

## Background & Motivation

- People can refer to the **shape** (geometry, topology) of an object to *distinguish* it among others objects.
- **Existing studies** rely primarily on properties like *color* and *spatial location* to refer to an object.
- **Existing studies** work explicitly with 2D images and are 'blind' to the *part-based* compositionality of 3D objects, or their *fine-grained* geometry.

## This Work

- Builds a large scale **multi-modal** dataset calibrated for shape-based reference, aka ShapeGlot!
- Introduces *novel speaker-listeners* considering the effect of using:
  - a) 2D & 3D object representations
  - b) context-based discrimination
  - c) neural word-attention
  - d) pragmatic referential reasoning
- Discovers a plethora of **surprising generalization** scenarios.

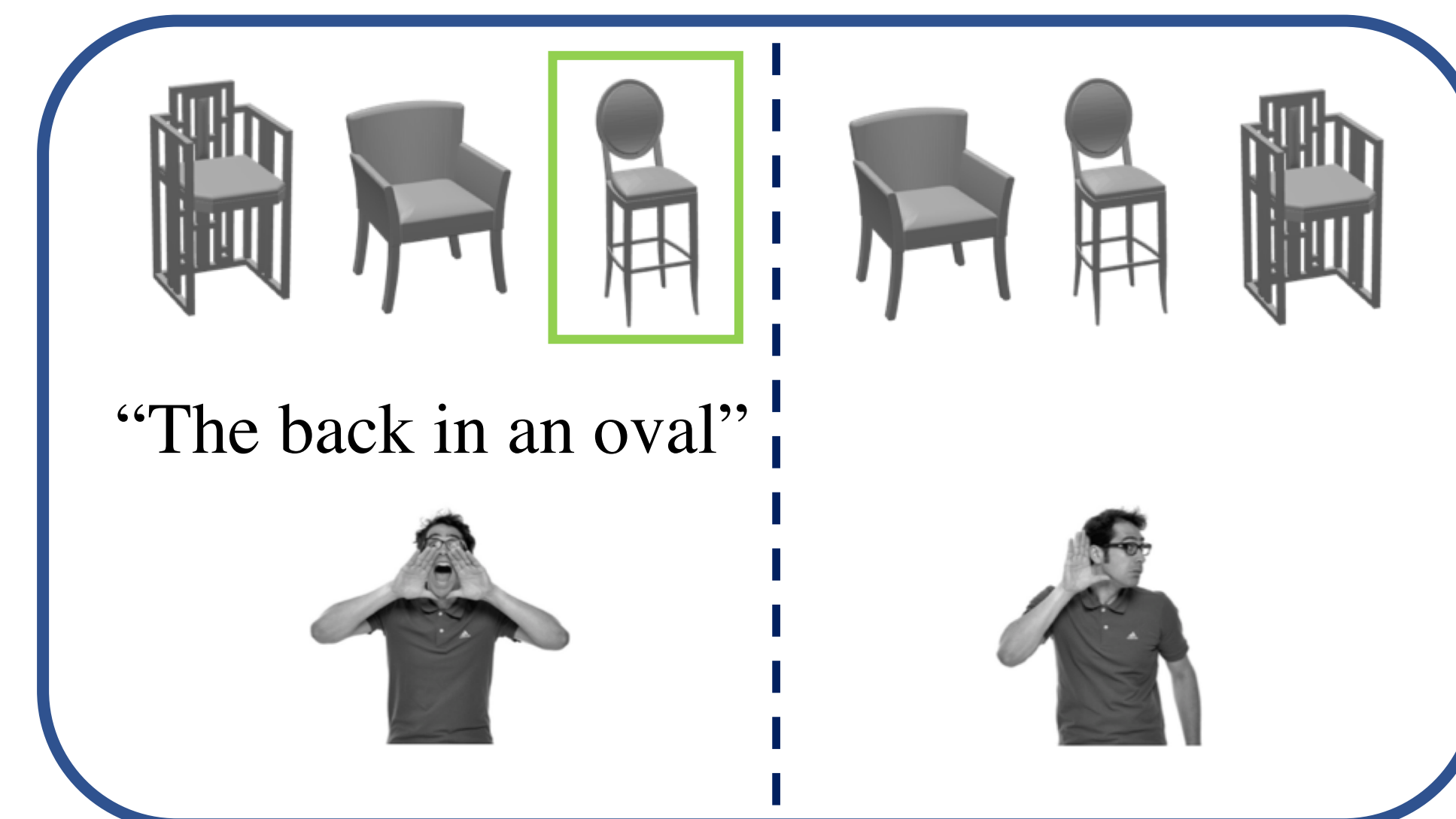
Code & Data



[www.bit.ly/shapeglot](http://www.bit.ly/shapeglot)

