

# Tracing the cultural evolution of meaning at the semantics-pragmatics interface

( – draft March 9, 2017— )

## Abstract

According to standard linguistic theory, the meaning of an utterance is the product of conventional semantic meaning and general pragmatic reasoning applied to the context of utterance. This implies that models of cultural evolution of meaning should likewise take into consideration that observable language use is a complex interaction of semantic representations and pragmatic use. To this end, we present a game theoretic model of cultural language evolution where communicative pressures work on abstract semantic representations and pragmatic patterns of use. Our model traces two evolutionary forces and their interaction: (i) fitness-based pressure towards communicative efficiency and (ii) systematic transmission perturbations when linguistic knowledge is transferred from one agent to another. We illustrate the model based on a case study showing that learning biases that favor simple semantic representations can prevent the lexicalization of pragmatic inferences as long as mutual reasoning counteracts a potential loss in expressivity from such semantics.

## 1 Introduction

What is conveyed usually goes beyond what is said. A request for a blanket can be politely veiled by uttering “I’m cold.” The temporal succession of events can be communicated by the order in which conjuncts appear as in “I traveled to Paris and got married.” An invitation can be declined by saying “I have to work.” An influential explanation of the relation between the literal meaning of expressions and what they are intended and interpreted to convey is due to Grice (1975), characterizing pragmatic use and interpretation as a process of mutual reasoning about rational language use. For instance, under the assumption that the speaker is cooperative and relevant, “I have to work” may be interpreted as providing a reason why the speaker will not be able to accept the invitation, going beyond its literal meaning. Some of these enrichments are rather *ad hoc*. Others show striking regularities, such as the use of ability questions for polite requests (“Could you please ...?”), or certain enrichments of lexical meanings such as *and* to mean *and then*.

A particularly productive and well studied class of pragmatic enrichments are scalar implicatures (Horn 1984, Hirschberg 1985, Levinson 1983, Geurts 2010). Usually, the utterance of a sentence like “I own some of Johnny Cash’s albums” will be taken to mean that the speaker does not own all of them. This is because, if the speaker had them all, he could have used the stronger word *all* instead of *some* in his utterance, thereby making a more informative statement. Scalar implicatures, especially the inference from *some* to *some but not all*, have been studied extensively, both theoretically (e.g., Sauerland 2004, Chierchia et al. 2012, van Rooij and de Jager 2012) as well as experimentally (e.g., Bott and Noveck 2004, Huang and Snedeker 2009, Grodner et al. 2010, Goodman and Stuhlmüller 2013, Degen and Tanenhaus 2015). While there is much dispute in this domain about many interesting details, a position endorsed by a clear majority

is that a word like *some* is underspecified to mean *some and maybe all* and that the enrichment to *some but not all* is part of some regular process with roots in pragmatics.

If this majority view is correct, the question arises how such a division of labor between semantics and pragmatics could have evolved. Models of language evolution abound. There are simulation-based models studying the evolution of language in populations of communicating agents (Hurford 1989, Steels 1995, Lenaerts et al. 2005, Steels and Belpaeme 2005, Baronchelli et al. 2008, Steels 2011, Spike et al. 2016) and there are mathematical models of language evolution, mostly coming from game theory (Lewis 1969, Wårneryd 1993, Blume et al. 1993, Nowak and Krakauer 1999, Huttegger 2007, Skyrms 2010). Much of this work has focused on explaining basic properties such as compositionality and combinatoriality (e.g., Batali 1998, Nowak and Krakauer 1999, Nowak et al. 2000, Kirby and Hurford 2002, Kirby 2002, Smith et al. 2003, Gong 2007, Kirby et al. 2015, Verhoef et al. 2014, Franke 2016), but little attention has been paid to the interaction between conventional meaning and pragmatic use. What is more, many mathematical models explain evolved meaning as a regularity in the behavior of agents which maps objective states of the world to observable signals. There is no room in such a purely extensional approach to address the semantics-pragmatics division directly. Instead, we would need to look at richer representations of cognizing agents and their communicative interaction.

To fill this gap, we spell out a model of the co-evolution of conventional meaning and pragmatic reasoning types. The objects of replication and selection are pairs of lexical meanings and general types of pragmatic behavior, which we represent using state-of-the-art probabilistic cognitive models of pragmatic language use (Frank and Goodman 2012, Franke and Jäger 2016, Goodman and Frank 2016). Replication and selection are described by the *replicator mutator dynamic*, a general and established model of evolutionary change in large and homogeneous populations (Hofbauer 1985, Nowak et al. 2000; 2001, Hofbauer and Sigmund 2003, Nowak 2006). The approach allows us to study the interaction between (i) evolutionary pressure towards communicative efficiency and (ii) possible infidelity in the transmission of linguistic knowledge, caused by factors such as inductive learning biases and general aspects of learnability. Considering transmission of linguistic knowledge is important because neither semantic meanings nor pragmatic usage patterns are directly observable. Instead, language learners have to infer these unobservables from the observable behavior in which they result. We formalize this process as a form of Bayesian inference. Our approach thereby contains a well-understood model of iterated Bayesian learning (Griffiths and Kalish 2007) as a special case, but combines it with functional selection, here formalized as the most versatile dynamic from evolutionary game theory; the replicator dynamic (Taylor and Jonker 1978). Section 2 introduces this model.

Section 3 applies this model to a case study on scalar implicatures. We discuss a setting in which the majority view of underspecified lexical meanings and pragmatic enrichments emerges if selection and transmission infidelity are combined. In particular, we show that inductive learning biases of Bayesian learners that favor simpler lexical meanings can lead to this outcome.

We see the main contribution of this work as conceptual and technical, not as a definite answer to the question why scalar implicatures emerge. It rather demonstrates how current probabilistic cognitive modeling of language use and evolutionary modeling can be fruitfully combined to study the co-evolution of semantics and pragmatics side-by-side. Reversely, the approach taken here may be seen as a first step towards giving an evolutionary rationale for empirically successful probabilistic models of language use that embrace the majority view of the division of labor between semantics and pragmatics. Section 4 elaborates on these points.

## 2 Model

### 2.1 Expressivity and learnability at the semantics-pragmatics interface

The emergence and change of linguistic structure is influenced by many intertwined factors. These range from biological and socio-ecological to cultural ones (Benz et al. 2005, 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. In the following, our focus lies on cultural aspects, wherein processes of linguistic change are viewed as shaped by language use and its transmission, i.e., as a result of a process of cultural evolution (Pagel 2009, Thompson et al. 2016).

The idea that language is an adaptation to serve a communicative function has played a pivotal role in synchronic and diachronic analyses at least since Zipf’s (1949) explanation of word frequency rankings as a result of competing hearer and speaker preferences (e.g., in Martinet 1962, Horn 1984, Jäger and van Rooij 2007, Jäger 2007, Piantadosi 2014, Kirby et al. 2015). If processes of selection, such as conditional imitation or reinforcement, favor more communicatively efficient types of behavior, languages are driven towards semantic expressivity (e.g., Nowak and Krakauer 1999, Skyrms 2010). But pressure towards communicative efficiency is not the only force that shapes language. Learnability is another, as natural languages need to be learnable to survive their faithful transmission across generations. Clearly, an unlearnable code will not make it past the one happy fellow who invented it. Importantly, even small biases implicit in acquisition can build up and have quite striking effects on an evolving language in a process of iterated learning (Kirby and Hurford 2002, Smith et al. 2003, Kirby et al. 2014).

While natural languages are pressured for both expressivity and learnability these forces may pull in opposite directions. Their opposition becomes particularly clear when considering the extreme (cf. Kemp and Regier 2012, Kirby et al. 2015). A language with a single form-meaning association is easy to learn but lacking in expressivity. Conversely, a language that associates a distinct form with all possible meanings a speaker may want to convey is maximally expressive but challenging to acquire.

