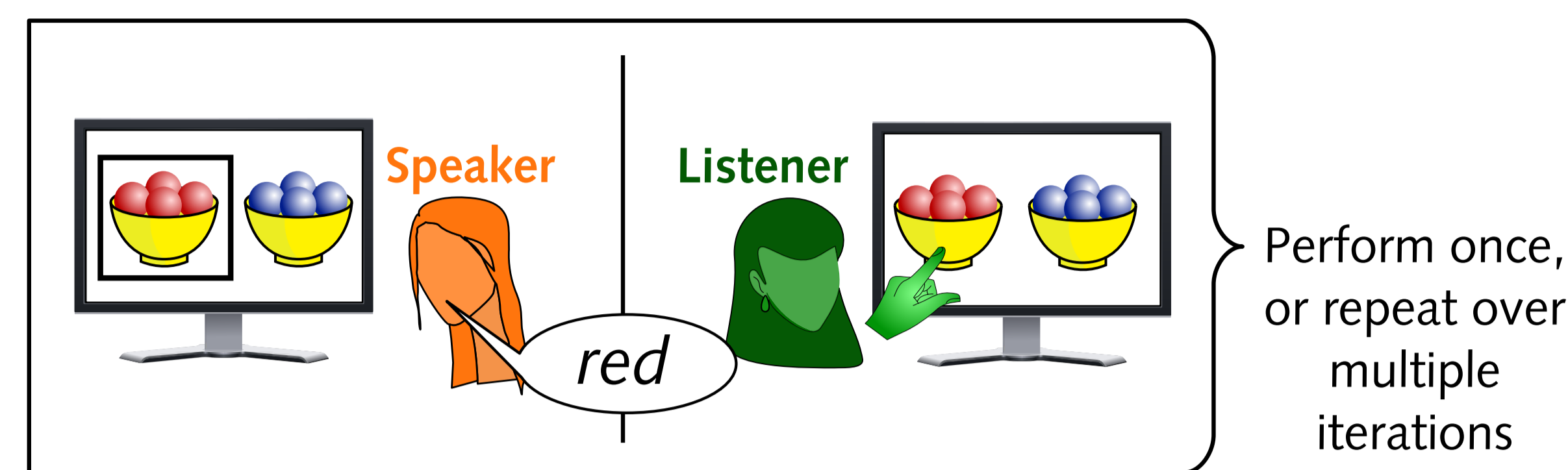# Learning and using language via recursive pragmatic reasoning about other agents
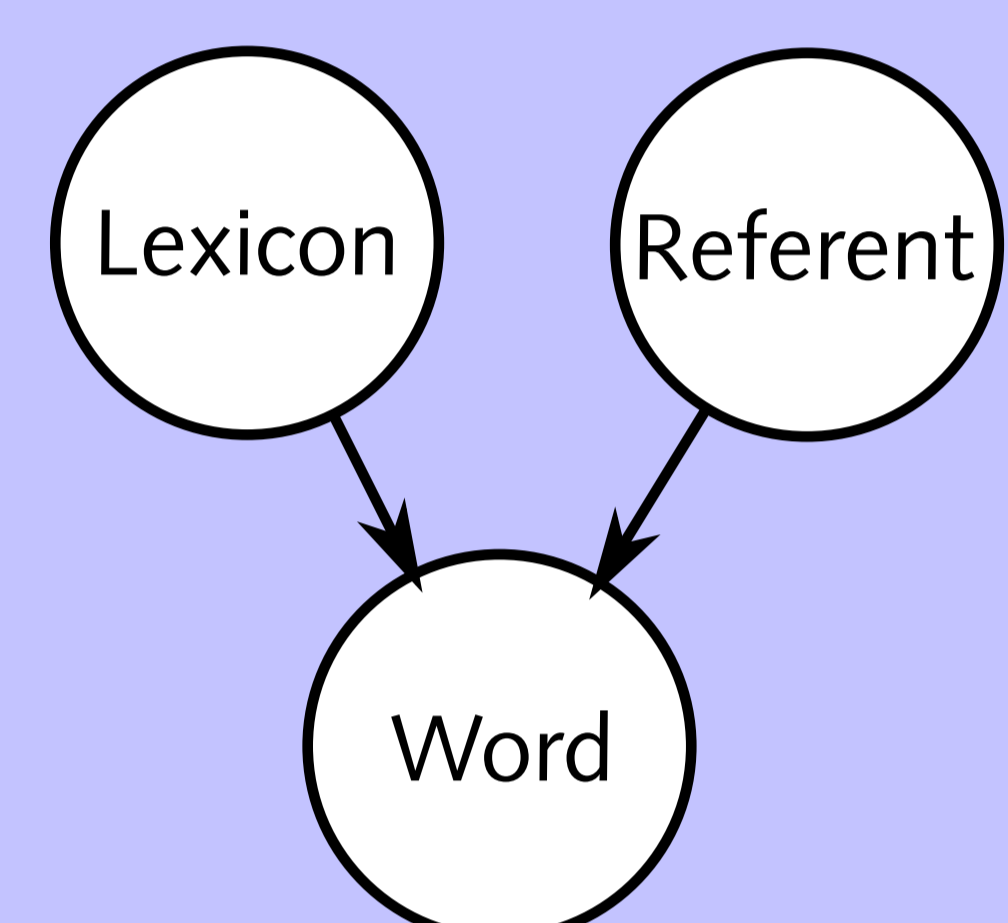
**Nathaniel J. Smith, U. Edinburgh**
<nathaniel.smith@ed.ac.uk>
**Noah D. Goodman, Stanford**
<ngoodman@stanford.edu>
**Michael C. Frank, Stanford**
<mcfrank@stanford.edu>