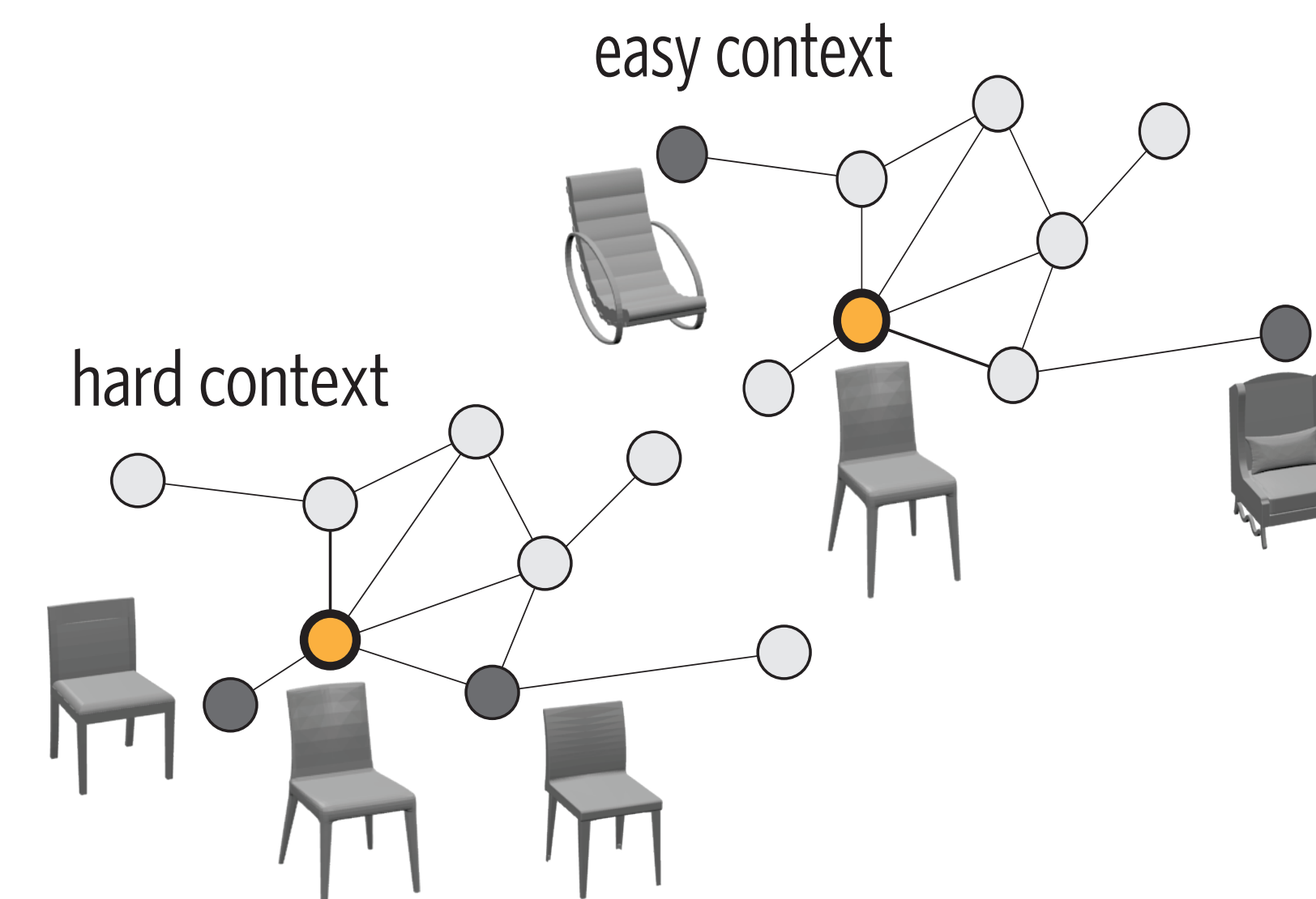
## Making ShapeGlot

- Tap on 'pure' 3D meshes to make a reference game