An elegant formal approach to capturing the interaction between expressivity and learnability is the *replicator mutator dynamic* (Hofbauer 1985, Nowak et al. 2000; 2001, Hofbauer and Sigmund 2003, Nowak 2006). In its simplest, discrete-time formulation, the RMD defines the frequency  $x'_i$  of each type  $i$  in a population at the next time step as a function of: (i) the frequency  $x_i$  of each type  $i$  before the step, (ii) the fitness  $f_i$  of each type  $i$  before the step, and (iii) the probability  $Q_{ji}$  that an agent who wants to imitate, adopt, or learn the type of an agent with type  $j$  ends up acquiring type  $i$ :

$$x'_i = \sum_j Q_{ji} \frac{x_j f_j}{\sum_k x_k f_k}. \quad (1)$$

The RMD consists of two components: fitness-based selection and transmission biases. This becomes most transparent when we consider an equivalent formulation in terms of a step-wise application of the discrete-time replicator dynamic (Taylor and Jonker 1978) on the initial population vector  $\vec{x}$  and its subsequent multiplication with a mutation matrix  $Q$ :

$$x'_i = (M(RD(\vec{x})))_i, \quad (2)$$

where

$$(\text{RD}(\vec{x}))_i = \frac{x_i f_i}{\sum_k x_k f_k} \quad \text{and} \quad (\text{M}(\vec{x}))_i = (\vec{x} \cdot Q)_i = \left( \sum_j x_j Q_{ji} \right)_i.$$

If the transmission matrix  $Q$  is trivial in the sense that  $Q_{ji} = 1$  whenever  $j = i$ , the dynamic reduces to the replicator dynamic. The replicator dynamic is a model of fitness-based selection in which the relative frequency of type  $i$  will increase with a gradient proportional to its average fitness in the population. This dynamic is popular and versatile because it can be derived from many abstract processes of biological and cultural transmission and selection (for overview and several derivations see Sandholm 2010), including conditional imitation (e.g., Helbing 1996, Schlag 1998) or reinforcement learning (e.g., Börgers and Sarin 1997, Beggs 2005). If fitness  $f_i$  is the same for all types  $i$ , the replicator step is the identity map  $(\text{RD}(\vec{x}))_i = x_i$  and the dynamic reduces to a process of iteration of the transmission bias encoded in  $Q$ . In this way, the process in (1), equivalently (2), contains a model of iterated learning (Griffiths and Kalish 2007). [MF: should we include a simple example here? I have an example from a lecture ready at hand; it's a simple coordination game in a one-population setting.] [TB: Yes, I think most readers would appreciate it. We can always cut it out again if the reviewers prefer more compression.]

Where our goal is an application of this dynamic to the case of co-evolution of semantic meaning and pragmatic use, we need to fix what the relevant types are, how fitness is measured and how the mutation matrix is computed. These issues will be addressed, one by one, in the following.

## 2.2 Types: Lexica and linguistic behavior

Types are what cultural evolution operates on. In standard applications of evolutionary game theory, types correspond to ways of acting in a game, e.g., either cooperating or defecting in a prisoner's dilemma. [MF: maybe good to refer back to the example from before if there was one?] [TB: Yes!] For our present purpose, types are identified by their cognitive make-up. Since we are interested in the question under which conditions processes of cultural evolution will favor specific divisions of labor between lexical meaning and pragmatic use, a type is a pair consisting of a lexicon and a pragmatic strategy.

Lexica codify the truth-conditions of expressions. A convenient way to represent lexica is by  $(|S|, |M|)$ -Boolean matrices, where  $S$  is a set of states (meanings) and  $M$  a set of messages (forms available in the language). For example, suppose that there are two relevant world states  $S = \{s_{\exists-\forall}, s_{\forall}\}$ . In state  $s_{\exists-\forall}$  Chris owns some but not all of Johnny Cash's albums while in  $s_{\forall}$  Chris owns them all. Suppose that there are two messages  $M = \{m_{\text{some}}, m_{\text{all}}\}$  where  $m_{\text{some}}$  is short for a sentence like *Chris owns some of Johnny Cash's albums* and  $m_{\text{all}}$  for the same sentence with *some* replaced by *all*. [TB: I left  $s_{\emptyset}/m_{\text{none}}$  out because they are not relevant to the strengthening of  $m_{\text{some}}$ , but that means that it needs to be introduced a bit later. Feel free to add it if you think that's better, exposition-wise.] Lexica for this case would assign a truth value for each state-message pair. The following two lexica exemplify the distinction between a lexicalized upper-bound for *some* in  $L_{\text{bound}}$  and the widely assumed logical semantics with only a lower-bound in  $L_{\text{lack}}$ .

$$L_{\text{bound}} = \begin{matrix} & m_{\text{some}} & m_{\text{all}} \\ \begin{matrix} s_{\exists-\forall} \\ s_{\forall} \end{matrix} & \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{matrix} \quad L_{\text{lack}} = \begin{matrix} & m_{\text{some}} & m_{\text{all}} \\ \begin{matrix} s_{\exists-\forall} \\ s_{\forall} \end{matrix} & \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \end{matrix}$$

We distinguish between two kinds of pragmatic behavior. *Literal interlocutors* produce and interpret messages literally, being guided only by their lexica. *Pragmatic interlocutors* instead engage in mutual reasoning to inform their choices. Recent probabilistic models of rational language use (Franke 2009, Frank and Goodman 2012, Franke and Jäger 2016, Goodman and Frank 2016) capture different types of pragmatic behavior in a reasoning hierarchy. The hierarchy’s bottom, level 0, corresponds to literal language use. Pragmatic language users of level  $n + 1$  act (approximately) rational with respect to level- $n$  behavior of their interlocutors. (3) and (4) define probabilistic behavior of literal hearers and speakers respectively: [MF: is it a problem that  $S$  denotes speakers and the set of states?] [TB: I think that’s fine.  $S_n(\cdot \mid \cdot)$  is not used beyond this section and the set  $S$  only by mention of its members. Otherwise we can, e.g., relabel  $M$  as  $F$ (orms) and use  $M$ (eanings) instead of  $S$ (tates). I’m fine with either alternative.]

$$H_0(s \mid m; L) \propto pr(s) L_{sm} \quad (3)$$

$$S_0(m \mid s; L) \propto \exp(\lambda L_{sm}) \quad (4)$$

According to (3), a literal hearer’s interpretation of a message  $m$  as a state  $s$  depends on her lexicon and her prior over states,  $pr \in \Delta(S)$ . For simplicity, in the following this prior is assumed to be uniform. A literal interpreter with lexicon  $L_{\text{bound}}$  from above would assign  $s_{\exists-\forall}$  a probability of  $H_0(s_{\exists-\forall} \mid m_{\text{some}}; L_{\text{bound}}) = 1$  after hearing  $m_{\text{some}}$ , while a literal interpreter with lexicon  $L_{\text{lack}}$  would assign  $s_{\exists-\forall}$  probability  $H_0(s_{\exists-\forall} \mid m_{\text{some}}; L_{\text{lack}}) = 0.5$ . As usual in probabilistic pragmatics models, speaker behavior 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 behavior that is increasingly in line with expected utility maximization. Expected utility of a message  $m$  in state  $s$  for a level  $n + 1$  speaker is here defined as  $H_n(s \mid m; L)$ , the probability that the hearer will assign to or choose the correct meaning. Put differently, how likely the hearer is thought to infer  $s$  when sent  $m$ . The formulation in (4) also allows false messages to be sent with a small positive probability in analogy to the definition of pragmatic speakers in probabilistic pragmatic models. This is necessary to guarantee a mutation matrix with only positive entries (see below). If  $\lambda = 1$  a literal speaker with either lexicon  $L_{\text{bound}}$  or  $L_{\text{lack}}$  produces  $m_{\text{some}}$  in  $s_{\exists-\forall}$  with probability  $S_0(m_{\text{some}} \mid s_{\exists-\forall}; L_{\text{bound}, \text{lack}}) \approx .73$ . The probability of producing  $m_{\text{some}}$  in  $s_{\forall}$  is  $S_0(m_{\text{some}} \mid s_{\forall}; L_{\text{bound}}) \approx .27$  for  $L_{\text{bound}}$  and  $S_0(m_{\text{some}} \mid s_{\forall}; L_{\text{lack}}) \approx .5$  for  $L_{\text{lack}}$ . By contrast, if  $\lambda = 20$ , a literal speaker with either lexicon produces  $m_{\text{some}}$  in  $s_{\exists-\forall}$  with probability  $S_0(m_{\text{some}} \mid s_{\exists-\forall}; L_{\text{bound}, \text{lack}}) \approx 1$ . But the probability of producing  $m_{\text{some}}$  in  $s_{\forall}$  is  $S_0(m_{\text{some}} \mid s_{\forall}; L_{\text{bound}}) \approx 1$  for  $L_{\text{bound}}$  whereas it is still  $S_0(m_{\text{some}} \mid s_{\forall}; L_{\text{lack}}) \approx .5$  for  $L_{\text{lack}}$ .

