# Communicative pressures at the semantics-pragmatics interface:
# Learning biases may prevent the lexicalization of pragmatic inferences

( – draft September 26, 2016— )

**Abstract**

   Certain lexical meanings enable for pragmatic enrichments in a notably productive fashion. This raises the challenge to justify their regular selection, in particular, over alternatives that codify semantically what is conveyed pragmatically. To address this challenge, we propose a general model that integrates iterated Bayesian learning in the replicator-mutator dynamics. This model allows for population-level analyses of the effects of linguistic pressures on probabilistic language users with varied degrees of pragmatic sophistication and distinct languages. We showcase the model's use and predictions in a case study on the (lack of) lexicalization of scalar implicatures. The results suggest simpler semantic representations to be 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.

## 1    The semantics-pragmatics divide

   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 pragmatic enrichments.

   In recent years, similar questions about the emergence and change of linguistic features have lead to a surge in models to address them (see Steels 2015 and Tamariz and Kirby 2016 for recent overviews). Our starting point is given by the overarching argument that has crystalized from 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 as well as learnability.

   We build on these insights by modeling these pressures using the replicator-mutator dynamics (see Hofbauer and Sigmund 2003 for an overview). This allows for the inspection of their interaction 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 is made precise by considering probabilistic models of rational language use in populations with distinct lexica (Frank and Goodman 2012, Franke and Jäger 2014, Bergen et al. 2016).

# 2   Simplicity, expressivity, and learnability

The emergence and change of linguistic structure is influenced by many intertwined factors, ranging 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 their means of fulfillment. In the following, our focus lies on the latter, cultural factor, wherein processes of linguistic change are understood as shaped by its use and transmission. That is, as a result of cultural evolution.

At latest since Zipf's (1949) rationalization of the observation that word frequency rankings can be approximated by a power law distribution as competing hearer and speaker preferences, the idea that linguistic change is influenced by communicative pressures has played a pivotal role in synchronic and diachronic analyses (e.g. Martinet 1962, Horn 1984, Jäger and van Rooij 2007, Jäger 2007, Piantadosi 2014, Kirby et al. 2015).

As noted above, expressivity and learnability are two major competing pressures. Their opposition becomes particularly clear when considering their consequences in the extreme (cf. Kemp and Regier 2012, Kirby et al. 2015). On the one side, a language with a single form is easy to learn but lacking in expressivity for most purposes. 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.

More important for our purposes is the problem 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. In the following, we argue an integral part of the answer to be that learners are a priori biased towards simpler, more compressed, lexical representations. This corresponds to the argument that rational learners should prefer simpler over more complex explanations of data (Feldman 2000, Chater and Vitányi 2003, Piantadosi et al. 2012a, Kirby et al. 2015, Piantadosi et al. under review). 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 remainder of this section introduces the individual components of the model in more detail, as well as the assumptions underlying them. These are: (i) languages and their use, (ii) pressures towards expressivity and learnability, regulated by the replicator and mutator dynamics, respectively, as well as (iii) a bias towards simpler semantic representations, codified as a language learner's prior. After laying out the model, we discuss its application to the lack of lexicalization of scalar implicatures.

## 2.1 Languages and linguistic behavior

Lexica codify the truth-conditions of a language's expressions, i.e., its semantics. A convenient way to represent such lexica 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 instance, the following two lexica fragments determine the truth-conditions of two messages, $m_1$ and $m_2$, for two states, $s_1$ and $s_2$:

$$
L_a = \begin{array}{c} \\ s_1 \\ s_2 \end{array} \begin{array}{cc} m_1 & m_2 \\ \left( \begin{array}{cc} 1 & 0 \\ 1 & 1 \end{array} \right) \end{array}
\qquad\qquad
L_b = \begin{array}{c} \\ s_1 \\ s_2 \end{array} \begin{array}{cc} m_1 & m_2 \\ \left( \begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array} \right) \end{array}
$$

In words, according to lexicon $L_a$, $m_1$ is true of state $s_1$ as well as of $s_2$. In contrast, message $m_1$ is only true of $s_1$ in $L_b$. Otherwise, the two languages are truth-conditionally equivalent.

