

A Model, Training, and Dataset Details

All models are trained end-to-end with the Gumbel-Softmax trick with the Adam (Kingma and Ba, 2014) optimizer with learning rate 0.0001.

A.1 ShapeWorld

Model. f_{θ}^T and f_{ϕ}^S are 4-layer convolutional neural networks, each consisting of a 64-filter 3x3 convolution, batch normalization, ReLU nonlinearity, and 2x2 max-pooling layer, as used in the few-shot learning literature (Snell et al., 2017). RNN encoders and decoders are single layer Gated Recurrent Units (GRUs) (Cho et al., 2014) with hidden size 512 and embedding size 300. We train with batch size $B = 32$.

Data. As aforementioned, there are 312 total concepts (see Appendix B), of which 80% are reserved for training and 20% are reserved for test. From the training concepts, we sample 10,000 base games to use as our training set, each with 50 positive and negative targets each. At training time, we perform augmentation by randomly selecting 10 positive and negative examples given to both teacher and student, meaning that the total set of games is combinatorially large (over $\binom{50}{10}$ combinations for reference games, and $\binom{50}{10}^2$ combinations for setref and concept) and overfitting is highly unlikely. We set up validation and test datasets with 2000 games each, divided among seen and unseen concepts, with no augmentation performed. We train over epochs (defined by a single pass through 10,000 augmented games) until average performance on the validation set is maximized.

A.2 Birds

Model. f_{θ}^T and f_{ϕ}^S is an ImageNet (Russakovsky et al., 2015)-pretrained ResNet-18 (He et al., 2016); similar results were observed for models trained from scratch. RNN encoders and decoders are single layer GRUs with hidden size 1024 and embedding size 300. We train with batch size $B = 16$ and preprocess images with ImageNet mean normalization.

Data. From the 100 training classes, we sample games dynamically by randomly selecting 5 positive targets from the class and 5 negative targets randomly. Like in ShapeWorld, this makes the number of possible training games combinatorially large. We set up validation and test datasets

with 400 games each divided among seen and unseen concepts. We define an epoch as a single pass through 1,000 augmented training games, and like before, select the model with the highest performance on the validation set.

B ShapeWorld Concepts

The 312 ShapeWorld concepts are either:

1. A single *primitive* shape (*triangle, square, circle, ellipse, rectangle*) or color (*red, blue, green, yellow, white, gray*), possibly negated (e.g. *not gray*);
2. A disjunction of two (possibly negated) primitives (e.g. *blue or yellow, circle or not red*);
3. A conjunction of two (possibly negated) primitives (e.g. *red and triangle, red and not triangle*).

We enumerate all (boolean-equivalent) possible formulas, then discard formulas which are tautologically true (e.g. *not yellow or not red*) or unsatisfiable (e.g. *circle and square*).

For each concept, sampling positive and negative shapes uniformly often results in games that do not specifically test the concept. For example, for the concept *gray and not circle*, there may not be any negative *gray circles*, so the agent could just infer the concept *gray*. To ensure that concepts are fully tested, for disjunctive concepts, we sample 1/3 targets that satisfy *only* the *left* side of the disjunction; 1/3 that satisfy *only* the *right* side; and 1/3 that satisfy both. For conjunctions, we sample 1/3 *distractors* that only fail to satisfy the left side of the disjunction; 1/3 that only fail to satisfy the right side; and 1/3 that fail to satisfy both sides.

Code used for generating the dataset is available at <https://anonymized>.

C Experiments with Traditional Reference Games

We presented an atypical formulation of reference games as consisting of multiple targets, with student decisions made independently:

$$p^S(Y^S | X^S, m) = \prod_i p^S(y_i^S | x_i^S, m), \quad (1)$$

where students are trained with the binary cross entropy loss, defined for a single game as

$$\mathcal{L}_{\text{BCE}}(S) = - \sum_i \log p^S(y_i^S | x_i^S, m). \quad (2)$$

This was done to keep training objectives and models as consistent as possible, and to keep the amount of training data consistent (i.e. there are exactly the same number of targets and distractors seen by each agent across training).