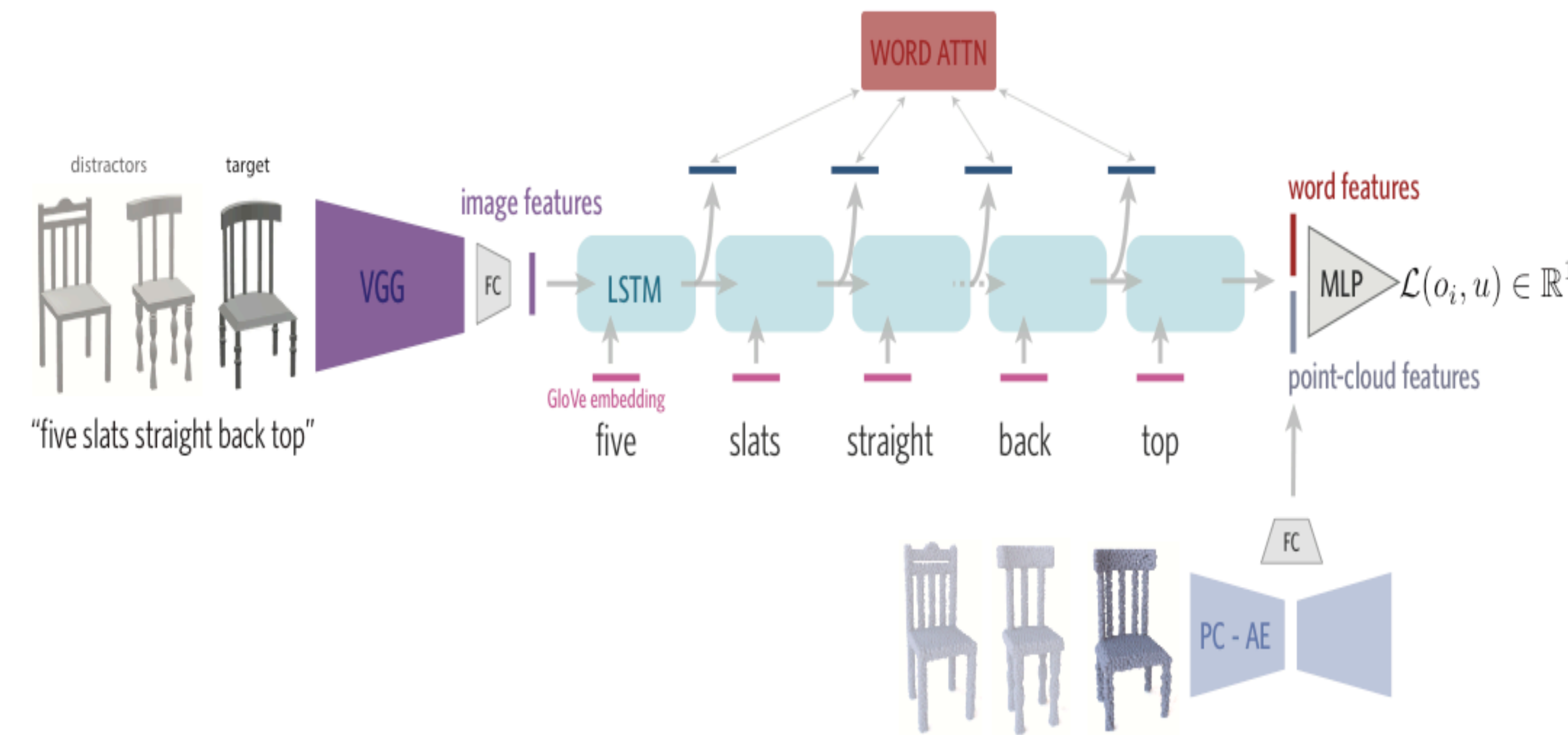
- ✓ Free-form language
- ✓ By construction: **only** about shape

