

# Co-evolution of lexical meaning & pragmatic use

## **Revision Cover letter**

We would like to thank the reviewers for their helpful comments and suggestions which have certainly helped improve the manuscript. The main changes in response to the reviewers comments are the following.

- We have tried to improve the exposition of the model, both on the technical side and on the way the model can be interpreted (re reviewer 1’s concerns about how to interpret particular formulas and the model as a whole). There are now three very short appendix chapters, each with formal detail about how to better understand the abstract model. There is also the new section 2.3 dedicated to an explicit discussion of how to interpret the replicator mutator dynamic and a formal derivation of one particular agent-level update process which gives rise to this dynamic at the population level (appendix B).
- Where previously we only had one simulation-based case study, we now complement these results with a more in-depth case study, presented in the new section 3.1, which focuses on a subset of types present in the larger case study presented in sections 3.2 and 3.3. This helps tie in our results closer to known results from the literature (re reviewer 2’s comment 2) and also helps understand the inner workings of the model better for the case at hand. It also makes clearer in which sense we may speak of “co-evolution of semantics and pragmatics” (in response to reviewer 1’s comment).
- To compensate for the added material, we also tried to shorten the paper. We tried to avoid repetitive formulations, delegated more detailed explanations to footnotes and omitted certain passages that we feel are interesting, but are possibly also too distracting in a paper that is already quite long and also partly rather technical.

In what follows we address the reviewers’ main comments and questions one by one, with pointers to respective changes effected on the manuscript.

---

**Reviewer 1 comment 1***(co-evolution)*

the co-evolutionary dynamic mentioned in the title could be better justified [...] In co-evolutionary dynamics, the evolution of one factor (here, the lexicon) should have an impact on the evolution of another (here, pragmatic strategy), and vice versa; the paper does not show evolutionary dynamics (change over time) for either of these factors, not does it show how evolution of a single factor favours the evolution of the other. [...] Please justify your characterization of the process modeled as co-evolution or change to adaptation.

This is an important point that we now stress and highlight more throughout the paper. In particular, section 3.1 illustrates how one factor has an impact on the evolution of the other. In a nutshell: a pragmatic strategy (level-1 reasoning) favors the evolution of underspecified semantics of the target kind ( $L_{\text{lack}}$ ) as it allows for the maintenance of simpler lexical representations that are easier to learn without functional disadvantages otherwise incurred in communication (e.g., by level-0 reasoners). In the other direction, underspecified semantics of the target kind also favors the evolution of a pragmatic disposition to act on them. When paired with other semantics, pragmatic refinement can instead be (close to a) neutral trait (e.g.,  $L_{\text{bound}}$  paired with high  $\lambda$ ) or it can even lead to functional disadvantages when compared with level-0 reasoning (e.g.,  $L_{\text{bound}}$  with low  $\lambda$ ). In this way both features benefit from each other.

---

**Reviewer 1 comment 2***(modeling choices)*

the mapping between the model and the example of choice (scalar implicatures) could be better motivated

As we discuss in Section 1, our case study focuses on scalar implicatures because they are well-studied systematic pragmatic inferences that have received much attention, both theoretically and experimentally. This makes them suitable candidates for the study of the evolution of regular pragmatic inferences; past research can guide our modeling choices and aid in evaluating our application to the study of these inferences. This ties in with our

particular choices to model the pragmatic use of scalar expressions. As we now make more explicit in §2.3.1, these draw from a growing game-theoretic and Bayesian tradition (e.g., [Franke 2009](#), [Franke and Jäger 2016](#), [Goodman and Frank 2016](#)).

---

Reviewer 1 asks for a clearer description of the model. As noted earlier, we have taken care to make the exposition more accessible and to guide the reader more through the technical details of Section 2. In particular, the newly added sections 2.3, 3.1, and the appendix sections should help readers unravel the model and understand its conceptual interpretation and inner mechanisms much better than before. We used the following questions from Reviewer 1, briefly answered below, to guide the changes made in this new version.

How is generation transmission implemented?