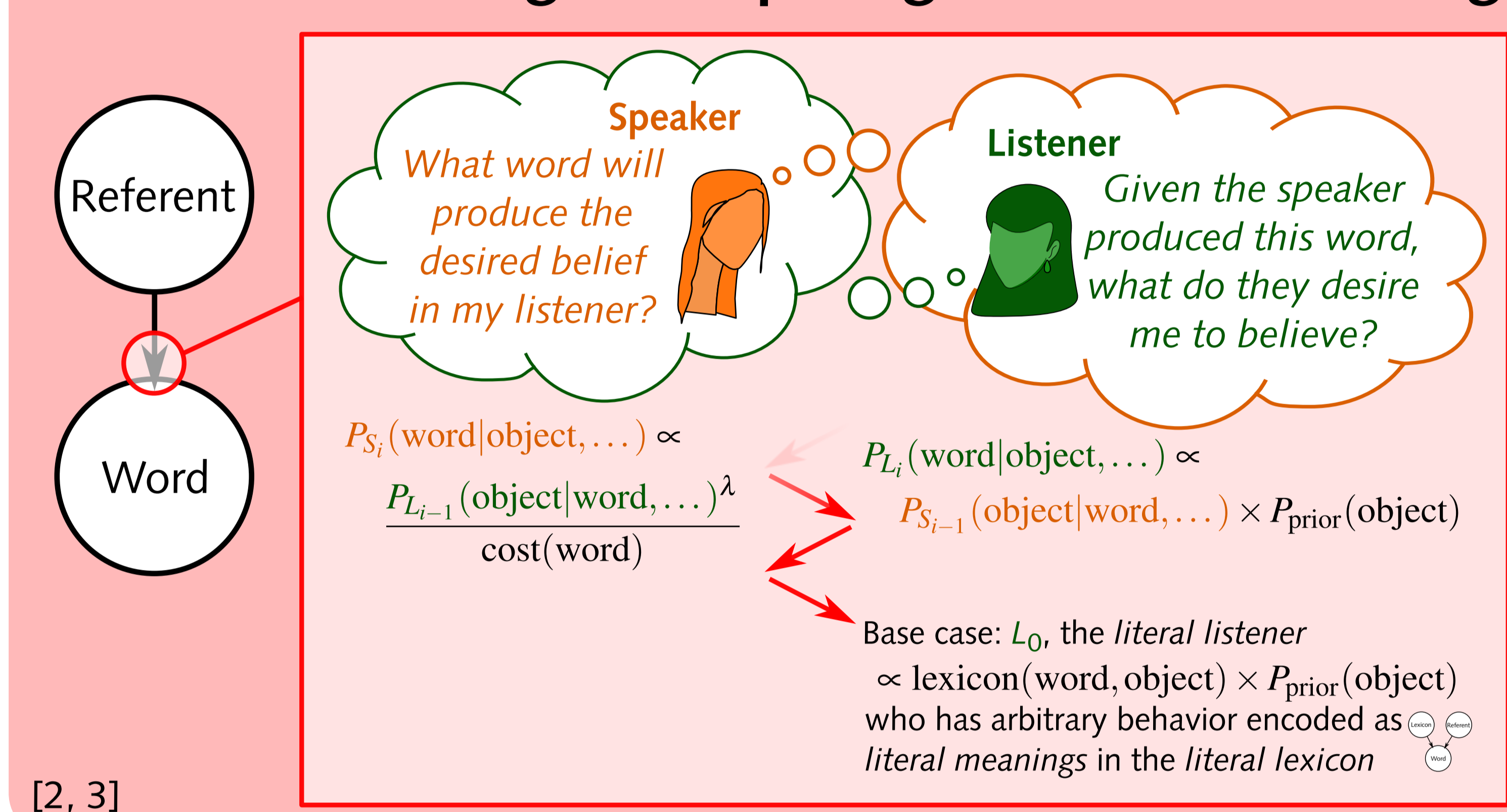## Our domain: We simulate language usage in simple referential games