- Use *latent*-based context formation



- ✓ 80K Utterances, 4K contexts

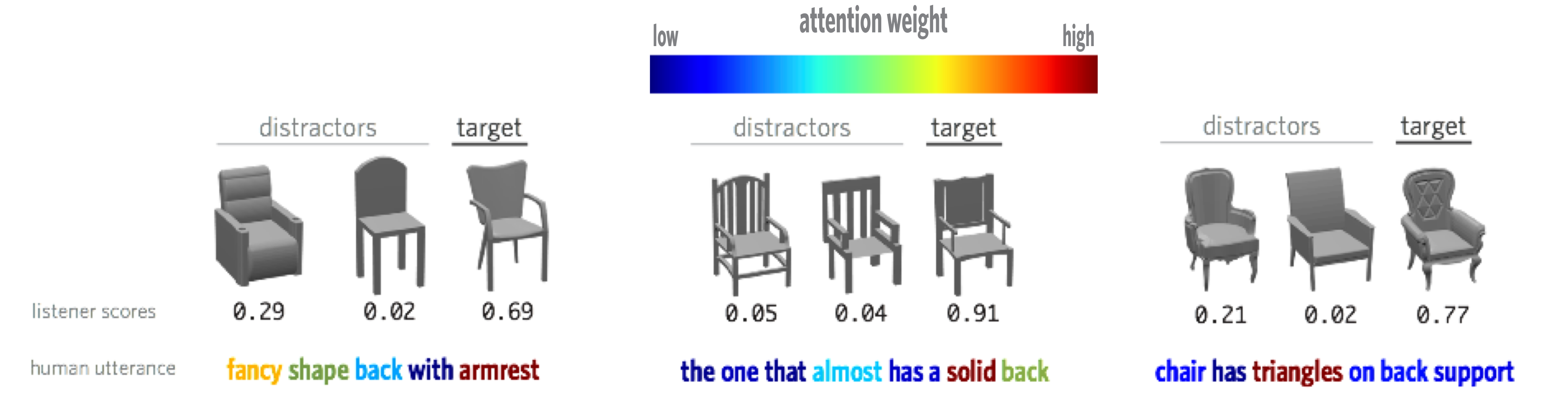
## Multi-modal Attentive Neural Listeners



## Pragmatic Neural Speakers

	distractors	target	distractors	target	distractors	target
listener scores	0.29	0.20	0.51	0.00	0.14	0.86
pragmatic speaker	it has rollers on the feet		square back, straight legs		thin-est seat	
listener scores	0.55	0.16	0.29	0.05	0.85	0.10
literal speaker	the one with the circle on the bottom		the one with the thick-est legs		the chair with the thin-est legs	
listener scores	0.19	0.24	0.57	0.19	0.32	0.49

