on a type’s learnability, this quantity is proportional to the probability of i producing the learning data and that of inferring j given the data.

The elements of the set of learning data D are k -length sequences of state-message pairings. That is, a sequence of observations of language use. Put differently, a datum $d \in D$ contains k members of the set $\{\langle s_i, m_j \rangle \mid s_i \in S, m_j \in M\}$. D is the set of all such sequences. With the at hand D , the learning matrix is constructed as follows:

$$Q_{ij} \propto \sum_{d \in D} P(d|t_i) F(t_j, d),$$

where $F(t_j, d) \propto P(t_j|d)^l$ and $P(t_j|d) \propto P(t_j)P(d|t_j)$. Given a type i , $P(d|t_i)$ can be straightforwardly computed based on t_i ’s production behavior.

1.4 Summary

We argued that expressivity, learnability and simplicity are central pressures that apply on the cultural evolution of language. They are respectively modelled as communicative efficiency-relative replication, iterated Bayesian learning, and a prior that biases learners for compressed lexical meanings. Taken together their dynamics can be described by the replicator-mutator dynamics (Hofbauer and Sigmund 2003):

$$\hat{x}_i = \sum_j Q_{ji} \frac{x_j f_j}{\Phi}$$

The units that the dynamics operate on are a combination of a lexicon and linguistic behavior which determines its use, i.e., a type. A type’s expressivity depends on its communicative efficiency within a population while its learnability depends on the fidelity by which it is inferred by new generations of naïve learners.

2 Scalar implicatures

A particularly well-studied type of conventional pragmatic enrichments are so-called *scalar implicatures*. These inferences are licensed for groups of expressions ordered in terms of informativity, here understood as an entailment induced order. For instance, *some* is entailed by *all*. If it were true that ‘All students came to class’, it would also be true that ‘Some students came to class’. However, while weaker expressions such as *some* are truth-conditionally compatible with stronger alternatives such as *all*, this is not what their use is normally taken to convey. Instead, the use of a less informative expression when a more informative one could have been used can license a defeasible inference that stronger alternatives do not hold (cf. Horn 1972, Gazdar 1979). That is, a hearer who assumes the speaker to be able and willing to provide all relevant information can infer that, since the speaker did not use a stronger alternative, e.g. *all*, this alternative must not hold. In this way, ‘Some students came to class’ is strengthened to convey ‘Some but not all students came to class’. Conversely, speakers can rely on their interlocutors to draw this inference. In a nutshell, mutual reasoning about rational language use supplies a bound that rules out stronger alternatives pragmatically which is absent from the semantics of weak scalar expressions. **This entire par needs to be more reader friendly**

This kind of strengthening corresponds to our previous description of the pragmatic use of lexicon L_a . A pragmatic hearer who reasons about a speaker’s use of message m_1 will associate it more strongly with s_1 than with s_2 since the latter is already unambiguously associated with s_2 . Conversely, a pragmatic speaker will reason about her interlocutor’s interpretation and use the messages accordingly.

$$L_a = \begin{matrix} & m_1 & m_2 \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \end{matrix} \qquad L_b = \begin{matrix} & m_1 & m_2 \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \end{matrix}$$

Our initial question can now be narrowed in terms of scalar implicatures by asking for a justification for the lack of lexical upper-bounds in weak scalar alternatives. That is, why semantics such as those of message m_1 in L_a are regularly selected for over the alternative of codifying the bound semantically as in L_b . More poignantly, would it not serve language users better if weak(er) expressions such as *warm*, *or*, *some* and *big* were truth-conditionally incompatible with stronger alternatives such as *hot*, *and*, *all* and *huge*? This question is particularly striking considering the number of expressions that license such inferences across languages and thusly lack a lexicalized upper-bound (Horn 1972, Horn 1984:252-267, Traugott 2004, van der Auwera 2010).

We see two main explanations for the lack of upper-bounds in the lexical meaning of weak scalar expressions. The first is that their truth-conditional compatibility with stronger expressions endows them with broader applicability, allowing them to occur in contexts in which their upper-bounded reading is absent. This can happen when embedded in downward-entailing contexts, when the speaker is likely uncertain about whether the upper bounded reading is true, or when the distinction between an upper-bounded reading and the simple, only lower-bounded reading, is not relevant. For instance, if for all the speaker knows ‘Some students came’ but she doesn’t know whether ‘All came’ then the use of unbounded *some* succinctly conveys her uncertainty about the latter. This may suggest a functionalist argument for why upper-bounded meanings do not conventionalize: should contextual cues provide enough information to the hearer to identify whether a bound is intended to be conveyed pragmatically, then this is preferred over expressing it overtly through longer expressions. For example by saying *some but not all* explicitly. Importantly, although morphosyntactic disambiguation is dispreferred due to its relative length and complexity, it allows speakers to enforce an upper-bound and override

contextual cues that might otherwise mislead the hearer.