However, the typical reference game has a single target: instead of $Y^S \in \{0, 1\}^n$, we have a single target $t^S \in [1, n]$ denoting the index of the single positive example. Then the student probability that input i is the target is the softmax-normalized

$$p^S(i | X^S, m) = \frac{\exp(\text{RNN-ENCODE}(m) \cdot f_\phi^S(x_i^S))}{\sum_{i'} \exp(\text{RNN-ENCODE}(m) \cdot f_\phi^S(x_{i'}^S))} \quad (3)$$

and the training objective for a single game is

$$\mathcal{L}_{\text{XENT}}(S) = -\log p^S(t^S | x_i^S, m). \quad (4)$$

To ensure that our alternative formulation did not affect results, we ran 5 experiments with the standard reference game trained with cross entropy loss, with a single target and 10 distractors. Figure S1 summarizes the relevant statistics; besides slightly higher topographic ρ and AMI for the cross-entropy reference games for ShapeWorld, there are no qualitative differences compared to our reference game formulation and our conclusions are unchanged.

D Concept and Setref teachers evaluated on Reference games

While forcing teachers to speak about generalizations may necessarily increase systematicity of the resulting languages, here we examine whether speaking in generalizations also increases systematicity when producing referring expressions. We test this hypothesis by presenting reference games at test time to teachers trained in setref and concept games, without any training (i.e. zero shot evaluation). While setref and concept agents have seen examples of conjunctions (e.g. *red triangles*), they have never had a game with identical targets.

Figure S2 displays accuracy and systematicity measures in this setting, with reference game statistics provided for comparison. We note some qualitative differences: for example, there is no longer a significant difference in entropy of messages between setref and ref games for ShapeWorld; on the other hand, measures of topographic ρ are even higher for setref and concept games for ShapeWorld (as much as 0.5 edit ρ compared to 0.2 in

Figure 3), suggesting that concept/setref teachers are perhaps most systematic with simple conjunctions of shapes and color. Overall, training agents in setref and concept settings increases systematicity, even when producing singular referring expressions, since the models have already been biased to produce generic language.

E Additional plots of speaker messages

See Figure S3 for additional plots of teacher messages made for ShapeWorld and Birds games. Overall, the plots show a general reduction in language complexity from ref to setref to concept, although some quirks emerge: for example, some characters (e.g. *e* in concept) appear to be overloaded (across *green ellipse* and *red*), and concept uses similar language for *painted bunting* and *yellow warbler*. White gaps indicate end of sentence, so there are games where the speaker teacher utters nothing (e.g. *blue or not circle* concept; *not red* setref).

