

## Les Sikos, Noortje Venhuizen, Heiner Drenhaus, and Matthew W. Crocker (Saarland University)

#### Background

Tests of formalizations of Gricean maxims using web-based reference games have led to mixed results:

#### Frank & Goodman [2]

- One-shot paradigm
- 3 object displays in 7 visual context types
- Collected separate judgments from Speakers, Listeners, and for Salience

Speaker Task. Imagine you are talking to someone and you want him to refer to the middle object. Which word would you say, "green" or "circle"?

Listener Task. Imaging someone is talking to you and uses the word "green" to refer to one of these objects. Which object are they talking about?

Salience Task. Imaging someone is talking to you and uses a word you don't know to refer to one of the objects. Which object are they talking about?

• Rational Speech Act (RSA) model (Eq 1) closely predicted aggregate listener judgments



 This result was interpreted as indicating that participants reasoned pragmatically in this task

#### **Alternative Explanation**

 The reasoning required ranged from simple to more complex







Listener Listener

• The close fit of predicted to observed results might be driven by the simpler inferences

#### **Consistent With This Possibility**

- [3] attempted a close replication of [2], focusing on more challenging items and found that the basic RSA model was a poor predictor of their data
- [4] found that while listeners responded pragmatically in simpler contexts, they were at chance for more complex contexts
- To account for these results, [3] and [4] proposed various modifications to the RSA model (e.g., adding parameters for speaker/ listener degree of rationality)
- → Participants rarely go beyond the literal meanings of words in such studies

#### Goals

- Do Listeners in su presumed?
- Or can a simpler ra better explain hum

#### **Participants**

- 4642 participants i Amazon Mechanic
- 1137 excluded for non-native English
- 118 excluded for ir attention question

#### Procedure



**Eq 1.** RSA model for inferring the speaker's intended referent  $r_s$  in context C, given speaker's uttered word *w*. The Listener model combines a Speaker model (i.e. the likelihood that speakers use a particular word to refer to the target) with empirically measured salience.

# **Reevaluating Pragmatic Reasoning in Web-based Language Games**

#### Goal and Methods

						5				
ch task	s reason a	as pragmaticall	y as							
ather t man be	han more havior th	e complex mod an RSA?	el		ć	objects that	color ha	iape etst	iave jolor	
• Remaining 3 recruited via randomly as			387 were		Ann per of	aree confection	ent Ischetit	o`, l'so <sup>sent</sup> E (targ	Exampl set in mi	e iddle
identifying as		Speaker N Listener N	Speaker $N = 1143$ Listener $N = 1111$		1s1c	ds	dc			
ncorrec	ct	<ul><li>Salience N</li><li>Demographi</li></ul>	I = 1133	1s1c	ds	SC			9	
ıs (3%)		Gender	F 53% M 48%		1s1c	SS	dc			
Age 18-2 26-3			5 17% 5 40%		1s1c	SS	SC			
ıtil after cal Turk.		36-4 46-5 56-6	5       22%         5       13%         5       6%         5       2%	1s2c	ds	dc			ſ	
Next		over 6		1s2c	SS	dc			ſ	
Speaker T		<b>Fask</b> . Imagine		-	1s3c	ds	SC			
	Robert ar him to pic	nd you want ck out Item B. —	_ The target was always object B	_	1s3c	SS	SC			ſ
Next ntil after ical Turk.	If you can only use one word, which word would you say,		Order of options	2s1c	ds	dc			9	
	Listener	<b>Fask.</b> Robert	<ul> <li>were counterbalanced for shape/color first</li> </ul>		2s1c	ds	SC		1	2
Next ntil after nical Turk.	of the ob but he ca	jects below n only say one	Given words were	2s2c.a	ds	dc			9	
	Word. He says, "green". — Which object do you think he is talking		<ul> <li>counterbalanced</li> <li>for shape/color</li> </ul>		2s2c.b	ds	dc			9
Next	about: A	, B, or C? Task. Robert			2s3c	ds	SC			٢
ntil after Ical Turk.	of the ob	jects below, o background		_	3s1c	SS	dc			
Robert	understand what he said. Which object do			-	3s1c	SS	SC			•
Finish HIT	likely talk A, B, or C	ing about: ?		-	3s2c	SS	dc			ſ
					3s3c	SS	SC			
model Sa	alience	Listener model	Salience		$\Delta$ - mode respo	els make d nse, giver	ifferent p a color (	redictior c) or sha	າs for Lis pe (w) v	stene vord