In a nutshell, this explanation posits that scalar implicatures fail to lexicalize because, all else being equal, speakers prefer to communicate as economically as possible and pragmatic reasoning enables them to do so. Compare this with a hypothetical language that lexicalizes two expressions for each English scalar expression; one with and one lacking an upper-bound. Along this functionalist explanation, we see four conditions that may pressure languages for English-like semantics over this alternative. First, contextual cues are strongly reliable. Second, morphosyntactic disambiguation is seldom necessary. Third, morphosyntactic disambiguation is only marginally dispreferred. Fourth, larger lexica are costly. Overall, neither condition is convincing in its role as a central explanatory device for such a wide-spread phenomenon. The first two put a heavy burden on the ability to retrieve contextual cues to a degree that seems unlikely to undercut the benefit of unambiguous communication. It is likely that human language users are very good at retrieving cues from contexts, but to stipulate that they are so good as to undercut the benefit of safe communication provided by our hypothetical alternative strikes us as too strong of an assumption. As for the third and fourth condition, these seem mostly like technical solutions without a proper empirical basis.

All the rest of this section needs to be tidied up Instead, in what follows we investigate the hypothesis that the lack of lexicalization of scalar inferences is driven by the advantage in compression that lexical meanings lacking an upper-bound have over those that explicitly codify it. Note however that we do not represent this contrast in compression between lexical meanings explicitly in lexica. Instead, the bias towards a lack of upper-bounds in weak scalar alternatives is directly encoded in the learners' prior over types.

In principle this difference could be made precise with an adequate representational language, e.g., through measures over representational complexity such as minimal description length. There is a growing effort to develop such empirically testable representational languages. For instance, the so-called language of thought has been put to test in various rational probabilistic models that show encouraging results (see e.g. Katz et al. 2008, Piantadosi et al. under review; 2012b and references therein). We think that our assumption is well-warranted as a working hypothesis and decide against such an enrichment given that the introduction of a larger framework would also require further assumptions and justifications.¹

In sum, while we do not want to argue that functionalist pressure may not play a role, we do see a clear benefit in exploring whether matters of learnability would not give us additional leverage.

2.1 Setup

We consider populations with two signaling behaviors, either literal or pragmatic. As mentioned earlier, the former correspond to level 0 reasoners who only take their lexica into consideration. In the following, pragmatic types correspond to level 1 reasoners since higher level reasoning is not required to derive scalar implicatures from the lexica fragments we consider here. The prior over states is assumed to be uniform for all behaviors.

The lexica are listed in Table 1. As in our previous examples, they are (2, 2)-Boolean matrices, i.e., $|M| = |S| = 2$. Lexica L_1 to L_3 are not optimal for communication because they assign the same meaning to all their messages. They were included to showcase the selection process for a larger hypothesis space. L_4 and L_5 are our target lexica, codifying upper-bounded semantics for message m_1 (the first column in a lexicon) and a lack thereof, respectively. Lastly, L_6 is similar to L_5 in that two messages are true of the same state but differs from it in assigning upper-bounded

¹T: We could possibly show the length difference of the lexical meanings with and without an upper-bound using a LOT grammar in the appendix. I don't know if it is worth it though.

$$L_1 = \begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix} \quad L_2 = \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix} \quad L_3 = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$L_4 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad L_5 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \quad L_6 = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

Table 1: Space of possible lexicon fragments considered.

semantics to m_1 . Combining a signaling behavior with one of these 6 lexica yields a total of 12 distinct types.²

We focus our analysis on the contrast between literal and pragmatic types using lexica L_4 and L_5 . Note in particular that a type that has conventionalized upper and lower bounds to realize a (quasi-)partition of the relevant semantic space, such as L_4 , will produce speaker behavior that is *almost* indistinguishable from that of a language with only the respective upper bounds, but with Gricean speakers, such as L_5 . Almost, because there may be slight differences between the probability with which speakers would (erroneously) use a semantically false description and the probability with which speakers would (erroneously) use a pragmatically suboptimal description. This further contrasts with literal L_5 , a type lacks means to convey an upper-bound with m_1 . Due to the possible marginal difference in signaling behavior between pragmatic L_4 and L_5 , the selection of one type over the other is expected to mainly depend on the learning bias. Things are less clear for literal L_5 contrasted with literal/pragmatic L_4 . The former has a learning advantage but is expected to fare worse in terms of communicative fitness in virtue of ambiguous m_1 .

