

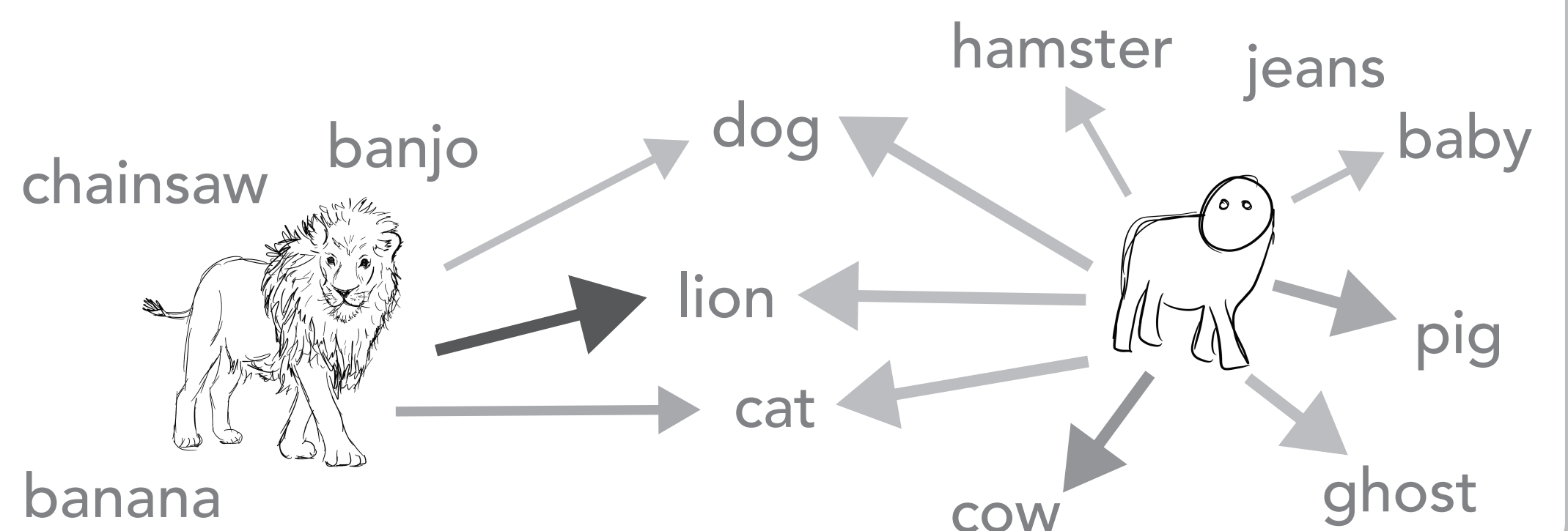
Evaluating machine comprehension of sketch meaning at different levels of abstraction



⁵ Kushin Mukherjee¹ Xuanchen Lu² Holly Huey² Yael Vinker³ Rio Aguina-Kang² Ariel Shamir⁴ Judith E. Fan^{2,5}
¹University of Wisconsin-Madison, ²University of California San Diego, ³Tel-Aviv University, ⁴Reichman University, ⁵Stanford University

QUESTION

How well do vision models exhibit human-like understanding of sketches that vary in semantic ambiguity?