To make the distinction betweeen 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 by their lexica only. 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 should use $m_1$ to signal state $s_1$ given that $s_2$ can unambiguously be conveyed with $m_2$. Analogously, should rational hearers expect their interlocutors to reason along these lines, they will interpret ambiguous $m_1$ accordingly. Note in particular that according to this strenghtening of $m_1$, $L_a$ is indistinguishable from $L_b$ in terms of expressivity if its users are pragmatic reasoners.

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 hierarchy over reasoning types. 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. The behavior of literal and pragmatic hearers of a language $L$ is given by their respective selection functions in (1) and (3). Mutatis mutandis for the speaker functions in (2) and (4).

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

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

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

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

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 choice in (2) is regulated by a soft-max parameter $\lambda$, $\lambda \geq 1$ (Luce 1959, Sutton and Barto 1998). As $\lambda$ increases, choices made in production are more rational in that higher values lead to more deterministic in line with expected utility maximization.

For the most part, pragmatic behavior mirrors its literal counterpart. As described above, their difference lies in that level $n + 1$ speakers/hearers reason about level $n$ hearer/speaker behavior instead of solely relaying 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. Additionally, pragmatic production is further regulated by a parameter $\alpha$ 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 particular level of linguistic sophistication, yields a type $t$. Types are the basic units on which our population dynamics operate.

## 2.2   Replication & expressivity

Communicative efficiency, or 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 success in communication confers it a higher fitness relative to less successful ones. As a consequence they replicate more than other types, increasing their proportion in the population. This association of a type's communicative success within a population with changes in the 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 captured by a vector $x$, where $x_i$ is type $i$'s proportion in the population. The fitness of a type $i$, $f_i$, is given by its expected utility in this population, $f_i = \sum_j x_j \mathrm{EU}(t_i, t_j)$. That is, its fitness is the sum of its expected communicative success with other types weighted by the latter type's population share. The expected utility of $i$ and $j$ is obtained by considering the expected utility of speaker $i$ interacting with hearer $j$, and vice versa: $\mathrm{EU}(t_i, t_j) = [U_S(t_i, t_j) + U_R(t_i, t_j)]/2$. $U_S(x, y)$ and $U_R(x, y)$ are respectively $\sum_s P(s) \sum_m S_n(m|s; L) \sum_{s'} R_o(s'|m; L)\delta(s, s')$ and $U_S(y, x)$ for $n$ and $o$ being the reasoning level of $x$ and $y$, and $\delta(s, s') = 1$ iff $s = s'$ and 0 otherwise.[1] This quantity is symmetric, reflecting the probability of two types' mutual understanding. Lastly, the average fitness of the population is captured by $\Phi$, $\Phi = \sum_i x_i f_i$. 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. In analogy to this kind of replication, many aspects of natural language are subject to processees of transmission and change across varied time-spans. For example, the replicator equation can be understood as a learning 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 adapt their lexica and their use to that which works best within their population. It should be stressed, however, that the model itself is compatible with either view.

Nowak et al. did not only consider replication, but also recognized the important role of the variation that is introduced by a language's transmission across generations, construed as mutation. Due to this process, the offspring of a type may end up adopting a different type than that of its parent. Crucially, in this work mutation rates were modelled as begin independent from a type. This means that the variation introduced by generational turnovers did not depend on factors such as the relative learnability of a type. To address the issue of selecting a particular type over (near) functional equivalents, we turn to a different strand of research in cultural evolution: *iterated learning*.

## 2.3   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. For linguistic purposes this process can be thought of as 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 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

---

[1]Note that the definition of $U_R(\cdot, \cdot)$ implies equal sender and receiver payoff in an interaction. This need not be so in the general case but suffices for our application.

type to adopt from the resulting posterior distribution.

Due to the pressure towards learnability it exherts, iterated learning generally leads to simpler and more regular languages (surveys of empirical data and models are given in Kirby et al. 2014 and Tamariz and Kirby 2016). Importantly, experimental and mathematical results suggest the results of this process to reflect learners' learning biases, codified in the following as a prior $P \in \Delta(\mathcal{T})$. A way to think about this bias is as the amount of data a learner would require in order to adopt a language – or, in our case, a combination of a lexicon and a signaling behavior (cf. Griffiths and Kalish 2007:450). Crucially, 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, from an individual's perspective, the likelihood of learning a language 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 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 in the range between posterior sampling and MAP, $l \geq 1$. When $l = 1$ learners sample from the posterior. As $l$ increases towards infinity, the learners' tendency maximize the posterior increases.