As implicit in the discussion above, one may think of s_1 as a “some but not all”-state and s_2 for an “all”-state. The literal meaning of weak scalar expressions such as *some* then corresponds to a message true of both s_1 and s_2 in these fragments. Following our assumption for a preference for simple lexical representations, the prior biases learners against lexica in which a message holds true only of the former and not the latter. All other semantics are a priori equally probable. This prior is given by $P(t_i) \propto n - c \cdot r$, where n is the total number of states and r is the number of messages only true of s_1 in t_i ’s lexicon, $c \in [0, 1]$. In sum, an increase in c brings about a stronger learning bias against languages that lexicalize upper-bounds, i.e., L_2, L_4 and L_6 .

2.2 Analysis

The dynamics are initialized with an arbitrary distribution over types, constituting the population’s first generation. If not stated otherwise, the results for a given parameter setting were obtained from 1000 independent runs. Each run consisted of 20 generations. This corresponds to a developmental plateau after which no noteworthy change was registered. As specified in ??, the learning transition matrix Q can be obtained by considering all possible state-message sequences of length k . Given that this is intractable for large k , matrices with $k > 5$ were approximated by using Monte Carlo with 10 sequences sampled from each type’s production probabilities and a type’s children being exposed only to this subset. The model’s parameters are summarized in Table 2.

According to our hypothesis, functional pressure on successful communication combined with learning pressures in the form of a bias against upper-bounds may lead to the selection of L_5 -like

²While there is a total of 16 possible $(2, 2)$ -matrices, a number of them are identical both in terms of expressivity and the learning bias against lexical upper-bounds. The competition between such types is determined by their proportions in the initial population but this fact can be obscured when averaging across simulations. We therefore focus only on the subset that exhibits the properties we set out to explain.

parameter	explanation	locus
$\lambda \geq 1$	rationality parameter	$S_{n+1}(m s; L) \propto \exp(\lambda H_n(s m; L)^\alpha)$
$\alpha \in [0, 1)$	semantics-pragmatics tension	$S_{n+1}(m s; L) \propto \exp(\lambda H_n(s m; L)^\alpha)$
$ D $	learning data produced per parent type	$P(d t_j)P(t_i d)$
$k = d $	number of observations per datum	$P(d t_j)P(t_i d)$
$l \geq 1$	posterior parameter from sampling to MAP	$P(t_i d) \propto [P(t_i)P(d t_i)]^l$
$c \in [0, 1]$	learning bias for lack of upper-bounds	$P(t_i)$

Table 2: Summary of model parameters.

semantics. However, it is instructive to first inspect the effect of these pressures in isolation. For this purpose, we focus on three pragmatic types.³ Users of L_3 , representing a type that is lacking in expressivity but is a prior preferred for its lack of upper-bounds. Users of L_4 , a type that is functionally advantageous but biased against. And users of L_5 , combining the virtues of the latter two.

Expressivity only. Recall that the outcome of the replicator dynamics are influenced by λ and α as they have a bearing on a type’s fitness. Low α disadvantages types that rely on pragmatic reasoning for more deterministic signaling behavior to the gain of those that codify this information semantically. The rationality parameter λ has a similar effect for different reasons. Less utility maximizing behavior decreases the association of an ambiguous message with a single state, even when other states are uniquely associated with a different one. That is, λ regulates the strength by which users of L_5 associate non-upper-bounded m_1 exclusively with the “some”-state s_1 over the “all”-state s_2 .

The effect of the rationality parameter using only the replicator dynamics is shown in the left-hand side of Figure 1. As expected, the less expressive type using L_3 fares worse and shows little variation across λ -values. Crucially, values of $\lambda \leq 10$ lead to an increase in L_4 and a decrease in L_5 . As the rationality parameter increases, the functional difference between L_4 and L_5 is levelled. Overall, the final populations that result only from a pressure towards expressively approximate an even share of pragmatic L_4, L_5 and L_6 types. The latter follows the same trajectory as L_4 in Figure 1.

Learnability only. To get an impression of the effect of iterated learning without a pressure towards expressivity, we first consider (relatively) deterministic parental data production ($\lambda = 20, \alpha = 1$). In this way small data sequences suffice for learners to differentiate types that produce strongly diverging expressions. In line with Griffiths and Kalish’s (2007) analysis, under these conditions the resulting populations converge to the learner’s prior distribution when sampling from the posterior. This is shown in the right-hand side of Figure 1, which directly reflects the prior distribution for each value of c .