As it is unlikely that lexical representations and/or pragmatic reasoning are inherited genetically, we model the transmission of linguistic knowledge across generations as iterated Bayesian learning. Technical details and motivations are given in §2.3 and appendix B.

Do agents in a population die little by little and are replaced by new agents? Or all at once?

This question clearly demonstrates a failure on our part to make clear enough how the replicator mutator dynamic can be interpreted. To answer the question posed here directly: all of the above possibilities *could* give rise to the replicator dynamic. It is a deliberate choice of abstract evolutionary modeling to *not* to have to commit to any particular model, when it is also possible to answer what would happen under *many* different regimes of fitness-based selection. To give a concrete interpretation of the population-level dynamic nonetheless, appendix B sketches the mean-field derivation of one form of an agent-level conditional imitation scheme.

Do the models consider just a single teacher and a single learner at a time (I don't think this is the case, since the results refer speak of "majority type, i.e., the type with the highest proportion in the final population")?

The parenthetical remark is correct. The amount of teachers of a given type is proportional to the frequency of this type in the population after the replicator step. The population is infinite so there are technically infinitely many learners/teachers. This was implicit in the definitions of the replicator mutator dynamic in (1) and (2), and is now more explicitly stated in sections 2.2 and 2.3.

How many agents were there per population? And per generation?

The prose surrounding definitions (1) and (2) in Section 2.2 and 2.3 now makes explicit that the population is infinite. This allows us to track change that does not depend on varying population sizes nor their growth rate. This is why we speak of proportions of a type in a population, rather than of there being a particular number of agents of a type.

Within one run of the model, do all agents in a population have the same  $\lambda$ ,  $l$  and  $k$ ? If so, please this homogeneous population, and discuss the possible consequences of variation in the population in these parameters.

While our model is compatible with the assumption that every agent comes with her own parameter values, the reviewer is correct in pointing out that we assume they all to have the same values (explicit in §3.3). Variation across agents would mean having many more types: one for each combination of a possible lexicon, pragmatic strategy,  $\lambda$ -,  $l$ - and  $k$ -value. This would allow us to trace the evolutionary trajectory of not only lexica and pragmatic strategies, which is our goal here, but of all these factors combined. As mentioned in Section 4, this would be interesting future research, but we cannot go into this here for reasons of space and a huge increase of complexity in the resulting model.

Does the model assume that learners observe not only the behaviour (the message, e.g. agent A hearing agent B say: "I own some of JC's albums") but also the true state of the world (i.e. agent A also knows whether agent B owns all or some-but-not-all albums)? Or do they get

feedback on their communicative success? In other words, how do they learn the correct state-message mappings that form the lexicon?

Learners witness sequences of data  $d$ . Such a sequence is composed of message-state pairings. In other words, learners do witness the true state of the world (Section 2.4.3, *learnability*). This, however, is a widely made assumption in the literature on iterated learning.

Additionally, some interpretations of the replicator mutator dynamic would require agents to have some rudimentary access to how good their own and/or other agents' behavior is. For example, the conditional imitation scheme sketched in appendix B assumes that agents imitate with a probability that depends on how good the possibly to be imitated behavior is. A biological interpretation or an interpretation in terms of reinforcement learning, for example, would not require this. In other words, even if this assumption is felt to be too strong, it is not necessary to discredit the model proposed here, because it is compatible with many agent-level update schemes, some of which are very innocuous when it comes to what each agent knows or sees or learns.

How exactly are pragmatic strategies inherited/inferred by learners? I think all of this may be in the formulas, but it would be useful if it was expressed in prose too.

In short, a type is a combination of a lexicon and a disposition to act on it. Together, both ingredients define an agent's linguistic behavior. A learner faced with input  $d$  (see above) calculates  $P(\tau \mid d)$ , the probability that a type  $\tau$  produced input  $d$ . Crucially, neither linguistic strategies nor lexical meanings are directly observable, so they can only be faithfully recovered inasmuch as the *overt* behavior evidenced by a type – the data it produces – is (in tendency) attributable to only one, *covert*, pairing of a lexicon and a linguistic strategy.