Speaker    Listener
red

Perform once, or repeat over multiple iterations

## Our problem: What is a word's meaning?

**In existing learning models: a latent variable**

Lexicon    Referent
Word
[1]

**In existing pragmatics models: a communicative goal, requiring recursive reasoning**

Referent
Word

**Speaker**
*What word will produce the desired belief in my listener?*

**Listener**
*Given the speaker produced this word, what do they desire me to believe?*

$P_{S_i}(\text{word}|\text{object},\dots) \propto \dfrac{P_{L_{i-1}}(\text{object}|\text{word},\dots)^{\lambda}}{\text{cost}(\text{word})}$

$P_{L_i}(\text{word}|\text{object},\dots) \propto P_{S_{i-1}}(\text{object}|\text{word},\dots) \times P_{\text{prior}}(\text{object})$

Base case: $L_0$, the *literal listener* $\propto \text{lexicon}(\text{word},\text{object}) \times P_{\text{prior}}(\text{object})$ who has arbitrary behavior encoded as *literal meanings* in the *literal lexicon*.

[2, 3]

### A fundamental paradox for social learning

Word meaning, in the latent variable sense, does not appear in either Bayesian pragmatics models or the real world's actual generative process. Without uncertainty about a latent variable, learning is impossible. Yet humans both learn word meanings and perform pragmatic reasoning. **How can we reconcile these approaches?**

### Several obvious solutions don't work

Make $L_0$ or $S_1$ uncertain about the literal lexicon?
　But then actual speaker and listener must both marginalize out this uncertainty, so no-one's behavior is sensitive to the actual literal lexicon, so there is no data to learn it.
Make each recursive agent maintain uncertainty about their model of the next agent down?
　Arguably correct in theory, but would require actual speaker and listener to learn hyper-hyper-...-hyper-distributions, which is probably impossible even in principle due to data sparsity.

## Our solution: Assume conventionality + knowledgeable peers

Each agent assumes:
　(a) There is a specific, "conventional" literal lexicon that everyone is supposed to be using,
　(b) and everyone else knows what this lexicon is, and believes that I know it as well,
　(c) but in fact I don't know it, and have to do my best to fake it.
Assumption (a) explains why naive language users will argue – falsely! – that words have objective meanings (the "lexicographer's illusion"). Assumption (b) means that data is available, but avoids the explosion of hyper$^n$-distributions (and is uncomfortably familiar). Assumption (c) means there's something to learn. When combined with standard Bayesian techniques, these assumptions give a **definition of what a "convention" is**, and provide a **mechanism for learning, using, and creating them**.