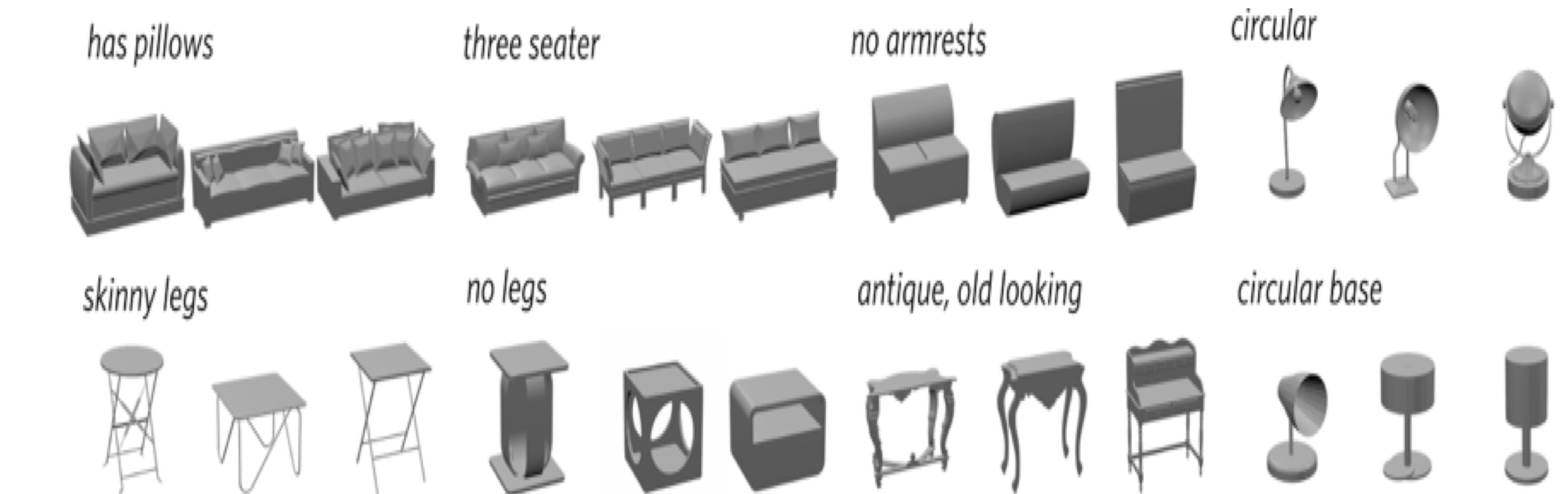
$$\text{Listener's "fit"} = \beta \log(P_L(t|U, O)) + \frac{(1-\beta)}{|U|^\alpha} \log(P_S(U|O, t)) \quad \text{Speaker's "fit"}$$



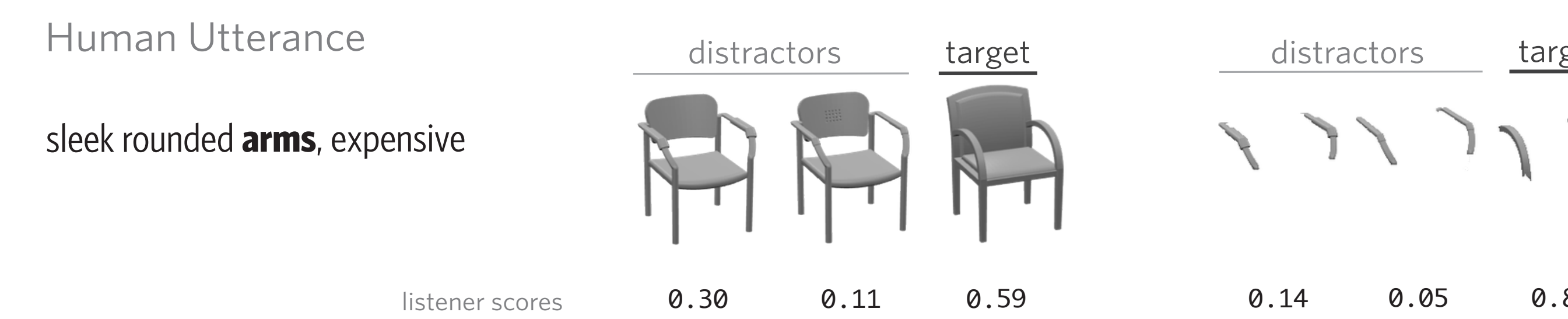
Neural attention provides intuitive model interpretation