Inspecting the effects of these dynamics separately not only gives some intuitions about the parameters’ influence, but also highlights some of their broader implications. First and foremost, neither dynamic comes close to converging to a monomorphic population under most parameter configurations. For instance, types using L_4 can come to occupy a large proportion of the final population. However, this holds only for a restricted range of low degrees of rationality. Apart from polymorphy, both dynamics make some undesirable predictions. A pressure only towards expressivity leads to the selection of types using L_4 to L_6 and to the ejection of L_1 to

³Pragmatic reasoning allows language users to refine their (possibly erroneous) choices. Therefore, it is advantageous even for those types that codify an upper-bound lexically.

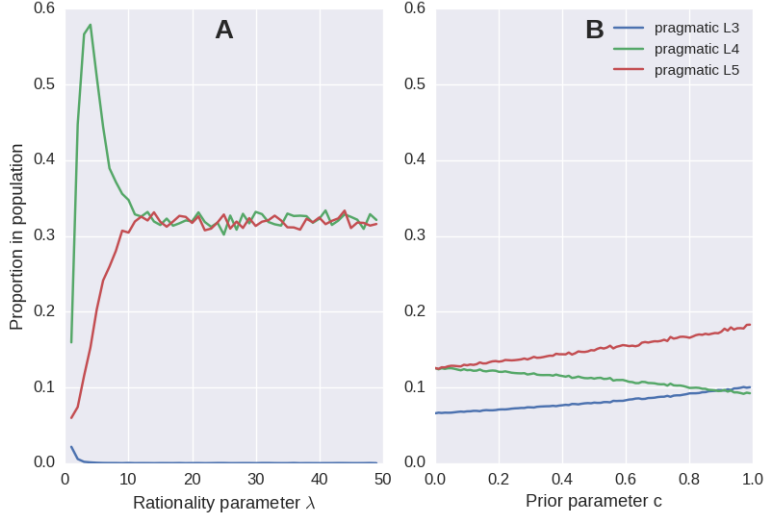


Figure 1: Mean proportions of target types after 20 generations in 1000 populations with only replication (A; $\alpha = 1$) and only mutation (B; $\alpha = 1, \lambda = 30, k = 5$).

L_3 . However, it can not explain the regular selection of L_5 -like semantics over either of these functionally similar alternatives. In contrast, a pressure only towards learnability has a modest but clear effect in differentiating L_5 from these alternatives but fails to rule out functionally suboptimal types such as tautological L_3 . This showcases that, while the bias plays a major role for the contrast between L_4 and L_5 , on its own it does not enable types that fail to convey an upper-bound to establish themselves in the population. In sum, neither dynamic on its own is a suitable candidate to provide a justification for the predicted prevalence of L_5 -like semantics.

Expressivity and learnability. Figure 2.2 illustrates the effect of the learning bias after 20 generations across values of c with $l = 1$ (A) and $l = 3$ (B). More detailed results for all types across a sample of c -values are presented in Table 3. Overall, these results suggest that in the present setup a weak bias is sufficient to lead to a selection of L_5 over L_4 . As in the simulations that only considered learnability, this effect increases with the bias’ strength provided L_5 users are pragmatic. Importantly, the addition of a pressure towards expressivity magnifies this effect by dampening the proliferation of functionally suboptimal types advantaged by the learning bias. As stressed above, this suggests that neither the learning bias nor functional pressure alone but their combination lead to the systematic lack of semantic upper-bounds in scalar expressions.

This data suggests that the proportion of scalar implicature users that is predicted primarily hinges on three aspects of the model. First, the degree to which linguistic behavior is deterministic, controlled by λ and α , plays a role both for expressivity as well as for producing data that allows learners to discriminate this type from others. Second, the learning bias c , leads learners to discriminate and prefer L_5 over L_4 and L_3 . Lastly, the posterior parameter l magnifies the effects of the learning bias in tandem with replication. This interaction is shown in 3. In the present setup posterior sampling can lead to the incumbency of pragmatic L_5 , but not even a strong favorable learning bias manages to completely drive out competing types (cf. 2.2.A). However, as posterior maximization increases, the range of bias values within which this type takes over the population increases drastically.