Could you express in prose how literal and pragmatic use work in practice? I am not sure I understand why you call the strategies 'pragmatic', given that pragmatics depends on the context, and there is no context here. I would like to see the use of the term 'pragmatic' justified further in the paper.

We follow the Gricean tradition of viewing pragmatic inference as effected by mutual reasoning about rational language use (Section 1). This may involve

the recruitment of contextual information, as the reviewer notes, but can also involve pragmatic enrichments that result purely from reasoning about linguistic alternatives (“why did the speaker say *some* instead of *all*”). It is in the latter sense that reasoning beyond level 0 is pragmatic (see §2.3.1). Classically, theoretical analyses of scalar implicatures have put their emphasis on reasoning about linguistic choice rather than on contextual information (e.g., Horn 1972, Gazdar 1979, Franke 2009, Goodman and Stuhlmüller 2013). We chose to do the same here.

I thought  $\lambda$  was the parameter that modeled communication, so it should not be present in the iterated-learning-without-communication models; why is  $\lambda$  set to 20 here, and why discuss lower levels of lambda?

This parameter regulates the linguistic production behavior of types, i.e., speaker behavior. Speaker behavior is relevant both for communication with hearers as well as for the learning input that a teacher produces. When considering only iterated Bayesian learning  $\lambda$  accordingly also plays a role because different  $\lambda$  values yield different likelihoods of producing particular data. Section 3.2 and appendix B zoom into this topic and make this point more explicit. We regret that this was not ideally transparent in the previous version.

Also [in the Subsection *learnability only*], why isn’t the proportion of  $L_{all}$  highest, given that it has the lowest complexity (highest prior)?

[MF: We need to possibly rewrite this.] Due to their stochastic speaker behavior, the data produced by  $L_{all}$ -teachers tends to be compatible with the behavior of many other types. Intuitively, the data such teachers produce are all over the place and do not do a good job in setting them apart from other types. Consequently, even if the prior favors  $L_{all}$ , this type is not transmitted very faithfully; learners reason that the input they get from  $L_{all}$ -teachers could also come from other types and may adopt those instead. If there are types that are transmitted more faithfully, the population will, over time, transition to these types rather than to  $L_{all}$ . As we put it in Section 4, iterated learning does not necessarily promote the *a priori* more likely type, but tends to promote a type  $t$  based on a gradient of how many other types might likely mutate into  $t$ , so to speak. We stress this in Section 3.2.2 and Section 4.

---

**Reviewer 1 comment 4***(LOT & model relation)*

This model rests heavily on the assumption that the representation or meaning of an expression is a logical formula like those in Table 2, and that these formulas yield a complexity hierarchy. This may be valid for scalar implicature, as you argue, but how realistic is it cognitively for other linguistic structures or for other socially learned items? Should we assume the non-parsimonious view that there are different complexity measures for different linguistic structures? It would be nice to see another linguistic or cultural example that could be fitted with this model.

Learnability [...] is operationalized as the complexity of the minimal logical formula for each state. Another possible measure of the difficulty/ease of learning would be comparing the regularity/systematicity of the lexicon types given in 3.1.1 (measured e.g. as mutual information (Cornish, Tamariz, Kirby 2009; using the mantel test, Kirby, Tamariz, Cornish, Smith 2015)). Yet another approach to complexity would be to assume that the meanings are mutually exclusive categories without internal structure, or to assume a preference or mutual exclusivity. In these cases, I expect L-bound (competitor), to be easier to learn than L-lack (target). Please justify your choice of learnability measure in the face of these alternatives.

First of all, the model does not rest on this assumption at all; only the particular application does. Following previous work in the iterated learning tradition, we do assume that (i) inductive learning biases can play a role in shaping culturally transmitted knowledge and that, as the reviewer notes, (ii) biases can be of manifold nature.