Pragmatic behavior of level- $n + 1$  is similar to its literal counterparts but uses the interpretation or production behavior of a level- $n$  player instead of the lexical meaning:

$$H_{n+1}(s \mid m; L) \propto pr(s) S_n(m \mid s; L) \quad (5)$$

$$S_{n+1}(m \mid s; L) \propto \exp(\lambda H_n(s \mid m; L)) \quad (6)$$

Particularly important to our purpose is the fact that players of a level greater than 0 using  $L_{\text{lack}}$  can *pragmatically* associate  $m_{\text{some}}$  with  $s_{\exists-\forall}$  over  $s_{\forall}$ , in contrast to their literal counterparts. Intuitively, this is so because hearers reason a speaker to convey  $s_{\forall}$  with  $m_{\text{all}}$ , as this has a higher chance of leading to mutual understanding than using underspecified  $m_{\text{some}}$ . Therefore they will expect  $m_{\text{some}}$  to be more strongly associated with  $s_{\exists-\forall}$  than with  $s_{\forall}$  as it is already possible to clearly convey the latter state with a different message. The converse reasoning pattern holds for the speaker.

### 2.3 Fitness & fitness-based selection based on expressivity

Most evolutionary dynamics assume that the proportion of type  $i$  in a population will increase or decrease as a function of its relative fitness  $f_i$ . In the context of language evolution, fitness is frequently associated with expressivity, i.e., the ability to successfully communicate with other language users from the same population (e.g., Nowak and Krakauer 1999, Nowak et al. 2000; 2002). Under a biological interpretation, the assumption is that organisms have a higher chance of survival and reproduction if they are able to share and receive useful information via communication with peers. Under a cultural interpretation, the picture is that agents themselves strive towards communicative success and therefore occasionally adapt or revise their behavior to achieve a higher communicative success (see Benz et al. 2005:§3.3 for discussion).

The replicator equation gives us the means to make the ensuing dynamics precise, without necessarily committing to a biological or cultural interpretation. As above, the proportion of types in a given population is codified in a vector  $\vec{x}$ , where  $x_i$  is type  $i$ 's proportion. The fitness of type  $i$  is its average expected communicative success, or *expected utility* (EU), given the frequencies of types in the current population:

$$f_i = \sum_j x_j \text{EU}(t_i, t_j).$$

The expected utility  $\text{EU}(t_i, t_j)$  for type  $i$  when communicating with type  $j$  is the average success of  $i$  when talking or listening to  $j$ . Assuming that agents are speakers half of the time this yields:

$$\text{EU}(t_i, t_j) = 1/2 \text{EU}_S(t_i, t_j) + 1/2 \text{EU}_H(t_i, t_j),$$

where  $\text{EU}_S(t_i, t_j)$  and  $\text{EU}_H(t_i, t_j)$  are the expected utilities for  $i$  as a speaker and as a hearer when communicating with  $j$ , defined as follows, where  $n_i$  and  $n_j$  are type  $i$ 's and  $j$ 's pragmatic reasoning types and  $L_i$  and  $L_j$  are their lexica:

$$\begin{aligned} \text{EU}_S(t_i, t_j) &= \sum_s P(s) \sum_m S_{n_i}(m \mid s; L_i) \sum_{s'} R_{n_j}(s' \mid m; L_j) \delta(s, s') \\ \text{EU}_H(t_i, t_j) &= \text{EU}_S(t_j, t_i) \end{aligned}$$

As usual for cooperative communication,  $\delta(s, s') = 1$  iff  $s = s'$  and 0 otherwise.

### 2.4 Learnability

Languages are shaped not only by functionalist forces towards greater expressivity. Another important factor is the fidelity by which language is transmitted. Among others, linguistic production can be prone to errors, states or messages may be perceived incorrectly, and multiple languages may be compatible with the data learners are exposed to. These sources of uncertainty introduce variation in their transmission from one generation to the next. In particular, learning biases in the iterated transmission process can influence language evolution substantially.

In biological evolution, where types are expressed genetically, transmission infidelity comes into the picture through infrequent and mostly random genetic mutations. However, an agent's lexicon and pragmatic reasoning behavior is not inherited genetically. They need to be learned from observation. Concretely, when agents of type  $j$  want to adopt or imitate the linguistic behavior of type  $i$ , they observe the linguistic behavior of type  $i$  and need to infer what their type is from that. Iterated learning is a process in which languages are learned repeatedly from the observation of linguistic behavior of agents who have acquired the language from observation and inference before as well. In the simplest case there is a single teacher and a single learner.

After sufficient training the learner becomes a teacher and produces behavior that serves as input for a new learner. 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).

Following Griffiths and Kalish (2007) we model language acquisition as a process of Bayesian inference in which learners combine the likelihood of a type producing the witnessed learning input with prior inductive biases. Experimental and mathematical results on iterated learning suggest that the outcome of this process reflects learners’ inductive biases (e.g., Kirby et al. 2014). In a Bayesian setting these biases can be codified in a prior  $P \in \Delta(T)$ , which reflects the amount of data a learner requires to faithfully acquire the language of the teacher (cf. Griffiths and Kalish 2007:450). The extent of the prior’s influence has been shown to heavily depend on the learning strategy assumed to underly the inference process. On the one hand, early simulation results suggested that weak biases could be magnified by exposing learners to only small data samples (e.g., in Brighton 2002). On the other, Griffiths and Kalish’s (2007) mathematical characterization showed that iterated learning converged to the prior, i.e., the resulting distribution over languages corresponds to the learners’ prior distribution and is not influenced by the amount of input given to them. This difference in predictions can be traced back to differences in the selection of hypotheses from the posterior. Griffith & Kalish’s convergence to the prior holds for learners that sample from the posterior. More deterministic strategies such as the adoption of the type with the highest posterior probability, so-called *maximum a posteriori estimation* (MAP), increase the influence of both the prior and the data (Griffiths and Kalish 2007, Kirby et al. 2007). In the following, we use a parameter  $l \geq 1$  to modulate between posterior sampling and the MAP strategy. When  $l = 1$  learners sample from the posterior. The learners propensity to maximize the posterior grows as  $l$  increases.

Let  $D$  be the set of possible data that learners may be exposed to. This set  $D$  contains all sequences of state-message pairs of length  $k$ , e.g.,  $\langle\langle s_1, m_1 \rangle, \dots, \langle s_k, m_k \rangle\rangle$ . As  $k$  increases, learners have more data to base their inference on and so tend to recover the true types that generated a given sequence with higher probability. The mutation matrix  $Q$  of the replicator mutator dynamics in (1) can then be defined as follows:  $Q_{ji}$  is the probability that a learner acquires type  $i$  when learning from an agent of type  $j$ . The learner observes a length- $k$  sequence  $d$  of state-message pairs, but the probability  $P(d | t_j)$  with which sequence  $d = \langle\langle s_1, m_1 \rangle, \dots, \langle s_k, m_k \rangle\rangle$  is observed depends on type  $j$ ’s behavior:

$$P(d = \langle\langle s_1, m_1 \rangle, \dots, \langle s_k, m_k \rangle\rangle | t_j) = \prod_{i=1}^k S_{n_j}(m_i | s_i; L_j),$$

where, as before,  $n_j$  is  $j$ ’s pragmatic reasoning type and  $L_j$  is  $j$ ’s lexicon. For a given observation  $d$ , the probability of acquiring type  $i$  is  $F(t_i | d)$ , so that:

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

The acquisition probability  $F(t_i | d)$  given datum  $d$  is obtained by probability matching  $l = 1$  or a tendency towards choosing the most likely type  $l > 1$  from the posterior distribution  $P(\cdot | d)$  over types given the data, which is calculated by Bayes’ rule:

$$\begin{aligned} F(t_i | d) &\propto P(t_i | d)^l \quad \text{and} \\ P(t_i | d) &\propto P(t_i) P(d | t_i). \end{aligned}$$

## 2.5 Model summary

Expressivity and learnability are central to the cultural evolution of language. These components can be modelled, respectively, as replication based on a measure of fitness in terms of communicative efficiency and iterated Bayesian learning. Their interaction is described by the discrete time replicator mutator dynamics in (1), repeated here:

$$x'_i = \sum_j Q_{ji} \frac{x_j f_j}{\sum_k x_k f_k}.$$

This equation defines the frequency  $x'_i$  of type  $i$  at the next time step, based on its frequency  $x_i$  before the step, its fitness  $f_i$  and the probability that a learner infers  $i$  when observing the behavior of a type- $j$  agent. Fitness-based selection can be thought of as biological (fitness as expected relative number of offspring) or cultural (fitness of likelihood of being imitated or repeated). The types that the dynamic operates on are pairs consisting of a lexicon and a pragmatic use pattern. 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 learners. The learners' task is consequently to perform a joint inference over types of linguistic behavior and lexical meaning.