	Input Modality	Language Task	Object Task
No Attention	Point Cloud	67.6 ± 0.3%	66.4 ± 0.7%
	Image	81.2 ± 0.5%	77.4 ± 0.7%
	Both	83.1 ± 0.4%	78.9 ± 1.0%
With Attention	Point Cloud	67.4 ± 0.3%	65.6 ± 1.4%
	Image	81.7 ± 0.5%	77.6 ± 0.8%
	Both	<b>83.7 ± 0.3%</b>	<b>79.6 ± 0.8%</b>

Listening Ablations



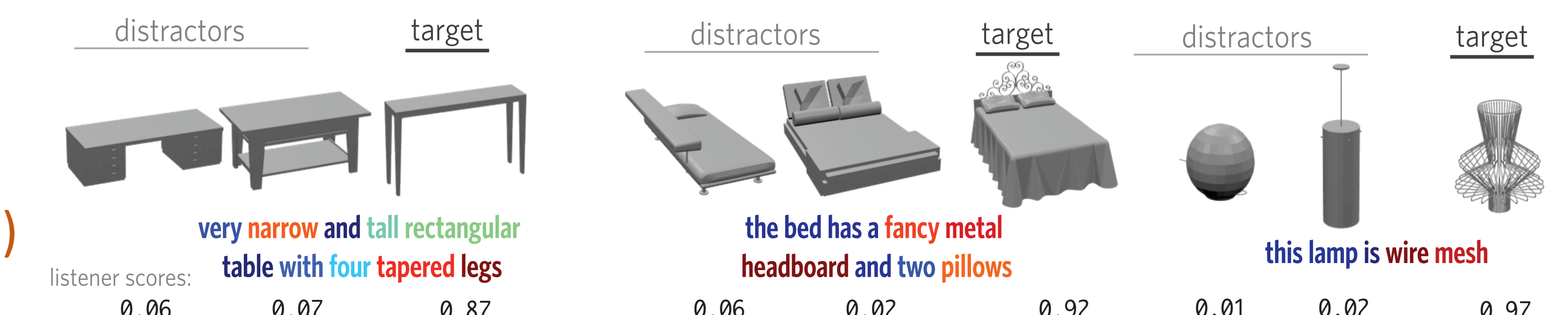
Listeners for novel shape-retrieval



	Single Part Lesioned	Single Part Present
Mentioned Part	42.8% ± 2.3	66.8% ± 1.4
Random Part	67.0% ± 2.9	38.8% ± 2.0

If used, object parts (e.g. visual 'arms') are necessary & sufficient for disambiguation

Zero-shot Listening (in unseen class & lang.)



Zero-shot Speaking (in catalogue models)



Speaker Architecture	Modality	Neural Listener	Human Listener
Context Unaware	Point Cloud	59.1 ± 2.0%	-
	Image	64.0 ± 1.7%	-
Literal	Point Cloud	71.5 ± 1.3%	66.2
	Image	76.6 ± 1.0%	68.3
Pragmatic	Point Cloud	90.3 ± 1.3%	69.4
	Image	<b>92.2 ± 0.5%</b>	<b>78.7</b>

Speaking Ablations

## Key Take Away Points

- Shape-based referential language is **robust** across classes (e.g. ZSL from 'chairs' to 'lamps').
- Language *alone* enables part-based **visual** reasoning.
- Pragmatic neural agents perform *significantly* better than literal ones.