On the other hand, it is true that the model itself poses constraints on what might count as a plausible inductive bias. One well-studied example of an inductive bias from the iterated learning tradition is the mutual exclusivity bias. This bias should be expected to work against targets and in favor of competitors, as the reviewer says. Yet in the context of our model a bias for mutual exclusivity might be hard to motivate, because the model itself has an explicit component of fitness-based selection for communicative efficiency. This component should take the part of selecting for mutual exclusivity. In other words, the division of labor is shifted here: not every constraint is as

easily justifiable; the model tries to give an explanation in terms of fitness-based selection for some of the effects that certain biases are meant to explain in systems which only look at iterated learning.

What counts for this model is chiefly a bias that is *not* motivated by the effects it has on communicative efficiency. Instead we only focus on the effects of a single plausible contributing factor to a preference of target lexica over competitors: a well-motivated bias favoring representational simplicity of single concepts, not lexica as a whole (Feldman 2000, Chater and Vitányi 2003, Piantadosi et al. 2012, Kirby et al. 2015, Piantadosi et al. view). How to spell out such a lexical simplicity bias is, we readily admit, an open question. As mentioned in Section 3.1.2, the complexity measure in terms of a Language of Thought is just a convenient operationalization of this particular hypothesis for the sake of a concrete working example.

We should stress that, just as models of iterated learning make different predictions under the assumption of different inductive biases, so does our model. The assumption and operationalization of a particular bias should, of course, always be seen critically. We criticize our own setup and assumptions in Section 4.

---

**Reviewer 1 comment 5**

(*transmission fidelity*)

What is your measure of fidelity/infidelity? Is it just 0 if the teacher and learner are of a different type and 1 if they are of the same? Could you have a more graded measure of fidelity, e.g. taking into account if two types are of the same 'kind' or not?

We understand transmission fidelity as the probability of acquiring type  $i$  when learning from type  $i$  (Section 2.4.3). This is the value of a cell in the learning matrix  $Q$ . In this case,  $Q_{ii}$ . As discussed in Section 2.2, if  $Q_{ii} = 1$  then  $i$  is always acquired when learning from type  $i$ . On the other extreme,  $Q_{ii} = 0$  means that  $i$  is never acquired from  $i$ .  $Q$  is a stochastic matrix and we usually see neither extreme but rather, as suggested by the reviewer, a graded notion of transmission fidelity (Section 2.2). The details on how  $Q$  is computed are given in Section 2.4.3. We have expanded this section to make these details clearer in the prose.

---



**Reviewer 2 comment 1***(Interpretation of results)*

The results are interesting, though somewhat difficult to interpret. First, pragmatic language use does not evolve in the absence of learnability pressures, due to its slightly lower fitness. Second, pragmatic language use can arise from learnability pressures alone, though only for sufficiently “rational” learners. Third, pragmatic language use can evolve given both fitness and learnability pressures, though only if both learners and language speakers are both sufficiently optimal. The first two results are fairly straightforward, but the third (which is the most important in the paper) is more puzzling. Why do learnability pressures not have the same effect in the joint learnability/fitness model, as they do in the lesioned model? Why is high speaker optimality necessary for scalar implicatures to evolve, when it does not have this effect in the fitness-only model?

We have thoroughly revised the manuscript to make the contrast between these three predictions clearer (see Section 3.3, particularly 3.3.3; but also Section 3.1, which illustrates how the dynamics behave in a smaller type space; as well as the discussion of iterated learning alone and its dependence on  $\lambda$  and  $l$  in Sections 2.4.3 and appendix C). These questions tie in with Reviewer 1’s request for a clearer exposition and the explanation of why we speak of co-evolution in this case (reviewer 1, comment 1). In a nutshell, pressure for learnability alone does not put the types in a population in competition. Instead, the population comes to be inhabited by types that are recovered from learning input of replicated types more often – of which there might be many (e.g., all variants of a kind are equally learnable; if there are many variants of a type that is, in tendency, acquired more often, then the population will stay polymorphic). What fitness-relative selection adds to this process is the competition that learnability alone lacks, independent of the type space we look at. Now, in the case of our type space we have a slight (modulo  $\lambda$ ) communicative disadvantage for targets relative to competitors, for instance. However, the fact that targets are inferred more often by naïve learners (by asymmetries in production likelihoods, see appendix C; a tendency which is also mediated by  $\lambda$  and  $l$ ), leads (i) to the gradual existence of less competitors, counteracting this functional disadvantage, as well as (ii) to the existence of only one variant of a kind (contra what we get in the *learning alone* condition).