More generally, 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 learnbility, this quantity is proportional to the probability of $i$ producing the learning data and that of $j$ given the data.

The elements of the set of learning data $D$ are sequences of length $k$ 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 $\{\langle s_i, m_j \rangle \,|\, s_i \in S, m_j \in M\}$ and $D$ is the set of all such sequences. Having fixed $D$,

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

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.

## 2.4   Summary

We argued for expressivity, learnability and simplicity as important 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 these evolutionary dynamics are described by the replicator-mutator dynamics (Hofbauer and Sigmund 2003):

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

The basic units that the dynamics operate on are a combination of lexica lexicon and their use; 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 infered by new generations of naïve learners.

In sum, the innovation of this model lies in its combination of functional pressure on successful communication, effects of learning biases on (iterated) Bayesian language learning Griffiths and Kalish 2007, and probabilistic models of language use in populations with distinct lexica (Frank and Goodman 2012, Franke and Jäger 2014, Bergen et al. 2016). In particular, this synthesis enables for the investigation of the effects of communicative pressures on the semantics-pragmatics interface. In doing so, it links the previously disconnected areas of rational probabilistic language use and cultural evoltion.

# 3 Lack of lexicalization of pragmatic inferences

With this model, we set out to investigate the prevalence of lexical meanings that allow for regular pragmatic enrichments over other alternatives. A particularly well-studied type of conventional pragmatic enrichment 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 necessarily what their use is 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'. Analogously, a speaker can rely on her interlocutor to draw this inference without having to express this upper-bound overtly, e.g. by stating *some but not all*. In other words, mutual reasoning about rational language use supplies a bound that rules out stronger alternatives pragmatically.

This corresponds to our previous description of the pragmatic use of lexicon $L_a$, repeated below for convenience. A pragmatic hearer who reasons about a speaker's use of message $m_1$ will associate it more strongly with $s_1$ than with $s_2$ given that the latter is unambiguously associated with $s_2$. The strength of this association dependens on the individuals' degree of rationality $\lambda$ and their prior over states. Conversely, a pragmatic speaker will reason about her interlocutor's interpretation and use the messages accordingly.

$$L_a = \begin{array}{c} \\ s_1 \\ s_2 \end{array} \begin{array}{cc} m_1 & m_2 \\ \left( \begin{array}{cc} 1 & 0 \\ 1 & 1 \end{array} \right) \end{array} \qquad\qquad L_b = \begin{array}{c} \\ s_1 \\ s_2 \end{array} \begin{array}{cc} m_1 & m_2 \\ \left( \begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array} \right) \end{array}$$

Our initial question can now be rephrased in terms of scalar implicatures by asking for justifiations 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 lexicalizing it as in $L_b$. More poignantly, would it not serve language users better if weak(er) expressions such as *warm*, *or*, *some* or *big* were truth-conditionally incompatible with stronger alternatives such as, respectively, *hot*, *and*, *all* and *huge*? This question is particularly striking considering the number of expressions that license such inferences across languages (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 expres-

sions endows them a broader range of application. Scalar expressions 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 *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 these means are preferred over expressing it overtly through a longer expression, e.g., by stating 'some but not all' explicitly. Importantly, although dispreferred due to its relative length and complexity, morphosyntactic disambiguation still allows speakers to enforce an upper-bound to 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 lexicalizes two expressions instead of a single scalar one; one with and one lacking an upper-bound.[2] Should the explanation rest on purely functional grounds, 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.

However, these conditions are not convincing in their role as central explanatory devices 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 . 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 the hypothetical alternative strikes us as too strong of an assumption. As for conditions (3) and (4), these seem mostly like technical solutions, without a proper empirical basis.

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.[3]

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

---

[2]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). To the best of our knowledge, this claim stands unchallenged.

[3]T: 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.

leverage.

## 3.1 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 semantics to $m_1$. Combining a signaling behavior with one of these 6 lexica yields a total of 12 distinct types.[4]

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

---

[4]While 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.

stronger learning bias against languages that lexicalize upper-bounds, i.e., $L_2, L_4$ and $L_6$.

## 3.2 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 2.3, 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 semantics. However, it is instructive to first inspect the effect of these pressures in isolation. For this purpose, we focus on three pragmatic types.[5] 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 $\lambda$ and $\alpha$ as they have a bearing on a type's fitness. Low $\alpha$ 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 $\lambda$ 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, $\lambda$ 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 $\lambda$-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.