## 3 Scalar implicatures

Scalar implicatures are a particularly well-studied type of conventional pragmatic inferences. They 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 'Chris owns all of Johnny Cash's albums', it would also be true that 'Chris owns some of Johnny Cash's albums'. 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 stronger alternatives do not hold because the speaker used a weaker alternative instead. In this way, 'Chris owns some of Johnny Cash's albums' is strengthened to convey that he owns some but not all albums. A bound that rules out stronger alternatives is thusly not codified in the lexical meaning of weak alternatives but instead pragmatically supplied.

As noted earlier, this kind of strengthening is captured by the linguistic behavior of pragmatic types introduced in §2.2: A pragmatic hearer who reasons about the use of a message involving a weak scalar alternative will associate it more with a state in which stronger alternatives do not hold. This is so because a rational speaker would use a more informative message when in such a state. Conversely, a pragmatic speaker will reason about her interlocutor's expected interpretation and use the messages at her disposition accordingly.

Our initial question about the division of labor between semantics and pragmatics can be narrowed to the case of scalar implicatures by asking for a justification for the lack of lexical upper-bounds in weak scalar alternatives. That is, we ask why lexical meanings that lack upper-bounds and convey it pragmatically are regularly selected for over alternatives such as that of codifying the bound semantically. 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 natural languages.



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 a broader range of applicability by 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 ‘Chris owns some of Johnny Cash’s albums’ but she does not know whether he owns all, then the use of *some* lacking an upper-bound succinctly conveys her uncertainty. 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, e.g., by saying *some but not all* explicitly. Importantly, although morphosyntactic disambiguation may be dispreferred due to its relative length and complexity (Piantadosi et al. 2012b), 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 weak scalar expression – one with and one lacking an upper-bound. We see four conditions along this functionalist explanation that may pressure languages for English-like semantics over this alternative. First, contextual cues are very reliable. Second, morphosyntactic disambiguation is seldom necessary. Third, morphosyntactic disambiguation is only marginally dispreferred. Fourth, larger lexica are costly. Overall, neither condition seems convincing as a pivotal explanatory device for such a widespread phenomenon. The first two conditions 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 context, but to stipulate that they are so good as to undercut the benefit of safe communication provided by this 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.

Instead, the systematicity and typological spread of scalar implicatures together with the observation that monomorphemic expressions that lexically rule out stronger alternatives are unattested across languages (Horn 1984:252-267, Horn 1972, Traugott 2004, van der Auwera 2010) suggests that other forces may be at play. In what follows we investigate the predictions of our model under the assumption that the lack of lexicalization of scalar inferences may be accounted for by the relative representational simplicity of lexical meanings lacking an upper-bound over those that explicitly codify it. This difference is reflected in a learning bias towards more compressed lexical representations. That is, in a preference of learners for 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). While we do not want to argue that functional aspects as the ones discussed above do not play a role, we do see a clear benefit in exploring whether matters of transmission biases would not give us additional explanatory leverage.

### 3.1 Analysis

We analyze the model’s predictions in populations of types with one of the two signaling behaviors introduced earlier; literal or pragmatic. The former correspond to level 0 reasoners and the latter to ones of level 1. Higher level reasoning is not required to derive scalar implicatures from the

intuitive name	$s_\emptyset$	$s_{\exists-\forall}$	$s_\forall$	least complex formula	complexity
“all”	0	0	1	$A \subseteq B$	3
“some but not all”	0	1	0	$A \cap B \neq \emptyset \wedge A \neq \emptyset$	8
“some”	0	1	1	$A \cap B \neq \emptyset$	4
“none”	1	0	0	$A \cap B = \emptyset$	4
“none or all”	1	0	1	$\neg(A \cap B \neq \emptyset \wedge A \neq \emptyset)$	10
“not all”	1	1	0	$\neg(A \subseteq B)$	5

Table 1: Available concepts and their minimal derivation length

lexica we consider here, nor do they leave room for substantial pragmatic refinement.

**Lexica, types, and transmission bias.** We consider a state space with three states  $S = \{s_\emptyset, s_{\exists-\forall}, s_\forall\}$ . This space can be thought of as a partition of possible worlds into cells where none, some or all of the  $A$ s are  $B$ s, for some arbitrary fixed predicates  $A$  and  $B$ . Eight concepts can be distinguished based on their truth or falsity in three world states, six of which are not contradictory or tautologous. These are listed with mnemonic names in Table 1.

A lexicon  $L$  is a mapping  $M \rightarrow C$  from messages to concepts. With three messages there are  $6^3 = 216$  possible lexica. Of these, three kinds of lexica are of particular relevance for our purpose. First, lexica that assign the same concept to more than message. Such lexica lack in expressivity under either signaling behavior but may be easier to acquire than others. Second, lexica that conventionalized upper-bounds to realize a (quasi-)partition of the relevant semantic space. As discussed in relation to  $L_{\text{bound}}$ , users of these lexica need not resort to pragmatic reasoning to convey an upper-bound with weak scalar expressions. Instead, this information is already given by the semantics of their lexicon, giving them a functional advantage over others. Third and lastly,  $L_{\text{lack}}$ -style lexica that, paired with mutual reasoning, can be used to convey an upper-bound pragmatically. Overall, there are **N** lexica of the second kind and 12 of the third. The following three lexica exemplify each kind, with  $L_{\text{all}}$  being a communicatively suboptimal lexicon that conventionalizes the meaning of “all” with each message (cf. Table 1).

	$\underline{L_{\text{all}}}$			$\underline{L_{\text{bound}}}$			$\underline{L_{\text{lack}}}$		
	$m_{\text{none}}$	$m_{\text{some}}$	$m_{\text{all}}$	$m_{\text{none}}$	$m_{\text{some}}$	$m_{\text{all}}$	$m_{\text{none}}$	$m_{\text{some}}$	$m_{\text{all}}$
$s_\emptyset$	0	0	0	1	0	0	1	0	0
$s_{\exists-\forall}$	0	0	0	0	1	0	0	1	0
$s_\forall$	1	1	1	0	0	1	0	1	1

Recall that types are a combination of a lexicon and a manner of language use: literal or pragmatic. Accordingly, there are 432 types in this model. [Here. Add footnote on redundancy and results holding for smaller spaces with 40 types or the original 12 from Br](#)

In our setup, pragmatic  $L_{\text{lack}}$ -style speakers have a functional disadvantage over  $L_{\text{bound}}$ -style ones as

[TB: Many of these assign the same concept to more than one message. For simplicity, we restrict attention to only those lexica which avoid expressive redundancy. Hardcoding such an *mutual exclusivity bias* (e.g. Clark 2009), there are 20 non-redundant lexica. Functional pressure towards efficient communication will evidently favor non-redundant lexica, so this restriction is fairly innocuous but practical.]

$C \rightarrow_2 C \wedge C$	$X \rightarrow_1 \{A, B\}$
$C \rightarrow_2 \neg C$	$X \rightarrow_1 X \cap X$
$C \rightarrow_1 X \subseteq X$	$X \rightarrow_1 X \cup X$
$C \rightarrow_1 X \neq \emptyset$	
$C \rightarrow_1 X = \emptyset$	

Table 2: Toy grammar in a set-theoretic "language of thought"

[TB: Note however that we do not represent the contrast between lexical representations explicitly. Instead, the bias is directly encoded in the learners' prior over types.<sup>1</sup>]

Following our assumption of an inductive bias that favors simpler lexical representations, the prior biases learners against lexica in which a message holds true only of the first row, i.e., against messages that lexicalize an upper-bound that rules out state  $s_\forall$ . All other semantics are assumed to be a priori equally probable. This is captured 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  $s_{\exists-\forall}$  in  $t_i$ 's lexicon,  $c \in [0, 1]$ . Increments in the value of  $c$  thereby bring about a stronger bias against languages that lexicalize upper-bounds, i.e.,  $L_2, L_4$  and  $L_6$ .