**Reviewer 2 comment 2***(Model predictions)*

I agree with its claim that the results are non-trivial; it is not clear a priori that there should exist any regimes where learnability and fitness pressures balance in order to produce scalar implicatures. The paper does not, however, address why the pressures do balance in this case. Are there general theorems (from previous work) about the replicator mutator dynamic that could be brought to bear here? How do learnability and fitness pressures trade-off in general? Is there anything more general that can be said about conditions under which pragmatic language use will evolve? Does the current example generalize to other cases of pragmatic language use? The paper does not need to answer these questions in an exhaustive manner, but greater scientific understanding of the model is desired.

As noted above, we have expanded our analysis of the model and put much more emphasis on how the components that feed it play a role in driving particular evolutionary outcomes. Section 3.1 can be seen a direct response to this comment, where we try to deduce from what is known from the literature as much insight as possible about the workings of the dynamic in this case. Unfortunately, owing to the complexity of the case study we are not able to offer a mathematical characterization of, say, the stable rest points of the dynamic. (Notice that for this purpose we should switch to the continuous time formulation, which would also extend the paper further.)

---

**Reviewer 2 comment 3**    *(co-evolution and pragmatic maintenance)*

It makes sense for agents to acquire their grammar/lexicon by learning, as these are social conventions. I believe, however, that the paper also assumes that an agent's pragmatic type (i.e. whether they are literal or pragmatic speakers) is also determined by learning. This is more difficult to interpret. It is reasonable for an agent to learn whether the other agents in a population are behaving pragmatically. Even if the other agents are not behaving pragmatically, however, it will be rational for the learner to behave pragmatically. This is a property of the recursive definition of the pragmatics model. The pragmatic speaker is defined so as to select utterances that will be correctly interpreted by the literal listener. If an agent learns that their interlocutors are all literal agents, then it would be rational for them to pragmatically reason about these

literal agents. It is therefore not clear why the agents should, as currently proposed, copy the inferred pragmatic type of their interlocutor.

It is not necessarily true that it will always be advantageous for agents to behave pragmatically (see the end of Section 2.4.1 but also the discussion surrounding Figure 3a). The main thing to note is that a pragmatic level-1 hearer using lexicon  $L_{\text{bound}}$  will reason about the behavior of a *soft-maximizing* level-0 speaker of  $L_{\text{bound}}$ . The lower  $\lambda$ , the more the stochasticity percolates from the level-0 speaker to level-1 pragmatic interpretation. By contrast, literal level-0 hearers of  $L_{\text{bound}}$  have a one-to-one form-meaning mapping from the start. Put differently, pragmatic reasoning can actually encumber some hearers, depending on their lexicon. Particularly if  $\lambda$  is low.

More to the point of the question, even if higher levels of pragmatic reasoning were always at least as good as the next lower level, the current population-level dynamic rests on *minimal* assumptions about agent-level rationality. We have expanded on this point, which was definitely unclear in the first version, in section 2.3. In particular footnote 1 is almost a direct reply to this excellent comment. So, our agents are deliberately *not* assumed to know the type of their interlocutor, but merely iterate pragmatic reasoning based on their own behavior. They would also not play a best response to the population average. The dynamic assumes that agents do not innovate strategies (like higher-order pragmatic reasoning). The dynamic is built on the idea that fitness-based selection is a non-innovative, gradient and proportional amplifier of whatever works well. Agents might never find out what works well, e.g., when fitness is just a measure of the proportion of offspring (a biological interpretation). Or they find out implicitly based on what works for them when they explore different ways of behavior (reinforcement learning). Or they occasionally realize that other agents are communicatively successful to a certain degree (e.g., conditional imitation based on success). Yet other interpretations are possible as well.