METHODS

1. Vision Model Selection

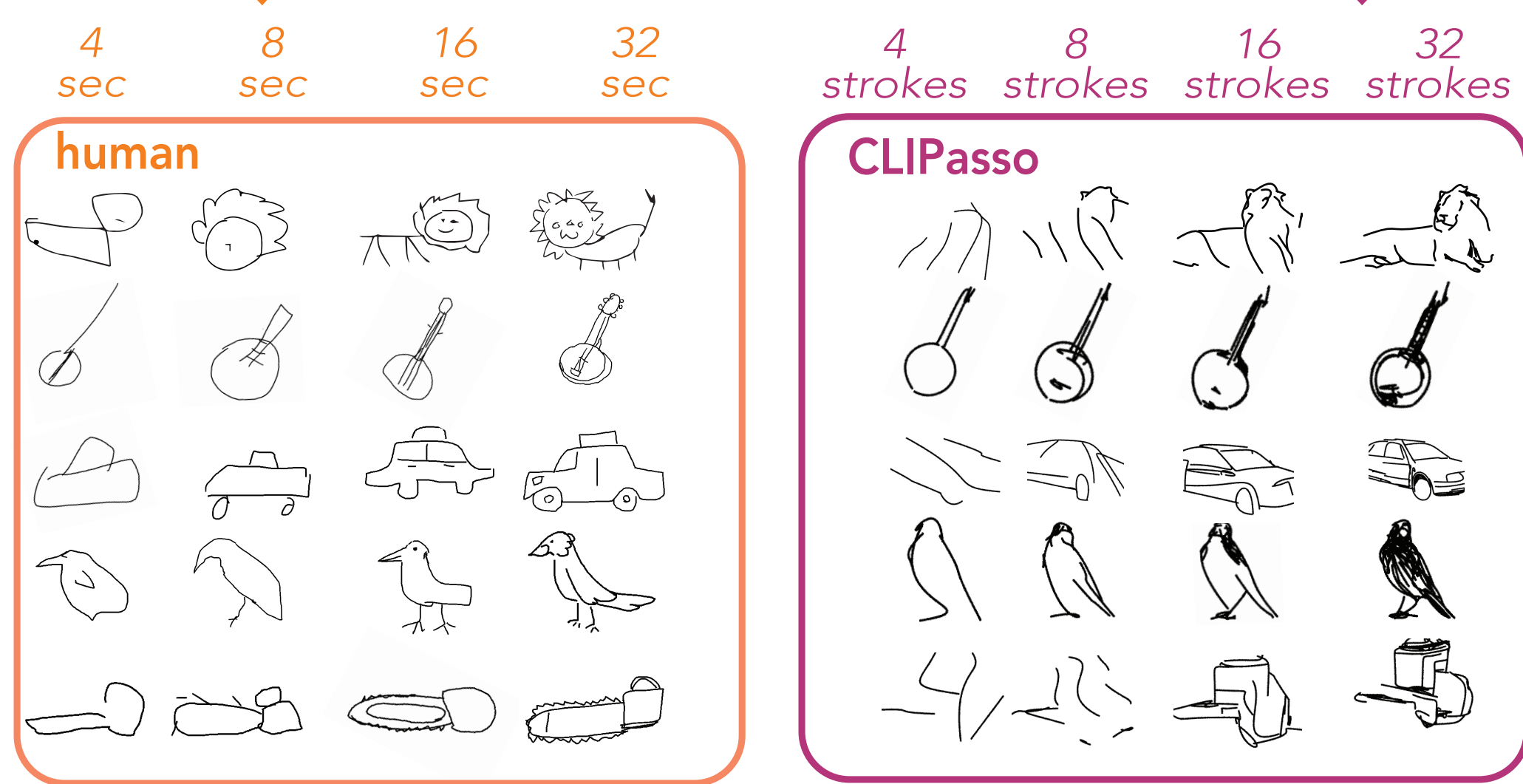
	supervised	self-supervised	semi-supervised
convnet	Inception-V3 VGG-19 ResNet-50 ECOSet CORNet-S		Noisy Student SWSL
transformer & MLP	ViT-B Swin-B Harmonization MLP-Mixer-B	MoCo-V3 DINO CLIP MAE	SimCLR IPCL

2. Sketch Dataset Generation

2048 photos from 128 classes from the THINGS database.