More discussion about parameter interaction: changes in sequence length influence the pop-

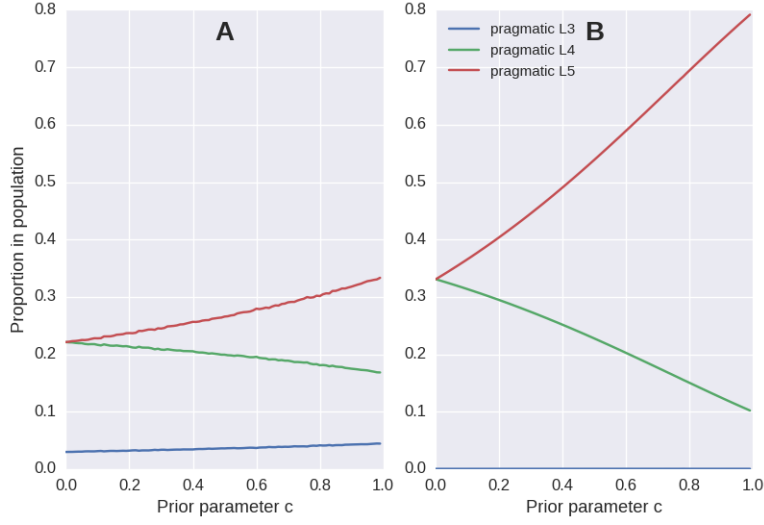


Figure 2: Mean proportions of target types after 20 generations in 1000 populations across bias values $c \in [0, 1]$ with $l = 1$ in A and $l = 3$ in B ($\alpha = 1, \lambda = 20, k = 5$).

$l = 1$						$l = 10$					
c	0	.01	.05	.1	.8	0	.01	.05	.1	.8	
lit. L_1	.03	.03	.03	.03	.04	€	€	€	€	€	
lit. L_2	.03	.03	.03	.03	.01	€	€	€	€	€	
lit. L_3	.03	.03	.03	.03	.04	€	€	€	€	€	
lit. L_4	.07	.07	.07	.07	.06	€	€	€	€	€	
lit. L_5	.04	.04	.05	.05	.06	€	€	€	€	€	
lit. L_6	.04	.04	.04	.04	.04	€	€	€	€	€	
prg. L_1	.03	.03	.03	.03	.04	€	€	€	€	€	
prg. L_2	.03	.03	.03	.03	.01	€	€	€	€	€	
prg. L_3	.03	.03	.03	.03	.04	€	€	€	€	€	
prg. L_4	.22	.22	.22	.22	.18	.33	.33	.32	.31	.15	
prg. L_5	.22	.22	.22	.23	.3	.33	.33	.35	.37	.7	
prg. L_6	.22	.22	.22	.22	.18	.33	.33	.32	.31	.15	

Table 3: Mean proportions of types in 1000 populations after 20 generations across bias values $c \in [0, 1]$ with $l = 1$ and $l = 3$ ($\alpha = 1, \lambda = 30, k = 5$), $\epsilon < 0.005$

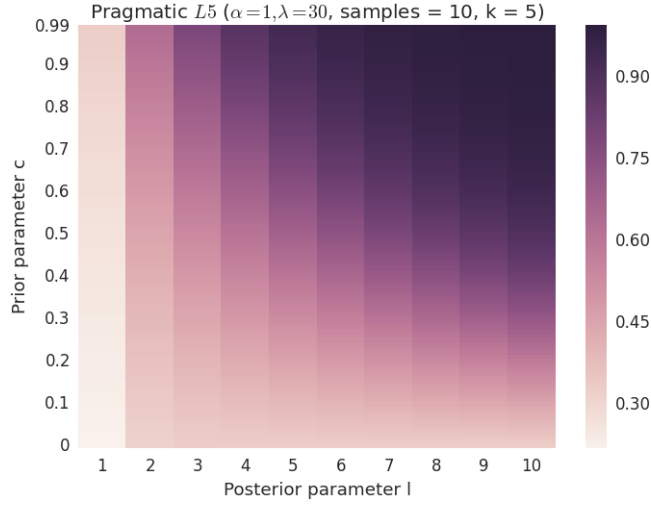


Figure 3: Mean proportion of pragmatic L_5 in 1000 populations after 20 generations ($\alpha = 1, \lambda = 30, k = 5$)

ulation in a predictable way: smaller values lead to more heterogeneous populations whereas larger ones lead to more pronounced differences. This is expected insofar as the likelihood that a sequence of length 1 was produced by any type is relatively uniform (modulo prior) whereas the likelihood of types with lexica $L_1 - L_3$ to produce, for instance, a sequence of 10 observations consistently with the same state-message combination is less likely than for pragmatic players using $L_4 - L_6$ or literal L_4 . Thus, while noteworthy, sequence length has no direct bearing on the main contrast of interest. Similar considerations hold for α and λ – set to 1 and 50 in the following. Overall, lower rationality in λ or more pragmatic violations in α lead to a higher selection of lexica with semantic upper-bounds. The fitness of pragmatic behavior increases with higher λ/α -values. In other words, these parameters level the functional contrast between L_4 and L_5 .