---

**Reviewer 2 comment 4**

(*agent simplicity*)

The paper assumes that the speaker (and listener) know the type of the player that they're interacting with. This does not seem like a natural assumption. Unless the speaker has repeated interactions with the listener (which is not assumed in the paper), the pragmatic type of the other agent is a latent property. One could imagine the speaker calcu-

lating expected utility with respect to the distribution on players that they expect (and speakers of different types choosing different strategies against this distribution).

Indeed, this was regrettably not clear enough in the previous version. As we now emphasize in Sections 2.3 and 2.4.1, we assume very little sophistication from our agents. In particular, they do not know the type of the player that they are playing against. They simply behave according to their subjective point of view in a boundedly rational fashion. For example, the linguistic choices of a level-1 speaker of lexicon  $L_{\text{some}}$  do not change depending on whom she interacts with. They are always defined as (boundedly) rational choice relative to the interpretative behavior of a level-0 user of  $L_{\text{some}}$  (see definitions (3)–(6) in §2.3.1).

---

**Reviewer 2 comment 5**

*(posterior type sampling)*

The learner in this model always selects a particular grammar/pragmatic type. An alternative would be for them to maintain uncertainty about the type (which is optimal behavior from a Bayesian perspective). I suspect that this model is intractable (or at least much less manageable) than the proposed one, but it would be worth noting that this is a substantive choice point in the model.

We now highlight that this is a design choice and that there are alternatives in Section 2.4.3, footnote 3.

---

**Reviewer 2 comment 6**

*(definitions in rational language use)*

The current pragmatic model is somewhat non-standard. In particular, neither the literal nor pragmatic speakers in Equations 3-6 use an information-theoretic utility function, as is standard in the literature. One consequence of this is that there is now an asymmetry between the literal speaker and listener: the literal listener never interprets utterances in a manner inconsistent with their literal meanings, while the literal speaker will sometimes use utterances in a non-literal way. I do not think this needs to be changed – the current utility function is perfectly sensible – but this is worth noting in the paper. It would also be

desirable to know (possibly in future work) whether the modeling results are robust to switching to information-theoretic utility functions.

A related issue is in the fitness definition in Section 2.3.2. If one adopts the alternate utility function above for the agents’ behavior, then it would be natural to have fitness scale in a logarithmic manner as well. In this case, the information-theoretic utility function would have the following interpretation: the fitness of a language is determined by how well it allows speakers to communicate their beliefs, rather than how often it leads the listener to make the correct guess about the world. Similar robustness questions apply here as above.

There are many alternatives when it comes to definitions of rational language use. The definitions that draw more strongly from a game-theoretic tradition standardly do not make use of an information-theoretic utility function (see [Qing and Franke 2015](#) for an overview and discussion up to 2015). As noted in our answer to Reviewer 1’s comment 2, we have taken care to highlight that our choices draw from this literature and to mention that there are other alternatives in the same spirit.

More concretely, there is a technical reason why we would like all speaker types to produce any message with positive probability: we do not want to make types too easily identifiable by even small sets of observed data, so that iterated learning has some work to do. Moreover, we believe that while a logarithmic definition of speaker utilities is just fine, when it comes to evolutionary fitness, a definition in terms of beliefs is not. Fitness should track genuine observable interaction with the world. Just holding a “good belief” gives no countable or observable advantage unless it is acted upon in some manner.

---

**Reviewer 2 comment 7***(Learnability)*

Throughout, the paper mostly equates the learnability of a grammar with its complexity/prior probability. In general, however, the relationship may be more complex. Certain grammars may require less data than others to be statistically identified. This is probably not an issue in the current case study, given the simplicity of the grammars and the large

amount of data relative to grammar size, but this is likely to be an issue in other cases.

The agents in this model do not actually learn the “grammar” of concepts which is used only as a tool for fixing one concrete measure of complexity of lexical concepts. It is these lexical concepts that are learned. We have tried to make this clearer by more careful reformulation.

---

**Reviewer 3 comment 1**

(*expected utility*)

what would be the impact of assuming that fitness is not contributed equally by being successful as speaker and being successful as a perceiver?