---

[5]Pragmatic reasoning allows language users to refine their (possibly erroneous) choices. Therefore, it is advantageous even for those types that codify an upper-bound lexically.
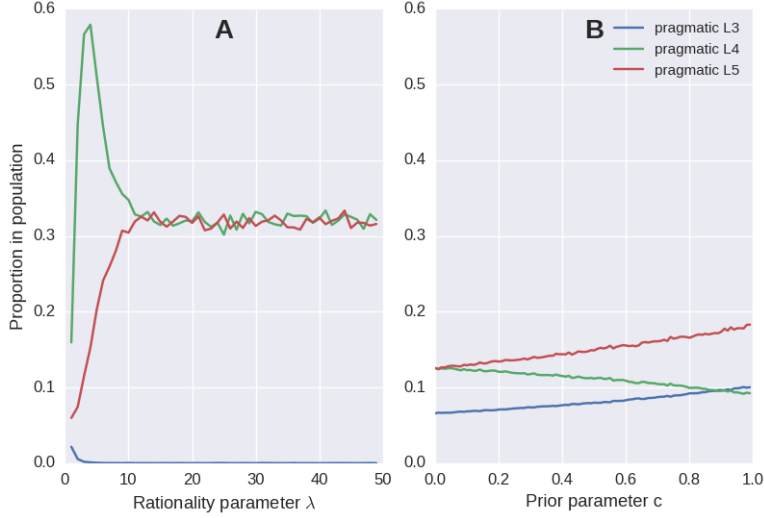
Figure 1: Mean proportions of target types after 20 generations in 1000 populations with only replication (A; $\alpha = 1$) and only mutation (B; $\alpha = 1, \lambda = 30, k = 5$).

**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 $L_3$. 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 3.2 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
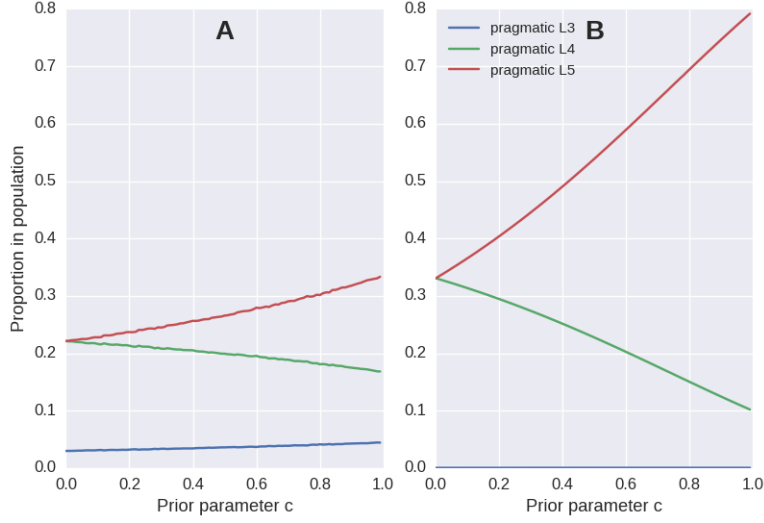
10

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 ($\alpha = 1, \lambda = 20, k = 5$).

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 $\lambda$ and $\alpha$, 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 $L_5$. 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. 3.2.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

## 3.3 Discussion

changes in sequence length influence the population 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 $\alpha$ and $\lambda$ – set to 1 and 50 in the following. Overall, lower rationality in $\lambda$ or more pragmatic violations in $\alpha$ lead to a higher selection of lexica with semantic upper-bounds. The fitness of pragmatic behavior increases with higher $\lambda$-/$\alpha$-values. In other words, these parameters level the functional contrast between $L_4$ and $L_5$.