2.3 Discussion

Broadly speaking these results suggest that a lack of semantic upper-bounds coupled with pragmatic reasoning can overcome selective pressures and stabilize in a population provided there is a bias for simpler representations. This outcome is particularly encouraging in light of the other advantages a lack of semantic upper-bounds may confer that were discussed above.

The model predicts this result to hinge on four conditions. First, types need to be pressured toward both expressivity and learnability. Second, language use can not be too unpredictable; low α or λ values render non-upper-bounded lexical meanings more prone to communicative failure and more difficult to learn. Third, learners should prefer simpler over more complex lexical representations. Fourth, the strength by which they need to be preferred depends on the inferential mechanism of learners. That is, to which degree learners maximize the posterior. Put differently, high rationality in learning and choice requires a weaker bias towards simpler representations. The selection of lexical meanings lacking upper-bounds is robust against parameter perturbations under these conditions.

A large range of parameter values lead to the stable incumbency of scalar implicatures.

However, while this is predicted by the literature, it is less clear to what proportion other types should be represented. On the one hand, there are reasons to expect functionally suboptimal types L_1 - L_3 to be largely ruled out. They fail to enable their users to communicate without substantial error both amongst themselves and with other types. In short, these lexica stand at odds with past theoretical and empirical observations. On the other, this is not so for L_4 .⁴ Notwithstanding, it is possible that natural language communities are composed of such mixed populations or that a single speaker may entertain L_4 -like semantics for one scalar expression and L_5 -like semantics for another **CITATIONS FOR THIS PURPOSE**.

Lastly, there are a number of extensions one may want to consider. For instance, the consideration of larger lexica that consider epistemic uncertainty. Another interesting addition relates to possible disadvantages of pragmatic reasoning and the effect of state frequencies on fossilization. We tacitly assumed pragmatic reasoning to come at no cost. However, there is experimental evidence that suggests that the pragmatic derivation of upper-bounds costs effort and takes additional processing time (cf. Neys and Schaeken 2007, Huang and Snedeker 2009). This raises the question at which point such usage-based cost undercuts the learnability advantage of simpler semantic representations. As noted above, this might also depend on a given scalar expression and its frequency. Frequently drawn scalar implicatures might fossilize to avoid cost, while infrequent ones could still be computed on-line.

A virtue of this model is that it allows for analysis specific modifications and extensions. Straightforward extensions include larger hypothesis spaces, as well as larger or different lexicon fragments. Another worthwhile modification relates to possible disadvantages of pragmatic reasoning. We tacitly assumed pragmatic reasoning to come at no cost. However, there is experimental evidence for the assumption that the pragmatic derivation of upper-bounds costs effort and takes additional processing time (cf. Neys and Schaeken 2007, Huang and Snedeker 2009). This raises the question at which point such usage-based cost undercuts the learnability advantage of simpler semantic representations.⁵

3 General discussion

We laid out a model that combines game-theoretical models of functional pressure towards efficient communication (Nowak and Krakauer 1999), effects of learning biases on (iterated) language learning (Griffiths and Kalish 2007), probabilistic speaker and listener types of varied degrees of pragmatic sophistication (Frank and Goodman 2012, Franke and Jäger 2014) as well as different lexica (Bergen et al. 2012; 2016). This model generates predictions about lexicalization patterns found in natural language and the failure to lexicalize of certain pragmatic inferences. While the puzzle raised by semantics is hard to explain by purely functional means, it suggests part of the answer to lie in learnability. Simpler semantic representations are more likely to be learned, and pragmatic reasoning can counteract functional disadvantages otherwise incurred. This result is of particular relevance for the longstanding assumption of a divide and interaction between semantics and pragmatics by offering an account of why (certain) pragmatic inferences are not

⁴ L_6 presents a special case. In our current setup, it mirrors L_5 in allowing for pragmatic strengthening of message m_2 rather than m_1 . However, its association of s_1 with m_1 and s_2 with m_2 under favourable parameter conditions is achieved by ruling out the “some but not all”-state s_1 and not, as with scalar implicatures, the “all”-state s_2 . L_6 speakers thusly have a “some”-message strengthened to convey “all but not some but not all”. The observation that monomorphemic expression that lexically rule out stronger alternatives are unattested across languages has received substantial argumentative support (most prominently in Horn 1984:252-267 but also e.g. in Horn 1972, Traugott 2004, van der Auwera 2010). However, our present setup is blind to these differences.