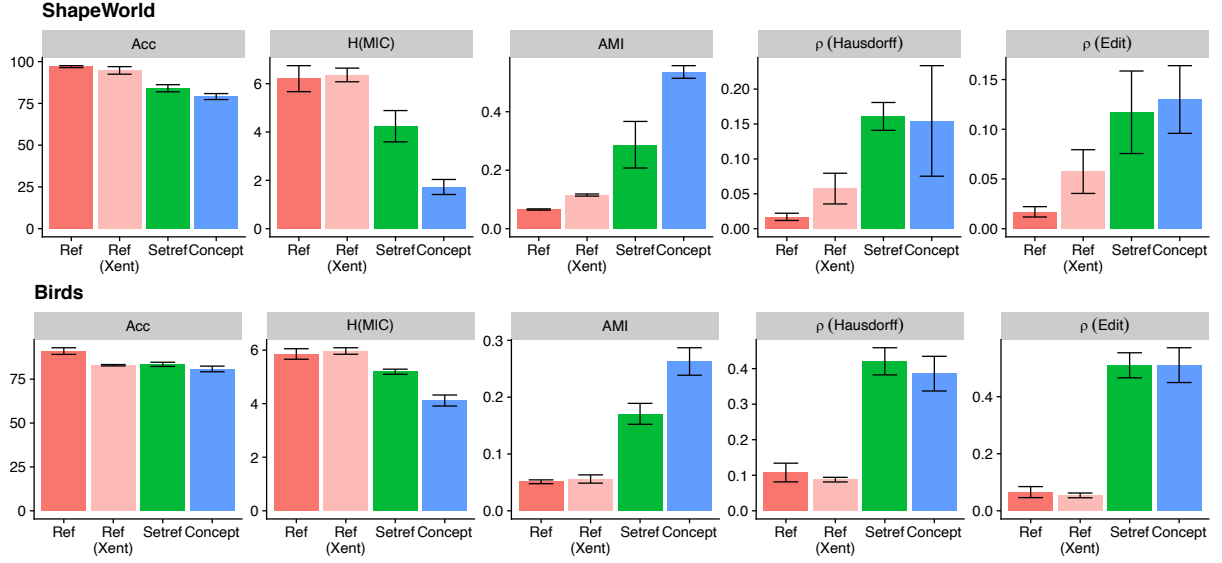


Figure S1: Accuracy and measures of language systematicity for reference games, setref games, and concept games, as well as reference games trained with the traditional cross entropy (xent) objective.

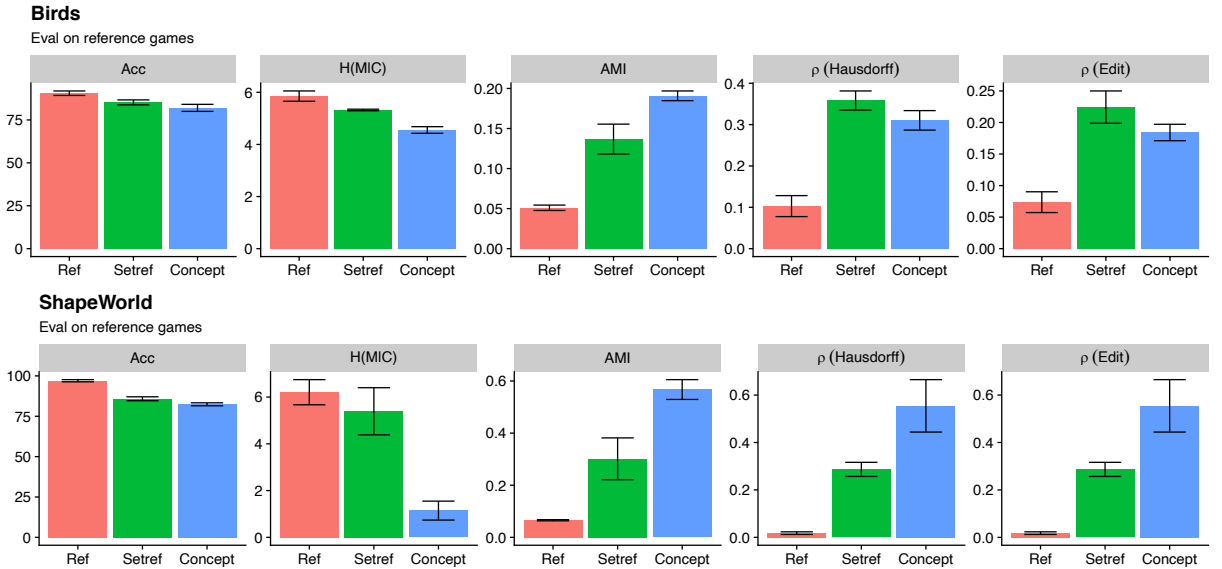


Figure S2: Accuracy and measures of language systematicity for setref games and concept games with zero-shot evaluation on reference games at test time. Reference game statistics provided for comparison.