11

| | $l=1$ | | | | | $l=10$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| c | 0 | .01 | .05 | .1 | .8 | 0 | .01 | .05 | .1 | .8 |
| lit. $L_1$ | .03 | .03 | .03 | .03 | .04 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| lit. $L_2$ | .03 | .03 | .03 | .03 | .01 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| lit. $L_3$ | .03 | .03 | .03 | .03 | .04 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| lit. $L_4$ | .07 | .07 | .07 | .07 | .06 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| lit. $L_5$ | .04 | .04 | .05 | .05 | .06 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| lit. $L_6$ | .04 | .04 | .04 | .04 | .04 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| prg. $L_1$ | .03 | .03 | .03 | .03 | .04 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| prg. $L_2$ | .03 | .03 | .03 | .03 | .01 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| prg. $L_3$ | .03 | .03 | .03 | .03 | .04 | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |
| 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 $l = 3$ ($\alpha = 1, \lambda = 30, k = 5$), $\epsilon < 0.005$
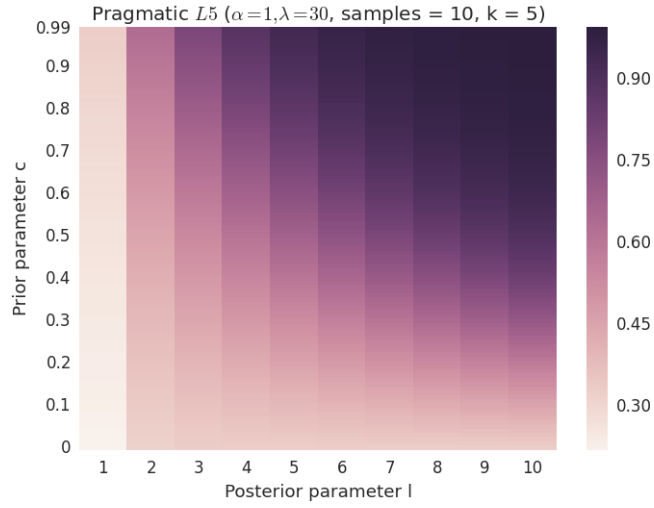


Figure 3: Mean proportion of pragmatic $L_5$ in 1000 populations after 20 generations ($\alpha = 1, \lambda = 30, k = 5$)

- There is experimental evidence that scalar implicature calculation costs effort and takes additional processing time (Schaeken, Snedeker, ...).

# 4 General discussion

The model combines:

- game-theoretical models of functional pressure towards efficient communication

- effects of learning biases on (iterated) language learning

- assumptions about particular learning biases (Piantadosi et al. 20XX)

- probabilistic speaker and listener types of various pragmatic sophistication (Frank & Goodman, 2012; Franke & Jger, 2014)

- speaker and listener types with different lexica (Bergen et al. 2012, to appear)

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.
**Extensions:**

**(I) Cost for pragmatic reasoning.** At least in the CogSci setup the effect of adding cost to pragmatic reasoning is unsurprising: High cost for pragmatic signaling lowers the prevalence of pragmatic types. Lexica that semantically encode an upper-bound benefit the most from this. However, the cost needed to be substantial to make the pragmatic English-like lexicon stop being the incumbent type (particularly when learning is communal).

**(II) Negative learning bias.** Instead of penalizing complex semantics (semantic upper-bounds) one may consider penalizing simple semantics (no upper-bounds). This is useful as a sanity check but also yields unsurprising results in the CogSci setup: The more learners are biased against simple semantics, the more prevalent are lexica that semantically encode upper-bounds.

**(III) Inductive bias.** A second learning bias that codifies the idea that lexica should be uniform, i.e. be biased towards either lexicalizing an upper-bound for all weaker alternatives in a scalar pair or for none.

**(IV) Uncertainty.** The other advantage of non-upper bounded semantics lies in being non-committal to the negation of stronger alternatives when the speaker is uncertain. Adding this to the model requires the most changes to our present setup and some additional assumptions about the cues available to players to discern the speaker's knowledge about the state she is in.

**(V) More scalar pairs.** Taking into consideration more than one scalar pair. Preliminary results suggest that this does not influence the results in any meaningful way without further additions, e.g. by (III).

**(VI) More lexica.** Not necessary. Preliminary results suggest that considering more lexica has no noteworthy effect on the dynamics (tested with all possible 2x2 lexica).

**(VII) State frequencies.** Variations on state frequencies. This may have an interesting interaction with (III).

**(VIII) Reintroduction of communal learning.** One possibility: The probably $N_{ij}$ with which a child of $t_i$ adopts $t_j$ could be the weighted sum of $Q_{ij}$ (as before) and a vector we get from learning from all of the population: $L_j = \sum_d P(d|\vec{p})P(t_j|d)$, where $P(d|\vec{p}) = \sum_i P(d|t_i)\vec{p_i}$ is the probability of observing $d$ when learning from a random member of the present population distribution.

# 5  Conclusion

# 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, Presuposition 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.