[TB: This entire paragraphs needs to add: mention to LOT, other possible priors in LOT and why we use ours] The prior probability of a type is just the prior probability of its lexicon. The prior of a lexicon is a function of the complexity of the concepts in its image set. Lexica that use simpler concepts are *a priori* more likely. This can be motivated by assuming that learners beam search for suitable concepts to map onto overt signals by (probabilistically) considering simpler concepts first. Many ways for defining complexity of a concept are conceivable. If strong empirical claims were at stake here, empirically motivated measures of complexity should be used. For the sake of a non-trivial example, we follow Piantadosi et al. (under review) and related work to define complexity of a concept as a function of its derivation cost in a (weighed or probabilistic) generative "language of thought". For concreteness of example, consider the toy grammar of concepts in Table 2. This grammar uses basic set-theoretic operations to form expressions which can be evaluated as true or false in our three world states. Applications of generative rules have a cost attached to them. (Alternatively, a probability.) Here we simply assume that Boolean combinations of concepts are more complex than "atomic" concepts and that otherwise each rule application adds the same cost unit. Table 1 lists, for each concept, the least complex formula derivable in this grammar that has the appropriate truth conditions. A simple way of defining priors over a lexicon is:

$$P(L) = \prod_{c \in Im(L)} P(c) \quad , \text{ with } \quad P(c) \propto \max_{c'} Compl(c') - Compl(c) + 1 ,$$

where  $Compl(c)$  is the complexity of each concepts.

The space of possible lexica is given in Table . These (2, 2)-Boolean matrices are the simplest ones that allow us to make the contrast between the presence or absence of an upper-bound and

<sup>1</sup>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; 2012a and references therein). At present we decide against such an enrichment in favor of a stronger focus on the general predictions of the model and the interaction of its components.

the use of scalar implicatures precise. As illustrated in Section 2.2, one may think of the state corresponding to the first row of any such lexicon as a “some but not all”-state,  $s_{\exists-\forall}$ , and the second as an “all”-state,  $s_{\forall}$ . The literal meaning of weak scalar expressions such as English *some* then corresponds to a message true of both rows in these fragments.

Lexica  $L_1$  to  $L_3$  are not optimal for communication because they assign all their messages to the same state(s). This failure to be able to associate a state to single a form inevitably leads to a communicative disadvantage in their use.  $L_4$  and  $L_5$  are our target lexica. The former assigns upper-bounded semantics to  $m_{\text{some}}$  (the first matrix’s column) whereas the latter does not. Lastly,  $L_6$  is similar to  $L_5$  in that one message is true of the same state but differs from it in assigning upper-bounded semantics to  $m_{\text{some}}$ .

**Instead of this: Give an example of one of our target lexica, one of our competing (functional) lexica and one that is functionally deficient** Combining a linguistic behavior with each of these 6 lexica yields a total of 12 distinct types. , such as  $L_4$ , will produce speaker behavior that is *almost* indistinguishable from that of a language that lacks upper-bounds, but with pragmatic users, 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. Due to this possibly marginal difference between pragmatic  $L_4$  and  $L_5$ , the selection of one type over the other is expected to mainly depend on how well each can be transmitted to new learners. Things are less clear for literal  $L_5$  contrasted with literal/pragmatic  $L_4$ . The former has a learning advantage when learners are biased against upper-bounds but is expected to fare worse in communicative terms.

The dynamic is initialized with an arbitrary distribution over types, constituting the population’s first generation. The results for each parameter setting were obtained from 1000 independent runs, each consisting of 20 generations. This corresponds to a developmental plateau after which no noteworthy change was registered. As specified in §2.4, the mutation matrix  $Q$  can be obtained by considering all possible state-message sequences of length  $k$ . Given that this is intractable for large  $k$ , matrices were approximated by sampling 10 sequences from each type’s production probabilities and a type’s children being exposed only to this subset.

Drawing from our preceding discussion, 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 our attention on three pragmatic types.<sup>2</sup> Pragmatic  $L_3$ , a type that is lacking in expressivity but is a priori preferred for its lack of upper-bounds. Pragmatic  $L_4$ , a type that is functionally advantageous but biased against. And pragmatic  $L_5$ , combining virtues of the latter two.

**Expressivity only.** The replicator dynamic is sensitive to  $\lambda$  and  $\alpha$  as both have a bearing on a type’s fitness. The influence of  $\lambda$  is depicted in Figure 1.A. The less expressive  $L_3$  speakers fare the worst and are influenced the least by change in  $\lambda$ . In contrast, low values of  $\lambda$  result in a higher proportion of  $L_4$  speakers relative to  $L_5$ . This is expected given the central role of rationality in producing more deterministic behavior in users of  $L_5$ -like languages. Consequently, as the rationality parameter increases the functional difference between  $L_4$  and  $L_5$  is leveled. Overall, the outcome from only a pressure towards expressivity approximates an even share of pragmatic  $L_4$ ,  $L_5$  and  $L_6$  types. The latter follows the same trajectory as  $L_5$  in Figure 1.A. This illustrates an important issue in the evolution of meaning at the semantics-pragmatics interface:

<sup>2</sup>Pragmatic reasoning allows language users to refine their (possibly erroneous) choices. Therefore, it is advantageous even for types that codify more lexically.

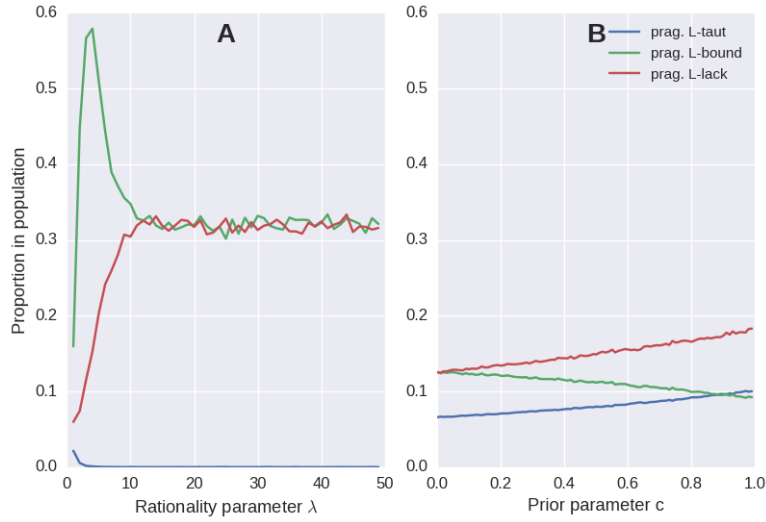


Figure 1: Mean proportions of target types after 20 generations in 1000 populations with only a pressure for expressivity in A ( $\alpha = 1$ ) and only for learnability in B ( $\alpha = 1, \lambda = 30, k = 5, l = 1$ ).

expressivity alone can not differentiate between (near-)functional equivalents to a degree that justifies the systematic prevalence of  $L_5$ -like semantics.

**Learnability only.** The effect of iterated learning with posterior sampling but without a pressure for expressivity is shown in Figure 1.B. In line with our expectations, the share of  $L_4$  speakers decreases as the bias against upper-bounds increases. In turn, this benefits  $L_3$  and, in particular,  $L_5$ . However, even a strong bias against lexical upper-bounds leads only to a moderate advantage of  $L_5$  over  $L_4$ . More importantly, a pressure only towards learnability can promote functionally defective languages such as  $L_3$ .

Inspecting these pressures separately not only showcases the contribution of the model’s components but also highlights some of their broader implications. First and foremost, neither dynamic on its own comes close to converging to a monomorphic population under most parameter configurations. For instance, while  $L_4$  speakers can come to take over a substantial proportion of the population, this only happens in a restricted range of low degrees of rationality. Apart from polymorphy, both pressures make undesirable predictions in isolation. A pressure only towards expressivity leads to the expulsion of communicatively suboptimal  $L_1$ ,  $L_2$  and  $L_3$  from the population but can not explain the regular selection of  $L_5$ -like semantics over either of its functionally similar alternatives. 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  $L_3$ .

