

# Supplementary materials

## Pragmatic interactions lead to efficient language structure and use

### 1 Rational Speech Act theory speaker and listener agents

RSA is a recursive Bayesian model of pragmatic language use, which can largely be seen as a mathematical formalization of essential Gricean principles. RSA has proven to be a productive framework for modeling a range of pragmatic phenomena in both language production and language use including hyperbole, metaphor, implicature and others (CITATIONS).

In the RSA framework, a "speaker agent" defines a conditional distribution, mapping meanings  $u \in U$  to utterances  $m \in M$ , written as  $S(u|m)$ . We consider a prior over utterances  $P(U)$  as well as a prior over meanings  $P(M)$ . A "listener agent" defines a conditional distribution mapping from utterances to meanings, written as  $L(m|u)$ . To capture recursive reasoning between interlocutors, each of these functions is described in terms of the other. That is,

$$S_i(u|m) \propto e^{-\alpha \times U(u;m)} \quad (1)$$

where

$$U(u; m) = -\log(L_{i-1}(m|u)) - \text{cost}(u) \quad (2)$$

and

$$L_{i-1}(m|u) \propto S_{i-1}(u|m) \times p(m) \quad (3)$$

Defining nested speaker and listener agents could, in principle, lead to infinite regress. RSA defines a \*literal listener\*, denoted  $L_0(m|u)$ , as a base-case. The literal listener does not reason about a speaker model, rather this agent considers the literal semantics of the utterance.

$$L_0(m|u) \propto \delta_u(m) \times p(m) \quad (4)$$

with

$$\delta_u(m) = \begin{cases} 1, & \text{if } m \in [[u]] \\ 0, & \text{o.w.} \end{cases} \quad (5)$$

### 2 Zipfian objective for linguistic system efficiency

#### 2.1 Basic objective derivation

Zipf (1949) proposed that the particular distributional properties found in natural language emerge from competing speaker and listener pressures. We operationalize this objective in equation (1) – the efficiency of a linguistic system  $\ell$  being used by speaker and listener agents  $S$

and  $L$  is the sum of the expected speaker and listener effort to communicate over all possible communicative events  $e \in E$ .

$$\begin{aligned} \text{Efficiency}(S, L, \ell) = & \mathbb{E}_{e \sim P(E)}[\text{speaker effort}] \\ & + \mathbb{E}_{e \sim P(E)}[\text{listener effort}] \end{aligned} \quad (1)$$

We assume that speaker effort is related to the surprisal of an utterance in a particular context<sup>1</sup> – intuitively, the number of bits needed to encode the utterance  $u$ . This particular formalization of speaker-cost is general enough to accommodate a range of cost instantiations, such as production difficult via articulation effort, cognitive effort related to lexical access, or others [BennettGoodman2015a].

$$\text{speaker effort} = -\log_2(p(u|c))$$

We assume listener effort is the semantic surprisal of a meaning given an utterance. This operationalization of listener effort is intuitively related to existing work in sentence processing in which word comprehension difficulty is proportional to surprisal [Hale2001a; Levy2008a].

$$\text{listener effort} = -\log_2(L(m|u, c; \ell))$$

Importantly, we assume that events  $e = \langle u, m, c \rangle$  are sampled according the to following generative model – some context occurs in the world with probability  $P(C = c)$ . Within this context, an object  $m$  occurs with probability  $p(m|c)$ . The speaker attempts to refer to that object by sampling from her conditional distribution  $S(u|m, c; \ell)$  (i.e.  $e \sim p(c)p(m|c)S(u|m, c; \ell)$ ) Rewriting (1) we have

$$\text{Efficiency}(S, L, C, \ell) = \mathbb{E}_{e \sim P(E)}[-\log_2(p(u|c))] + \mathbb{E}_{e \sim P(E)}[-\log_2(L(m|u, c; \ell))] \quad (2)$$

In general, these expectations are each taken over the set of possible communicative events  $e \in E$  weighted by their probability,  $P(E = e)$ . Recall this is the set of all utterance, meaning pairs  $\langle u, m \rangle = e \in E$ .

$$= \sum_{e \in E} p(e)[-\log_2(p(u|c))] + \sum_{e \in E} p(e)[-\log_2(L(m|u, c; \ell))] \quad (3)$$

We assume that the particular joint distribution over utterance-meaning-context  $\langle u, m, c \rangle = e$  triples follows from a simple generative model (ancestral sampling). First, some context is sampled with probability  $p(c)$ . Then some meaning is sampled with probability  $p(m|c)$ . Our speaker attempts to convey this intended meaning to a listener via an utterance  $u$  by sampling from the speaker conditional distribution  $S(u|m, c; \ell)$ . Combining these terms leads to the \*speaker’s joint distribution over utterance-meaning pairs which we can write as  $P(e) = S(u|m, c; \ell)p(m|c)p(c) = P_{\text{speaker}}(u, m, c; \ell)$ .