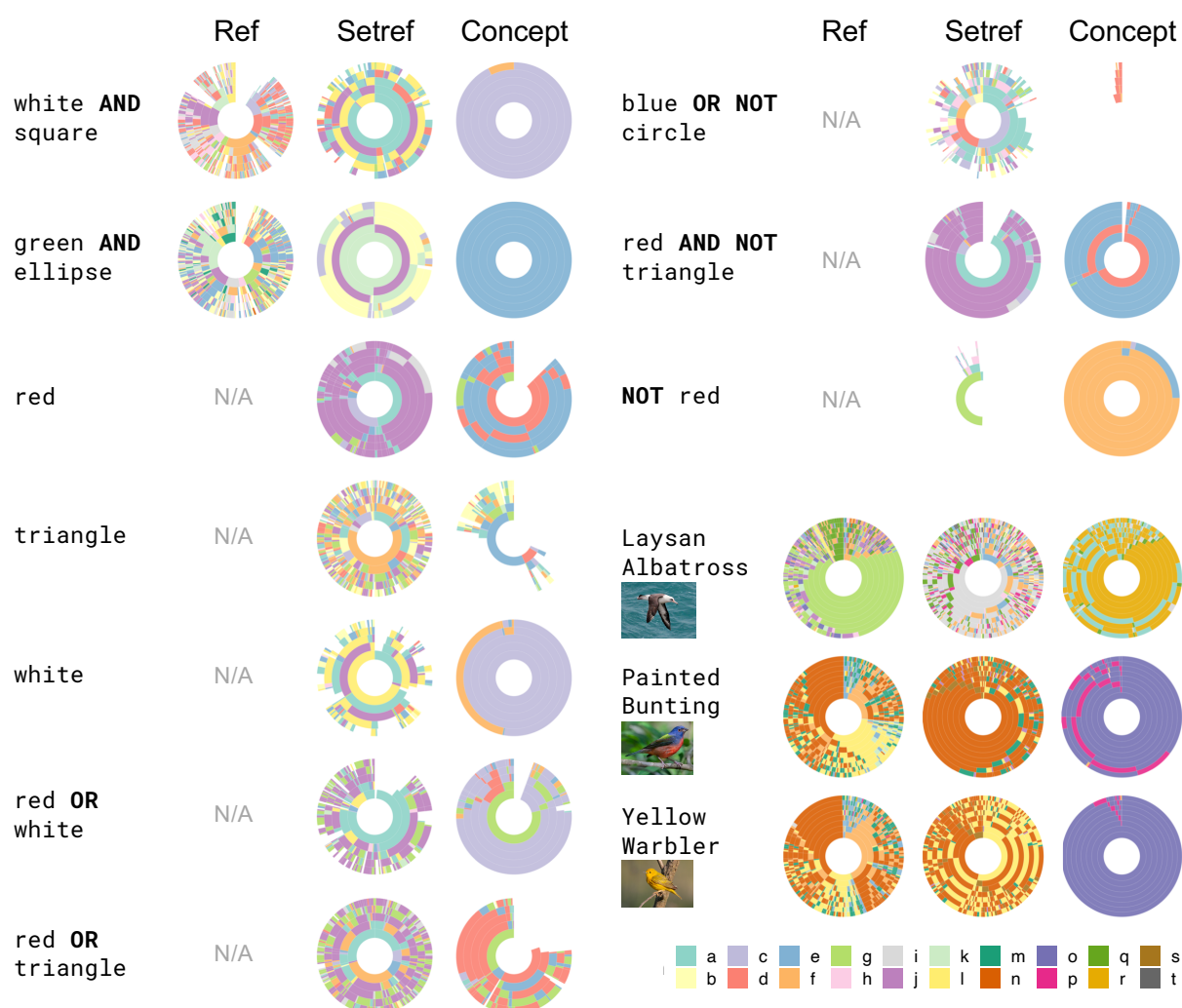


Figure S3: Additional plots of teacher messages for selected ShapeWorld and Birds games. Most ShapeWorld concepts are not tested in reference games, so those plots are not available.