**Expressivity and learnability.** Figure 2 illustrates the effect of the learning bias for posterior sampling (2.A) and slightly more MAP-like learning (2.B) when pressured for both expressivity and learnability. More detailed results for all types across a sample of  $c$ -values for  $l = 1$  and  $l = 5$  are presented in Table 3. These results show that a weak bias is sufficient to lead to a selection of  $L_5$  over  $L_4$ . As when only learnability is considered, this effect increases with the bias’ strength, provided  $L_5$  users are pragmatic. Importantly, the addition of a pressure towards

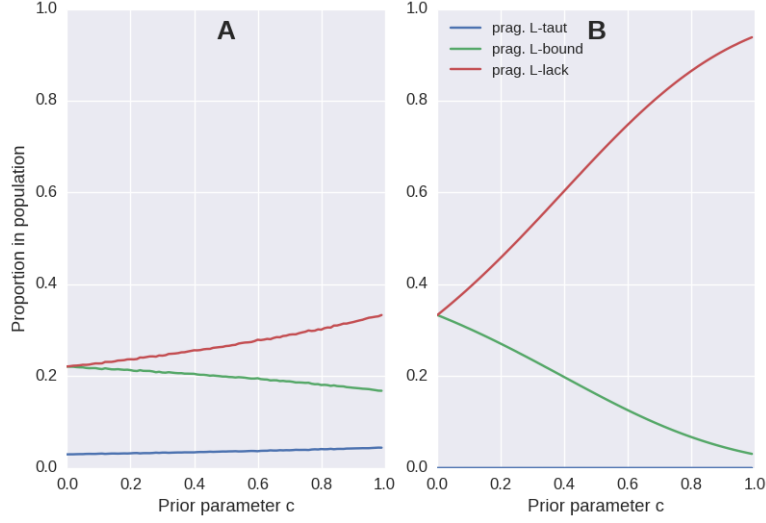


Figure 2: Mean proportions of target types after 20 generations in 1000 populations across bias values  $c \in [0, 1]$  with  $l = 1$  in A and  $l = 5$  in B ( $\alpha = 1, \lambda = 20, k = 5$ ).

expressivity magnifies this effect and dampens the proliferation of functionally suboptimal types advantaged by the learning bias. As stressed above, this indicates that neither a learning bias nor functional pressure alone but their combination may lead to the lack of upper-bounds in the lexical meaning of scalar expressions.

Other than the involvement of both pressures, the resulting proportion of pragmatic  $L_5$  speakers primarily hinges on three factors. First, the degree to which linguistic behavior is deterministic. This plays a role both for expressivity as well as in producing data that allows learners to discriminate this type from others. Second, the inductive bias, which controls the learners preference for simpler lexical representations. Lastly, the posterior parameter, as it magnifies the effects of the learning bias in tandem with replication.

As showcased by Figure 2.A, posterior sampling can lead to the incumbency of pragmatic  $L_5$ . However, not even a strong favorable learning bias combined with a pressure for expressivity completely drives out competing types under this inference strategy. This is not so for stronger posterior maximizing learning. As shown in Figure 3, the range of bias values within which  $L_5$  takes over the population increases with MAP-like learning. Put differently, the strength of the learning bias required for a given final proportion of  $L_5$  speakers strongly depends on learners' inferential strategy. As for the effect of the other parameters not mentioned so far, changes in sequence length influence the population in a predictable way: smaller values lead to more heterogeneous populations that reflect the learner's prior more faithfully. Larger ones lead to more pronounced differences among equally preferred types. This is due to the fact that the likelihood that a small sequence was produced by any type is relatively uniform (modulo prior) compared to that of types with lexica  $L_1 - L_3$  to produce larger sequences with the same state-message combination in contrast to pragmatic speakers of  $L_4 - L_6$ , or literal  $L_4$ . [TB: We need to discuss  $\alpha$  here as well if noisy perception is taken out]

	$l = 1$					$l = 5$				
$c$	0	.1	.5	.8	.9	0	.1	.5	.8	.9
lit. $L_1$	.03	.03	.04	.04	.04	0	0	0	0	0
lit. $L_2$	.03	.03	.02	.01	.04	0	0	0	0	0
lit. $L_3$	.03	.03	.04	.04	.04	0	0	0	0	0
lit. $L_4$	.07	.07	.06	.06	.05	0	0	0	0	0
lit. $L_5$	.04	.05	.05	.06	.06	0	0	0	0	0
lit. $L_6$	.04	.04	.04	.04	.03	0	0	0	0	0
prg. $L_1$	.03	.03	.04	.04	.04	0	0	0	0	0
prg. $L_2$	.03	.03	.02	.01	.04	0	0	0	0	0
prg. $L_3$	.03	.03	.04	.04	.04	0	0	0	0	0
prg. $L_4$	.22	.22	.2	.18	.17	.33	.30	.16	.07	.05
prg. $L_5$	.22	.23	.27	.3	.32	.33	.39	.68	.86	.90
prg. $L_6$	.22	.22	.2	.18	.17	.33	.30	.16	.07	.05

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

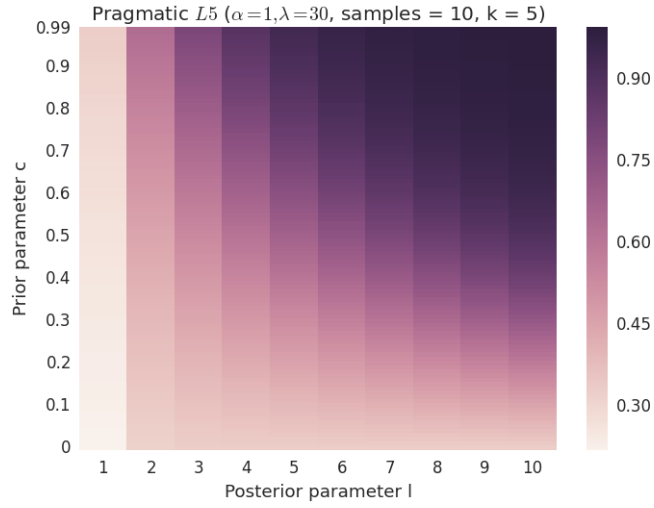


Figure 3: Mean proportion of pragmatic  $L_5$  in 1000 populations after 20 generations ( $\alpha = 1, \lambda = 30, k = 5$ )

### 3.2 Discussion

Under the assumption of a learning bias for simpler lexical representations, our results suggest that a lack of semantic upper-bounds coupled with pragmatic reasoning can overcome communicative pressures and stabilize in a population. This prediction hinges on three assumptions. First, that language is pressured toward both expressivity and learnability. Second, that language use is relatively deterministic. Lastly, that learners prefer simpler over more complex lexical representations. An important addendum to this third condition being that a combination of rationality in choice and maximization in learning requires a weaker bias towards simplicity. Under these conditions the selection of lexical meanings lacking upper-bounds in populations of pragmatic speakers is robust against parameter perturbations. This outcome is particularly encouraging in light of other advantages a lack of semantic upper-bounds may confer, as discussed at the beginning of Section 3.

A lack of upper-bounded in the lexical meaning of weak scalar expressions constitutes the majority view in the literature. However, it not clear to what extent other types should be present in the final population. It seems reasonable to expect functionally suboptimal types  $L_1$ ,  $L_2$  and  $L_3$  to be ruled out because they fail to enable their users to communicate effectively. However, this is not true of  $L_4$ .<sup>3</sup> The prediction that natural language communities are homogeneous or that a single speaker may entertain  $L_4$ -like semantics for one scalar expression and  $L_5$ -like semantics for another is not implausible (cf. Franke and Degen 2016). Alternatively, a stronger tendency for posterior maximization has to be assumed (see Figure 3). This empirical issue relates to other two aspects left undiscussed: disadvantages of pragmatic reasoning and the effect of state frequencies on the fossilization of pragmatic inferences. 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. Should cost play a role, then its effect is bound to depend on the frequency with which a given scalar expression is used. It is therefore plausible that frequently drawn scalar implicatures might fossilize to avoid cost, while infrequent ones are still derived on-line. This also opens a possible venue to address the preceding question about the expected presence of  $L_4$ -like semantics, but further empirical evidence is needed to assess these matters beyond speculation.

## 4 General discussion

We laid out a model that combines game theoretical models of functional pressure towards efficient communication (Nowak and Krakauer 1999), effects of transmission perturbations 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 and, more generally, effects of communicative pressures on the cultural evolution of language.