⁵In the present setup this modification has a straightforward effect: A penalty for pragmatic signaling lowers the fitness of pragmatic types, to the advantage of literal types. However, the penalty needs to be substantial to counteract the functional advantage pragmatic L_5 has over all but L_4 together with its learning advantage.

part of the literal meaning of expressions. It furthermore leaves the possibility of such inferences to fossilize open when they do not compete against a lexical simplicity bias.

Discussion about expressivity as external to learning (cf. Stadler, replicator-papers by Kenny Smith, Kirby et al 2015). Possibly add appendix with direct comparison between IL and RMD.

In contrast to past research using the replicator-mutator dynamics, the Bayesian setting afforded by iterated learning allows for a straightforward integration of a learning prior. Furthermore, less idiosyncratic assumptions about the variation introduced by learning without having to resort to similarity matrices between vocabularies (Nowak and Krakauer 1999) nor languages (Nowak et al. 2002).

4 Conclusion

Language change is affected by complex intertwined pressures. Drawing from past insights we put forward a model that allows for a general integration of core aspects involved in this process. In particular, the model combines functional pressure, iterated Bayesian learning, and probabilistic speaker and hearer models. We argued this combination of cultural evolution dynamics with models of probabilistic language use to be suitable for the analysis of longstanding issues concerning the semantics-pragmatics divide.

In a case-study on the lack of lexicalization of scalar implicatures, the model predicts that, when pressured for learnability and expressivity, the former force drives for simpler semantic representations inasmuch as pragmatics can compensate for them in language use. Consequently, semantic patterns observed in natural languages can be explained in virtue of the linguistic behavior exhibited by their users and their representational complexity. In particular, the relative ease of acquisition of simpler semantics may offer an answer to why natural languages do not lexicalize certain pragmatic inferences.

References

- Luc Steels. *The Talking Heads experiment: Origins of words and meanings*. Language Science Press, 2015.
- Monica Tamariz and Simon Kirby. The cultural evolution of language. *Current Opinion in Psychology*, 8:37–43, 2016.
- Josef Hofbauer and Karl Sigmund. Evolutionary game dynamics. *Bulletin of the American Mathematical Society*, 40(04):479–520, 2003.
- Thomas L. Griffiths and Michael L. Kalish. Language evolution by iterated learning with bayesian agents. *Cognitive Science*, 31(3):441–480, 2007.
- M. C. Frank and N. D. Goodman. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012.
- Michael Franke and Gerhard Jäger. Pragmatic back-and-forth reasoning. *Semantics, Pragmatics and the Case of Scalar Implicatures.*, pages 170–200, 2014.
- Leon Bergen, Roger Levy, and Noah D Goodman. Pragmatic reasoning through semantic inference. *Semantics and Pragmatics*, 2016.
- Luc Steels. Modeling the cultural evolution of language. *Physics of Life Reviews*, 8(4):339–356, 2011.