 $P(w|r_{\rm s},C)P(r_{\rm s})$ 

Listener model  

$$P(r_s | w, C) = \begin{cases} \frac{P(r_s)}{\sum_{r' \in R} P(r')} & \text{if } r_s \in R \\ 0 & \text{otherwise} \end{cases}$$
where:  $R = \{r \in C \mid w \text{ can refer to } r\}$ 

Eq 2. Literal Listener + Salience model. This model does not contain a speaker model. Instead, it provides a distribution over the set of referents in context C that can be referred to with word w, weighted based on empirically measured salience.

**Analysis** – We compared observed Listener responses to predictions from 4 models, as well as to observed Salience responses

- Basic RSA
- Basic RSA with uniform salience (RSA.us)
- Literal Listener + Salience (LL+S)
- LL with uniform salience (LL.us)
- Salience observations



ons for Listener



### **Predictions vs Observed over All Visual Context Types**



Correlations

ank	Model	R	Adj R sq	t	p	Rank	Model	summed KLD	mean KLD
1	LL+S	0.964	0.929	36.429	< .0001 ***	1	LL+S	0.954	0.028
2	RSA	0.954	0.909	31.805	< .0001 ***	2	RSA	1.474	0.043
3	LL.us	0.944	0.889	28.488	< .0001 ***	3	LL.us	4.166	0.123
4	RSA.us	0.918	0.84	23.085	< .0001 ***	4	RSA.us	5.466	0.161
5	Salience	0.506	0.248	5.862	< .0001 ***	5	Salience	6.418	0.189





Contexts for which

Models make

**Same Predictions** 

Contexts for which

Models make

**Different Predictions** 

r = -0.23, p = 0.32 0.75 0 25

RSA with uniform salience

#### Correlations

Rank	Model	R	Adj R sq	
1	RSA.us	0.988	0.976	
1	LL.us	0.988	0.976	
2	RSA	0.970	0.940	
2	LL+S	0.970	0.940	
3	Salience	0.448	0.190	
Rank	Model	R	Adj R sq	
1	LL+S	0.908	0.815	(
2	Salience	0.893	0.786	8
3	RSA	0.790	0.603	ļ
4	LL.us	0.313	0.048	
F		0 222	0 0 0 0	

#### **Discussion and Conclusions**

#### **Results were consistent with the alternative hypothesis**

- Although RSA had good fit to entire dataset (replicating [2]), LL+S performed better • When considering only contexts for which predictions from RSA and LL models
- differed (i.e. the more challenging inferences):
- (a) LL+S performed best (b) Salience alone was a better predictor than RSA • Comparing RSA and RSA.us models suggests salience essentially corrects for incorrect predictions in the basic RSA model
- Modified RSA models (*ala* [3, 4]) performed worse than the basic RSA model

**References** [1] Grice (1975). *Logic and conversation*. [2] Frank & Goodman (2012). Predicting pragmatic reasoning in language games. Science. [3] Qing & Franke (2015). Variations on a Bayesian theme: Comparing Bayesian models of referential reasoning. In *Bayesian natural language semantics and pragmatics*. [4] Frank, Emilsson, Peloquin, Goodman & Potts (2016). Rational speech act models of pragmatic reasoning in reference games. Preprint: osf.io/f9y6b.







#### Results

Kullback–Leibler Divergence (base 2)

#### Predictions vs Observed for Visual Contexts in which Model Predictions Differed ( $\Delta$ )



#### Conclusions

To the extent that one-shot web-based experiments accurately elicit the depth of pragmatic reasoning seen in typical human interactions, these findings indicate that a simpler model than RSA can better explain human behavior

#### Next Steps

- Increase engagement and cooperativity
- Allow multiple words and investigate ordering effects