This is a standard assumption in the literature (lacking reasons to believe that there is an asymmetry in how often, on average, agents use and hear a type of expression.) In the case of scalar implicatures, we see no reason to assume such an asymmetric contribution. As for concrete impact: The main contrast that we are interested in (competition between targets and competitors) does not hinge on agents being speakers half of the time. However, e.g., the contrast between pragmatic and literal competitor types does, as they differ only in their receiver behaviors. Assuming that receiver behavior contributes differently than sender behavior would therefore modulate differences between these types.

---

**Reviewer 3 comment 2**

(*message amount*)

what motivates modelling just three messages and not just two (“some” and “all”) or any  $n > 3$ ? One obvious prediction:  $n = 2$  will lead to larger prevalence of lack types, whereas increasing  $n$  beyond 3 will do the opposite. More concretely: is there any non-trivial relative frequency of some/all that could falsify the results obtained in the paper?

As now made explicit in Section 3.2, we chose three messages for illustrative purposes mainly. In particular, we wanted to inspect a type space where targets are not the most likely *a priori* as this may mislead readers into thinking that learnability alone, in particular the prior alone, drives the outcome. A larger type space also better allows us to showcase that multiple kinds of a type may exist and how this affects evolutionary outcomes. The reason why we did not go beyond 3 messages is for computational tractability. Calculating  $Q$ , even if approximated, is expensive. The question whether

increasing the type space will decrease the frequency to which we expect to see pragmatic *some* users is interesting. At present, we have no answer.

The new section 3.1 now also looks at the case with just two messages that the reviewer brings into play. It is useful to contrast exactly why the effects of iterated learning hinge on the set of all types (the point made in connection with appendix C).

---

**Reviewer 3 comment 3**

(*LOT & complexity*)

The discussion on derivation costs is not very satisfying on my opinion. For instance, the condition of "simple representations being favored over more complex ones" is compatible with many more cost specifications.

We have expanded this section to further motivate our choices (see also the discussion in Section 4). Relatedly, see our answer to Reviewer 1's comment 4 above.

## References

- Chater, N. and Vitányi, P. (2003). Simplicity: a unifying principle in cognitive science? *Trends in Cognitive Sciences*, 7(1):19–22.
- Feldman, J. (2000). Minimization of Boolean complexity in human concept learning. *Nature*, 407(6804):630–633.
- Franke, M. (2009). *Signal to Act: Game Theoretic Pragmatics*. PhD thesis, University of Amsterdam.
- Franke, M. and Jäger, G. (2016). Probabilistic pragmatics, or why Bayes' rule is probably important for pragmatics. *Zeitschrift für Sprachwissenschaft*, 35(1):3–44.
- Gazdar, G. (1979). *Pragmatics, Implicature, Presupposition and Logical Form*. Academic Press, New York.
- Goodman, N. D. and Frank, M. C. (2016). Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11):818–829.
- Goodman, N. D. and Stuhlmüller, A. (2013). Knowledge and implicature: Modeling language understanding as social cognition. *Topics in Cognitive Science*, 5:173–184.

- Horn, L. R. (1972). *On the Semantic Properties of Logical Operators in English*. Indiana University Linguistics Club, Bloomington, IN.
- Kirby, S., Tamariz, M., Cornish, H., and Smith, K. (2015). Compression and communication in the cultural evolution of linguistic structure. *Cognition*, 141:87–102.
- Piantadosi, S. T., Tenenbaum, J. B., and Goodman, N. D. (2012). Bootstrapping in a language of thought: A formal model of numerical concept learning. *Cognition*, 123(2):199–217.
- Piantadosi, S. T., Tenenbaum, J. B., and Goodman, N. D. (under review). Modeling the acquisition of quantifier semantics: a case study in function word learnability.
- Qing, C. and Franke, M. (2015). Variations on a Bayesian theme: Comparing Bayesian models of referential reasoning. In *Bayesian Natural Language Semantics and Pragmatics*, pages 201–220. Springer International Publishing.