- George Zipf. *Human behavior and the principle of least effort*. Addison-Wesley Press, 1949.
- André Martinet. *Functionalist View of Language*. Clarendon Press, Oxford, 1962.
- Laurence R. Horn. Toward a new taxonomy for pragmatic inference: Q-based and R-based implicature. In D. Schiffrin, editor, *Meaning, Form and Use in Context*, pages 11 – 42. Georgetown University Press, Washington, 1984.
- Gerhard Jäger and Robert van Rooij. Language structure: psychological and social constraints. *Synthese*, 159(1):99–130, 2007. doi: 10.1007/s11229-006-9073-5.
- Gerhard Jäger. Evolutionary game theory and typology: A case study. *Language*, 83(1):74–109, 2007. doi: 10.2307/4490338.
- Steven T Piantadosi. Zipfs word frequency law in natural language: A critical review and future directions. *Psychonomic bulletin & review*, 21(5):1112–1130, 2014. doi: 10.3758/s13423-014-0585-6.
- Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. Compression and communication in the cultural evolution of linguistic structure. *Cognition*, 141:87–102, 2015.
- C. Kemp and T. Regier. Kinship categories across languages reflect general communicative principles. *Science*, 336(6084):1049–1054, 2012. doi: 10.1126/science.1218811.
- Simon Kirby. Learning, bottlenecks and the evolution of recursive syntax. In Ted Briscoe, editor, *Linguistic Evolution Through Language Acquisition*, pages 173–204. Cambridge University Press (CUP), 2002. doi: 10.1017/cbo9780511486524.006.
- Jacob Feldman. Minimization of boolean complexity in human concept learning. *Nature*, 407(6804):630–633, 2000.
- Nick Chater and Paul Vitányi. Simplicity: a unifying principle in cognitive science? *Trends in Cognitive Sciences*, 7(1):19–22, 2003. doi: 10.1016/s1364-6613(02)00005-0.
- Steven T. Piantadosi, Joshua B. Tenenbaum, and Noah D. Goodman. Bootstrapping in a language of thought: A formal model of numerical concept learning. *Cognition*, 123(2):199–217, 2012a.
- Steven T. Piantadosi, Joshua B. Tenenbaum, and Noah D. Goodman. Modeling the acquisition of quantifier semantics: a case study in function word learnability, under review.
- Paul Grice. Logic and conversation. In *Studies in the Ways of Words*, chapter 2, pages 22–40. Harvard University Press, Cambridge, MA, 1975.
- Steven T. Piantadosi, Harry Tily, and Edward Gibson. The communicative function of ambiguity in language. *Cognition*, 122(3):280–291, 2012b. doi: 10.1016/j.cognition.2011.10.004.
- Anton Benz, Gerhard Jäger, Robert Van Rooij, and Robert Van Rooij, editors. *Game theory and pragmatics*. Springer, 2005a.
- Michael Franke. *Signal to Act: Game Theoretic Pragmatics*. PhD thesis, University of Amsterdam, 2009.
- Duncan R. Luce. *Individual choice behavior: a theoretical analysis*. Wiley, 1959.

- Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1998.
- M. A. Nowak and D. C. Krakauer. The evolution of language. *Proceedings of the National Academy of Sciences*, 96(14):8028–8033, 1999.
- Martin A. Nowak, Joshua B. Plotkin, and Vincent A. A. Jansen. The evolution of syntactic communication. *Nature*, 404(6777):495–498, 2000. doi: 10.1038/35006635.
- Martin A. Nowak, Natalia L. Komarova, and Partha Niyogi. Computational and evolutionary aspects of language. *Nature*, 417(6889):611–617, 2002. doi: 10.1038/nature00771.
- Anton Benz, Gerhard Jäger, Robert Van Rooij, and Robert Van Rooij, editors. *An Introduction to Game Theory for Linguists*. Springer, 2005b.
- Simon Kirby, Tom Griffiths, and Kenny Smith. Iterated learning and the evolution of language. *Current Opinion in Neurobiology*, 28:108–114, 2014. doi: 10.1016/j.conb.2014.07.014.
- Henry Brighton. Compositional syntax from cultural transmission. *Artificial Life*, 8(1):25–54, 2002. doi: 10.1162/106454602753694756.
- S. Kirby, M. Dowman, and T. L. Griffiths. Innateness and culture in the evolution of language. *Proceedings of the National Academy of Sciences*, 104(12):5241–5245, 2007. doi: 10.1073/pnas.0608222104.
- Laurence R. Horn. *On the Semantic Properties of Logical Operators in English*. Indiana University Linguistics Club, Bloomington, IN, 1972.
- Gerald Gazdar. *Pragmatics, Implicature, Presupposition and Logical Form*. Academic Press, New York, 1979.
- Elizabeth Closs Traugott. Historical pragmatics. In Laurence R. Horn and Gregory Wand, editors, *The Handbook of Pragmatics*, pages 538–561. Blackwell Publishing, 2004.
- Johan van der Auwera. On the diachrony of negation. In *The Expression of Negation*, pages 73–110. Walter de Gruyter GmbH, 2010. doi: 10.1515/9783110219302.73.
- Yarden Katz, Noah D Goodman, Kristian Kersting, Charles Kemp, and Joshua B Tenenbaum. Modeling semantic cognition as logical dimensionality reduction. In *Proceedings of Thirtieth Annual Meeting of the Cognitive Science Society*, 2008.
- Wim De Neys and Walter Schaeken. When people are more logical under cognitive load. *Experimental Psychology*, 54(2):128–133, 2007.
- Yi Ting Huang and Jesse Snedeker. Online interpretation of scalar quantifiers: Insight into the semantics–pragmatics interface. *Cognitive Psychology*, 58(3):376–415, 2009.
- Leon Bergen, Noah D Goodman, and Roger Levy. That’s what she (could have) said: How alternative utterances affect language use. In *Proceedings of Thirty-Fourth Annual Meeting of the Cognitive Science Society*, 2012.