### Our model

**Convention-based speaker $S$:**
Acts like $S_n$ + uncertainty about $L_0$:
$P_S(\text{word}|\text{object}, S\text{'s data}) \propto$
$\dfrac{\left(\sum_{\text{lexicon}} P_{L_{n-1}}(\text{object}|\text{word},\text{lexicon})P(\text{lexicon}|S\text{'s data})\right)^{\lambda}}{\text{cost}(\text{word})}$

Learning: uses $L_{n-1}$ as generative model

**Convention-based listener $L$:**
Acts like $L_{n-1}$ + uncertainty about $L_0$:
$P_L(\text{object}|\text{word}, L\text{'s data})$
$= \sum_{\text{lexicon}} P_{L_{n-1}}(\text{object}|\text{word},\text{lexicon})P(\text{lexicon}|L\text{'s data})$

Learning: uses $S_{n-2}$ as generative model

We set $n = 3$, $\lambda = 3$, and perform inference by importance sampling of lexicons from a Dirichlet prior.

**References:**
[1] M. C. Frank, N. D. Goodman, and J. B. Tenenbaum. Using speakers' referential intentions to model early cross-situational word learning. *Psychological Science*, 20:578–585, 2009.
[2] M. C. Frank and N. D. Goodman. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012.
[3] L. Bergen, N. D. Goodman, and R. Levy. That's what she (could have) said: How alternative utterances affect language use. In *Proceedings of the 34th Annual Conference of the Cognitive Science Society*, 2012.

## Phenomena

| | | [1]'s model | [2]'s model | [3]'s model | Our model |
|---|---|---|---|---|---|
| **Learning** | **Disambiguating new words using old words** — Click on the **dax**. **Learning new words using old words** Both humans and our model can learn words given only ambiguous input of this kind. | ✔ | | | ✔ |
| | **Disambiguation without learning** In some situations children select the correct object but then do not retain this mapping. Our model suggests an intriguing explanation: on single exposures, it selects the novel object but hypothesizes *dax* is a vague/generic term being strengthened by a specificity implicature. | | | ✔ | ✔ |
| **Pragmatics** | **Specificity implicature** *Using a vague word implies that a specific word does not apply.* Click on the face with **glasses**. Click on the bowl with **some** red balls. | ✔ | ✔ | ✔ | |
| | **Horn implicature** *Cheap utterances refer to common objects/situations; expensive utterances refer to unusual objects/situations.* Pretend you live in a world with these objects → Now do this task: Click the { wub. / flustergubbet. }  Black Bart (killed the sheriff / caused the sheriff to die). I (started the car / got the car started). | | ✔ | ✔ | |
| **Learning + Pragmatics** | **Learning scalar quantifiers like "some" and "all"** Training data: pragmatically strengthened uses "some" / "all" ×5 → Marginal literal belief "some" 0.78 / 0.22  *Dirichlet-multinomial: 0.83 / 0.17* | | | | ✔ |
| | **Emergence of novel & efficient conventions in interaction** When humans interact, novel and task-adapted communicative systems emerge. We simulate interactions between agents who begin with uniform priors over lexicons. | | | | ✔ |

**Emergence of novel & efficient conventions in interaction**

Mean communicative success over time
— 2x2 uniform prior
— Horn implicature
Dialogue turn

Example runs
Simple naming game
Horn implicature game
P(L understands)
Better run
Worse run
Dialogue turn

**Performance improves** as agents **converge on a shared lexicon**

Bias towards **sparse** lexicons

Bias towards **"Horn-compatible"** lexicons; implicature **shifts to become literal meaning**

→ **Biases** from "irrational" social anxiety assumption **systematically drive the communicative system towards greater efficiency**, and in the long run may leave their mark on structure of languages themselves.