We argued that the puzzle raised by semantics in light of pragmatics is hard to explain on purely functional grounds and that part of the answer may instead lie in the way transmission

---

<sup>3</sup> $L_6$  presents a special case. In our current setup, it mirrors  $L_5$  in enabling for the pragmatic strengthening of a message that does not codify an upper-bound lexically. However, this is achieved by ruling out  $s_{\exists \rightarrow \forall}$  and not, as with scalar implicatures,  $s_{\forall}$ .  $L_6$  speakers therefore strengthen a “some”-message to convey something paraphrasable as ‘some but not [some but not all]’. The current representation of lexica as Boolean matrices is blind to this anomaly.



shapes the outcome of cultural evolution in tandem with a pressure for successful information transfer. In the realm of inductive biases, we adopted the assumption that simpler semantic representations are more likely to be learned (cf. Chater and Vitányi 2003). Under this view, semantics and pragmatics play a synergic role in that representational simplicity is supplemented by pragmatic reasoning to counteract functional disadvantages otherwise incurred. As a consequence, iterated transmission and use of language lead to a regularization that may explain the lack of lexicalization of systematic pragmatic enrichments. This result is of particular relevance for the longstanding assumption of a divide and interaction between semantics and pragmatics. It offers an account of why (certain) pragmatic inferences fail to lexicalize. More generally, we showed that systematic noise in perception can produce outcomes that are similar from those generated by inductive biases.

The main innovations of the model are its modular separation of expressivity and learnability, allowing for their isolated and combined analysis, the learning process involving a joint inference over types of pragmatic behavior and lexical meaning, as well as in its accommodation of different transmission perturbations that go beyond learning biases. The goal to decouple but model both expressivity and learnability has also recently been addressed by Kirby et al. (2015). In contrast to our proposal, Kirby et al. model expressivity as exerting its force only in the production of learning data. This model’s expressivity parameter thereby fulfills a similar role to high values of  $\lambda$  in making speaker behavior more deterministic. In this way, it “favors” unambiguous languages. However, the degree of mutual understanding of interlocutors central to replication and to our notion of expressivity is not taken into consideration. That is, while our proposal combines bidirectional horizontal transmission with its vertical and unidirectional counterpart, Kirby et al.’s model only considers the latter’s influence. Our reasoning behind the inclusion of the former lies in the empirical and theoretical observation that learnability alone can lead the selection of functionally defective languages, as illustrated by the tautological language  $L_3$  in our analysis. This outcome has been reported in a number of laboratory experiments where the participants’ task was to learn and subsequently reproduce the language produced by a previous participant, leading to a proliferation of languages that associated a large number of meanings with a single form (see e.g. Silvey et al. 2014 and experiment 1 in Kirby et al. 2008). In contrast, experiments involving an interactive component have been found to foster languages that enable interlocutors to distinguish meanings more accurately (e.g. Fay and Ellison 2013; for a review of laboratory results under the iterated learning paradigm and further discussion see Kirby et al. 2015, Tamariz and Kirby 2016). It is not evident how to compare these empirical findings given that they consider distinct meaning spaces, modes of transmission, iterations and feedback given to participants. However, we take these results to suggest that there is an important difference between a language generating learnable linguistic data and its actual performance as a means of information transfer. The former solely depends on the mechanism by which speakers associate form and meaning. The latter additionally hinges on the addressee’s linguistic experience and her ability to interpret linguistic input based on this experience. In sum, we contend that successful information transfer in a linguistic community is central to the adoption of a communication system and that this measure is not adequately reflected by production alone.

The demonstration that noise can lead to regularized evolutionary outcomes that are similar to those generated by prior learning biases is relevant not for the case study at hand, but more so for the broader project of investigating the cultural evolution of language. On the one hand, the plurality of sources of transmission perturbations admitted by these models paints a cautionary tale for the design of studies that purport to provide explanatory accounts of linguistic phenomena. In particular when the outcome is interpreted as being informative about the perturbation assumed to generate it (cf. Tamariz and Kirby 2016). On the other, and most importantly, it showcases how regularities can arise as a byproduct of systematic noise rather

than from standardly assumed inductive biases.

## 5 Conclusion

The cultural evolution of language is influenced by intertwined pressures. We set out to investigate this process by putting forward a model that combines a pressure toward efficient and successful information transfer with perturbations that may arise in the transmission of linguistic knowledge in acquisition. Additionally, we argued for the necessity of considering the role of pragmatics in investigations on the cultural evolution of meaning. These components and their mutual influence were highlighted in a case study on the lack of lexical upper-bounds in weak scalar expressions that showed that, when pressured for learnability and expressivity, the former drives for simpler semantic representations inasmuch as pragmatics can compensate for lack of expressivity in use. That is, the relative learning advantage of simpler semantics in tandem with a functional pressure in use may offer an answer to why natural languages fail to lexicalize systematic pragmatic inferences.

We also considered an alternate instantiation of the model, which shows that systematic noise in state perception can give rise to evolutionary outcomes that are similar to those predicted by inductive biases. This stresses the fact that that learning and typology are not necessarily close reflections of each other (Bowerman 2010). In particular, language use and environmental factors can play an important role in language change, making them central variables in explanatory accounts of natural language properties.

## References

- 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.
- Andrea Baronchelli, Andrea Puglisi, and Vittorio Loreto. Cultural route to the emergence of linguistic categories. *PNAS*, 105(23):7936–7940, 2008. doi: 10.1073/pnas.0802485105.
- John Batali. Computational simulations of the emergence of grammar. In James R. Hurford, Michael Studdert-Kennedy, and Chris Knight, editors, *Evolution of Language: Social and Cognitive Bases*. Cambridge University Press, Cambridge, UK, 1998.
- A.W. Beggs. On the convergence of reinforcement learning. *Journal of Economic Theory*, 122(1):1–36, 2005. doi: 10.1016/j.jet.2004.03.008.
- Anton Benz, Gerhard Jäger, and Robert van Rooij. *An Introduction to Game Theory for Linguists*, pages 1–82. Palgrave, 2005.
- 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.
- Leon Bergen, Roger Levy, and Noah D Goodman. Pragmatic reasoning through semantic inference. *Semantics and Pragmatics*, 2016.
- Andreas Blume, Yong-Gwan Kim, and Joel Sobel. Evolutionary stability in games of communication. *Games and Economic Behavior*, 5(5):547–575, 1993. doi: 10.1006/game.1993.1031.

- Tilman Börgers and Rajiv Sarin. Learning through reinforcement and replicator dynamics. *Journal of Economic Theory*, 77(1):1–14, 1997. doi: 10.1006/jeth.1997.2319.
- Lewis Bott and Ira A. Noveck. Some utterances are underinformative: The onset and time course of scalar inferences. *Journal of Memory and Language*, 51(3):437–457, 2004.
- Melissa Bowerman. *Linguistic Typology and First Language Acquisition*. Oxford University Press, 2010. doi: 10.1093/oxfordhb/9780199281251.013.0028.
- Henry Brighton. Compositional syntax from cultural transmission. *Artificial Life*, 8(1):25–54, 2002. doi: 10.1162/106454602753694756.
- 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.
- Gennaro Chierchia, Danny Fox, and Benjamin Spector. Scalar implicature as a grammatical phenomenon. In Claudia Maienborn, Klaus von Stechow, and Paul Portner, editors, *Semantics. An International Handbook of Natural Language Meaning*, pages 2297–2332. de Gruyter, Berlin, 2012.
- Eve V. Clark. Lexical meaning. In Edith L. Bavin, editor, *Child Language*, pages 283–300. Cambridge University Press, New York, 2009. doi: 10.1017/CBO9780511576164.016.
- Judith Degen and Michael K. Tanenhaus. Processing scalar implicatures: A constraint-based approach. *Cognitive Science*, 39:667–710, 2015. doi: 10.1111/cogs.12171.
- Nicolas Fay and T. Mark Ellison. The cultural evolution of human communication systems in different sized populations: Usability trumps learnability. *PLoS ONE*, 8(8), 2013. doi: 10.1371/journal.pone.0071781.
- Jacob Feldman. Minimization of boolean complexity in human concept learning. *Nature*, 407(6804):630–633, 2000.
- Michael C. Frank and Noah D. Goodman. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012.
- Michael Franke. *Signal to Act: Game Theoretic Pragmatics*. PhD thesis, University of Amsterdam, 2009.
- Michael Franke. The evolution of compositionality in signaling games. *Journal of Logic, Language and Information*, 25(3–4):355–377, 2016. doi: 10.1007/s10849-015-9232-5.
- Michael Franke and Judith Degen. Reasoning in reference games: Individual- vs. population-level probabilistic modeling. *PLOS ONE*, 11(5):e0154854, 2016. doi: 10.1371/journal.pone.0154854.
- Michael Franke and Gerhard Jäger. Pragmatic back-and-forth reasoning. *Semantics, Pragmatics and the Case of Scalar Implicatures.*, pages 170–200, 2014.
- Michael Franke and Gerhard Jäger. Probabilistic pragmatics, or why bayes’ rule is probably important for pragmatics. *Zeitschrift für Sprachwissenschaft*, 35(1):3–44, 2016. doi: 10.1515/zfs-2016-0002.
- Gerald Gazdar. *Pragmatics, Implicature, Presupposition and Logical Form*. Academic Press, New York, 1979.