---

<sup>1</sup>In the current set of simulations we consider utterances costs as independent from context (ie.  $p(u|c)p(c) = p(u)p(c)$ ).

$$\begin{aligned}
&= - \sum_{u,m} P_{speaker}(u, m, c; \ell) [\log_2(p(u|c))] - \\
&\quad \sum_{u,m} P_{speaker}(u, m, c; \ell) [\log_2(L(m|u, c; \ell))]
\end{aligned} \tag{4}$$

$$= - \sum_{c \in C} p(c) \left( \sum_{u,m} P_{speaker}(u, m|c; \ell) [\log_2(p(u|c))] + \sum_{u,m} P_{speaker}(u, m|c; \ell) [\log_2(L(m|u, c; \ell))] \right) \tag{5}$$

$$= - \sum_{c \in C} p(c) \sum_{u,m} P_{speaker}(u|m, c; \ell) [\log_2(L(m|u, c; \ell))] \tag{6}$$

Simplifying we arrive at (8):

$$= - \sum_{c \in C} p(c) \sum_{u,m} P_{speaker}(u, m|c; \ell) [-\log_2(L(m|u, c; \ell)p(u|c))] \tag{7}$$

The inner summation is the cross-entropy between speaker and listener conditional distributions over utterance-meaning pairs.

$$= - \sum_{c \in C} p(c) H_{Cross}(P_{speaker}(u, m|c; \ell), P_{listener}(u, m|c; \ell)) \tag{8}$$

This is an expectation of speaker-listener cross-entropy over contextualized language use.

$$= \mathbb{E}_{c \sim P(C)} [H_{Cross}(P_{speaker}(u, m|c; \ell), P_{listener}(u, m|c; \ell))] \tag{9}$$

Note that in the case that  $|C| = 1$ , our objective simplifies to simply the cross-entropy between speaker-listener joint distributions over utterance-meaning pairs.

From an information-theoretic perspective this objective is intuitive. Cross-entropy gives us a measure of dissimilarity between two distributions – the average number of bits required to communicate under one distribution, given that the “true” distribution differs. In our case, this is the difference between the joint distribution assumed by the speaker  $P_{speaker}$  and listener  $P_{listener}$ . A good language  $\ell$  used by a set of a pair of speaker-listeners will have properties which minimize this objective.

## 2.2 Baseline model objectives

For comparison, we also examine properties of optimal languages under two additional objectives. Zipf [-@Zipf1949a] proposed that the optimal speaker language  $\ell_{speaker}^*$  should only optimize speaker effort. We operationalize this using the \*first half\* of equation (1).

$$\ell_{speaker}^* = \operatorname{argmin}_{\ell \in L} \mathbb{E}_{P_{speaker}(u, m; \ell)} (-\log_2(p(u))) \tag{1}$$

The optimal listener language  $\ell_{listener}^*$ , by contrast, should only optimize listener effort. We operationalize this using the \*second half\* of equation (1).

$$\ell_{listener}^* = \operatorname{argmin}_{\ell \in L} \mathbb{E}_{P_{speaker}(u, m; \ell)} (-\log_2(L(m|u; \ell))) \tag{2}$$

## 3 Simulation 2

### 3.1 Updates to RSA speaker-listeners

We consider the same model of basic speaker and listener ( $S_{vanilla}$ ,  $L_{vanilla}$ ) models in Section 1 in conjunction with discourse aware speaker-listeners ( $S_{discourse}$ ,  $L_{discourse}$ ) who can use the history of utterances (the discourse  $D$ ) to infer the topic of conversation ( $c \in C$ ):

$$S_{discourse}(u|m, c, D) \propto e^{\alpha U(u, c; m, D)}$$

$$U(u, c; m, D) = -\log_2(L_{discourse}(m, c|u, D)p(c|D)) - cost(u)$$

where

$$p(c|D) \propto p(c) \prod_{i=0}^{|D|} S_{vanilla}(u_i|m_i)p(m_i|c)$$

and

$$L_{discourse}(m, c|u, D) \propto S_{vanilla}(u|m)p(m|c)p(c|D)$$

Note that  $p(M|C = c)$  is simply the particular prior over meanings dictated by a topic  $c$ .

### 3.2 Language used in Simulation 2

We conduct  $N = 600$  simulations, generating discourses of length  $|D| = 30$  utterances with three different speaker models ( $n = 200$  each). We consider a single language  $\ell$  with  $|U| = 6$  and  $|M| = 4$  specified by the boolean matrix below. (Note that use of this particular language is not essential – the results are broadly generalizable languages that contain ambiguity.)

	$m_1$	$m_2$	$m_3$	$m_4$
$u_1$	1	0	0	0
$u_2$	0	1	0	0
$u_3$	0	0	1	0
$u_4$	0	0	0	1
$u_5$	1	1	0	0
$u_6$	0	0	1	1

We assume that  $p(u_5) = p(u_6) > p(u_1) = \dots = p(u_4)$ . That is, the two ambiguous utterances ( $u_5$  and  $u_6$ ) are less costly than the non-ambiguous utterances.

## References