- Bart Geurts. *Quantity Implicatures*. Cambridge University Press, Cambridge, UK, 2010.
- Tao Gong. *Language Evolution from a Simulation Perspective: On the Coevolution of Compositionality and Regularity*. Chinese University of Hong Kong, 2007.
- Noah D. Goodman and Michael C. Frank. Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11):818–829, 2016. doi: 10.1016/j.tics.2016.08.005Au.
- Noah D. Goodman and Andreas Stuhlmüller. Knowledge and implicature: Modeling language understanding as social cognition. *Topics in Cognitive Science*, 5:173–184, 2013. doi: 10.1111/tops.12007.
- Paul Grice. Logic and conversation. In *Studies in the Ways of Words*, chapter 2, pages 22–40. Harvard University Press, Cambridge, MA, 1975.
- Thomas L. Griffiths and Michael L. Kalish. Language evolution by iterated learning with bayesian agents. *Cognitive Science*, 31(3):441–480, 2007.
- Daniel J. Grodner, Natalie M. Klein, Kathleen M. Carbary, and Michael K. Tanenhaus. “some,” and possibly all, scalar inferences are not delayed: Evidence for immediate pragmatic enrichment. *Cognition*, 166:42–55, 2010.
- Dirk Helbing. A stochastic behavioral model and a ‘microscopic’ foundation of evolutionary game theory. *Theory and Decision*, 40(2):149–179, 1996. doi: 10.1007/BF00133171.
- Julia Hirschberg. *A Theory of Scalar Implicature*. PhD thesis, University of Pennsylvania, 1985.
- Josef Hofbauer. The selection mutation equation. *Journal of Mathematical Biology*, 23:41–53, 1985.
- Josef Hofbauer and Karl Sigmund. Evolutionary game dynamics. *Bulletin of the American Mathematical Society*, 40(04):479–520, 2003.
- Laurence R. Horn. *On the Semantic Properties of Logical Operators in English*. Indiana University Linguistics Club, Bloomington, IN, 1972.
- 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.
- Yi Ting Huang and Jesse Snedeker. Online interpretation of scalar quantifiers: Insight into the semantics–pragmatics interface. *Cognitive Psychology*, 58(3):376–415, 2009.
- James R. Hurford. Biological evolution of the saussurean sign as a component of the language acquisition device. *Lingua*, 77(2):187–222, 1989.
- Simon M. Huttegger. Evolution and the explanation of meaning. *Philosophy of Science*, 74:1–27, 2007. doi: 10.1086/519477.
- Gerhard Jäger. Evolutionary game theory and typology: A case study. *Language*, 83(1):74–109, 2007. doi: 10.2307/4490338.
- 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.

- 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.
- C. Kemp and T. Regier. Kinship categories across languages reflect general communicative principles. *Science*, 336(6084):1049–1054, 2012. doi: 10.1126/science.1218811.
- 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.
- S. Kirby, H. Cornish, and K. Smith. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. *Proceedings of the National Academy of Sciences*, 105(31):10681–10686, 2008. doi: 10.1073/pnas.0707835105.
- 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.
- Simon Kirby and James R. Hurford. The emergence of linguistic structure: An overview of the iterated learning model. In A. Cangelosi and D. Parisi, editors, *Simulating the Evolution of Language*, pages 121–148. Springer, 2002.
- 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.
- Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. Compression and communication in the cultural evolution of linguistic structure. *Cognition*, 141:87–102, 2015.
- Tom Lenaerts, Bart Jansen, Karl Tuyls, and Bart De Vylder. The evolutionary language game: An orthogonal approach. *Journal of Theoretical Biology*, 235:566–582, 2005.
- Stephen C. Levinson. *Pragmatics*. Cambridge University Press, Cambridge, UK, 1983.
- David Lewis. *Convention: A Philosophical Study*. Harvard University Press, Cambridge, 1969.
- Duncan R. Luce. *Individual choice behavior: a theoretical analysis*. Wiley, 1959.
- André Martinet. *Functionalist View of Language*. Clarendon Press, Oxford, 1962.
- Wim De Neys and Walter Schaeken. When people are more logical under cognitive load. *Experimental Psychology*, 54(2):128–133, 2007.
- 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. *Evolutionary Dynamics: Exploring the Equations of Life*. Harvard University Press, 2006.
- 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. Evolution of universal grammar. *Science*, 291(5501):114–118, 2001. doi: 10.1126/science.291.5501.114.

- 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.
- Mark Pagel. Human language as a culturally transmitted replicator. *Nature Reviews Genetics*, 10:405–415, 2009. doi: 10.1038/nrg2560.
- Steven T Piantadosi. Zipf’s 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.
- 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, 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.
- 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.
- Robert van Rooij and Tikitou de Jager. Explaining quantity implicatures. *Journal of Logic, Language and Information*, 21(4):461–477, 2012. doi: 10.1007/s10849-012-9163-3.
- William H. Sandholm. *Population Games and Evolutionary Dynamics*. MIT Press, Cambridge, MA, 2010.
- Uli Sauerland. Scalar implicatures in complex sentences. *Linguistics and Philosophy*, 27:367–391, 2004. doi: 10.1023/B:LING.0000023378.71748.db.
- Karl H. Schlag. Why imitate, and if so, how? *Journal of Economic Theory*, 78(1):130–156, 1998. doi: doi:10.1006/jeth.1997.2347.
- Catriona Silvey, Simon Kirby, and Kenny Smith. Word meanings evolve to selectively preserve distinctions on salient dimensions. *Cognitive Science*, 39(1):212–226, 2014. doi: 10.1111/cogs.12150.
- Brian Skyrms. *Signals: Evolution, Learning, and Information*. Oxford University Press, Oxford, 2010.
- Kenny Smith, Simon Kirby, and Henry Brighton. Iterated learning: A framework for the emergence of language. *Artificial Life*, 9:371–386, 2003.
- Matthew Spike, Kevin Stadler, Simon Kirby, and Kenny Smith. Minimal requirements for the emergence of learned signaling. *Cognitive Science*, 2016. doi: 10.1111/cogs.12351.
- Luc Steels. A self-organizing spatial vocabulary. *Artificial Life*, 2(3):319–332, 1995. doi: 10.1162/artl.1995.2.3.319.
- Luc Steels. Modeling the cultural evolution of language. *Physics of Life Reviews*, 8(4):339–356, 2011.
- Luc Steels and Tony Belpaeme. Coordinating perceptually grounded categories through language: A case study for color. *Behavioral and Brain Sciences*, 28(4):469–529, 2005. doi: 10.1017/S0140525X05000087.

- Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1998.
- Monica Tamariz and Simon Kirby. The cultural evolution of language. *Current Opinion in Psychology*, 8:37–43, 2016.
- Peter D. Taylor and Leo B. Jonker. Evolutionary stable strategies and game dynamics. *Mathematical Bioscience*, 40(1–2):145–156, 1978.
- Bill Thompson, Simon Kirby, and Kenny Smith. Culture shapes the evolution of cognition. *Proceedings of the National Academy of Sciences of the United States of America*, 113(16):4530–4535, 2016. doi: 10.1073/pnas.1523631113.
- Elizabeth Closs Traugott. Historical pragmatics. In Laurence R. Horn and Gregory Wand, editors, *The Handbook of Pragmatics*, pages 538–561. Blackwell Publishing, 2004.
- Tessa Verhoef, Simon Kirby, and Bart de Boer. Emergence of combinatorial structure and economy through iterated learning with continuous acoustic signals. *Journal of Phonetics*, 43: 57–68, 2014. doi: 10.1016/j.wocn.2014.02.005.
- Karl Wärneryd. Cheap talk, coordination, and evolutionary stability. *Games and Economic Behavior*, 5(4):532–546, 1993. doi: doi:10.1006/game.1993.1030.
- George Zipf. *Human behavior and the principle of least effort*. Addison-Wesley Press, 1949.