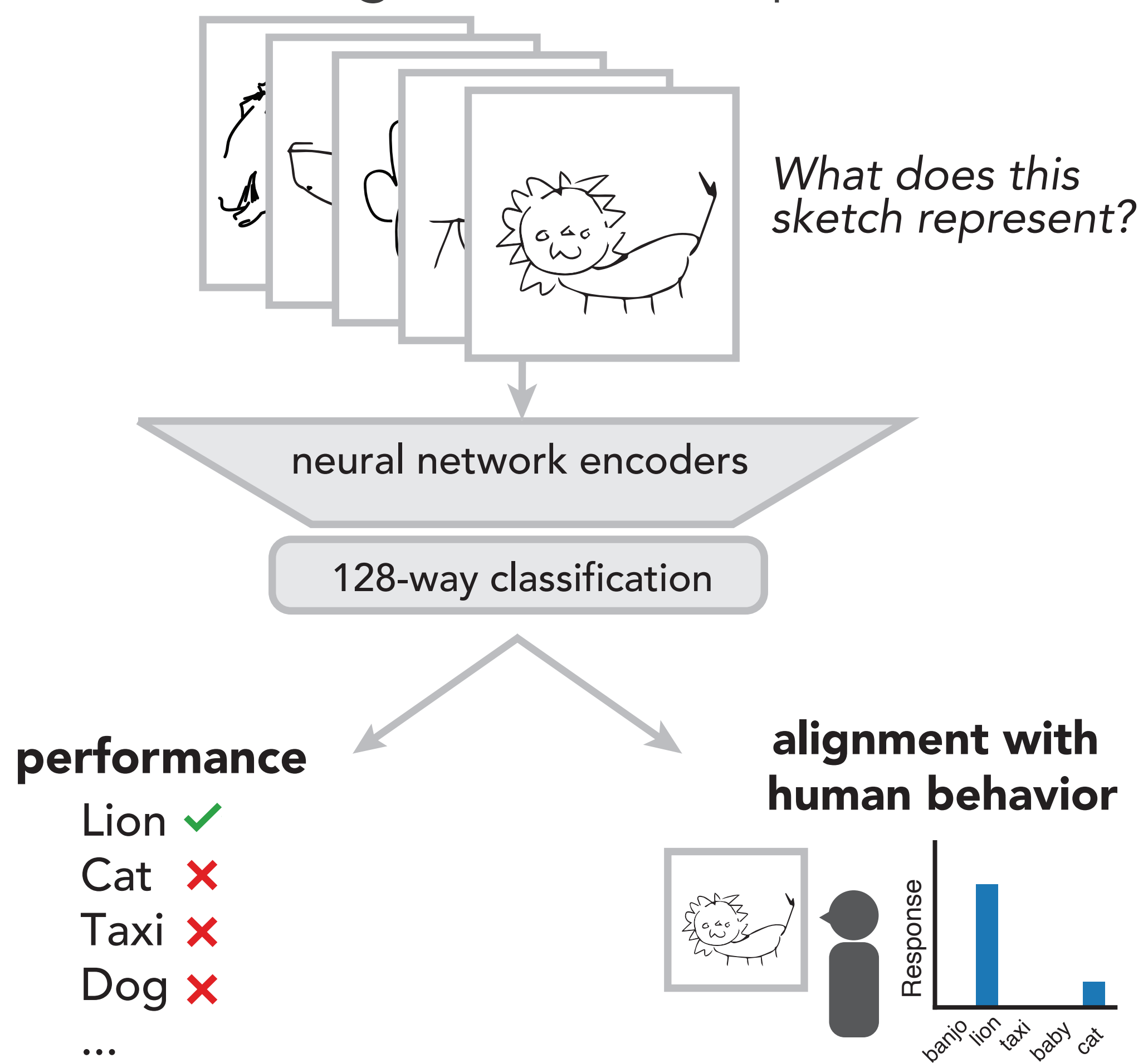


N = 5,563



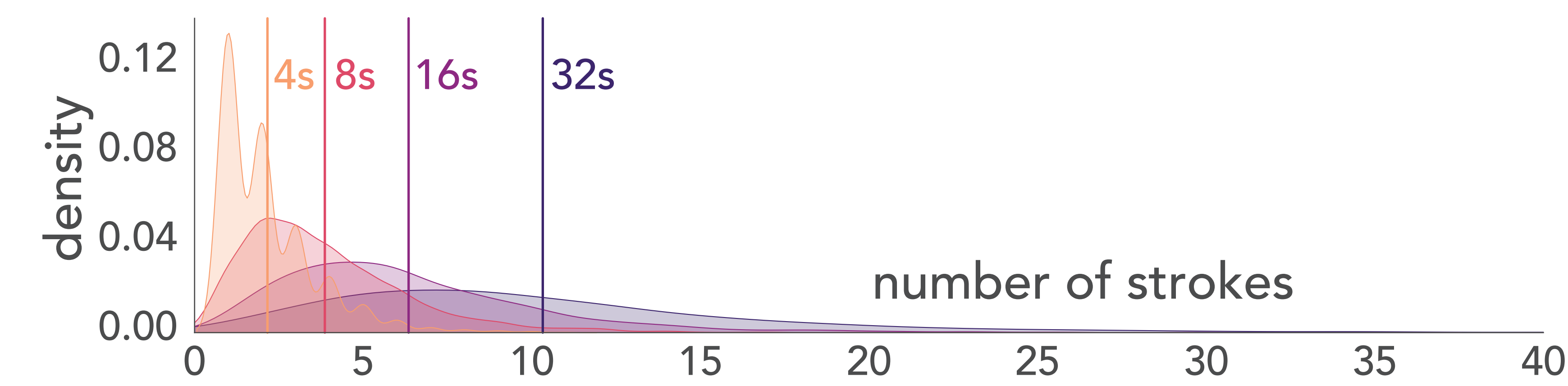
Over 90K sketches including sketches of each concept at 4 levels of detail made by humans and CLIPasso.

3. Measuring Sketch Comprehension



RESULTS

Humans produce sparser sketches under stronger time constraints.

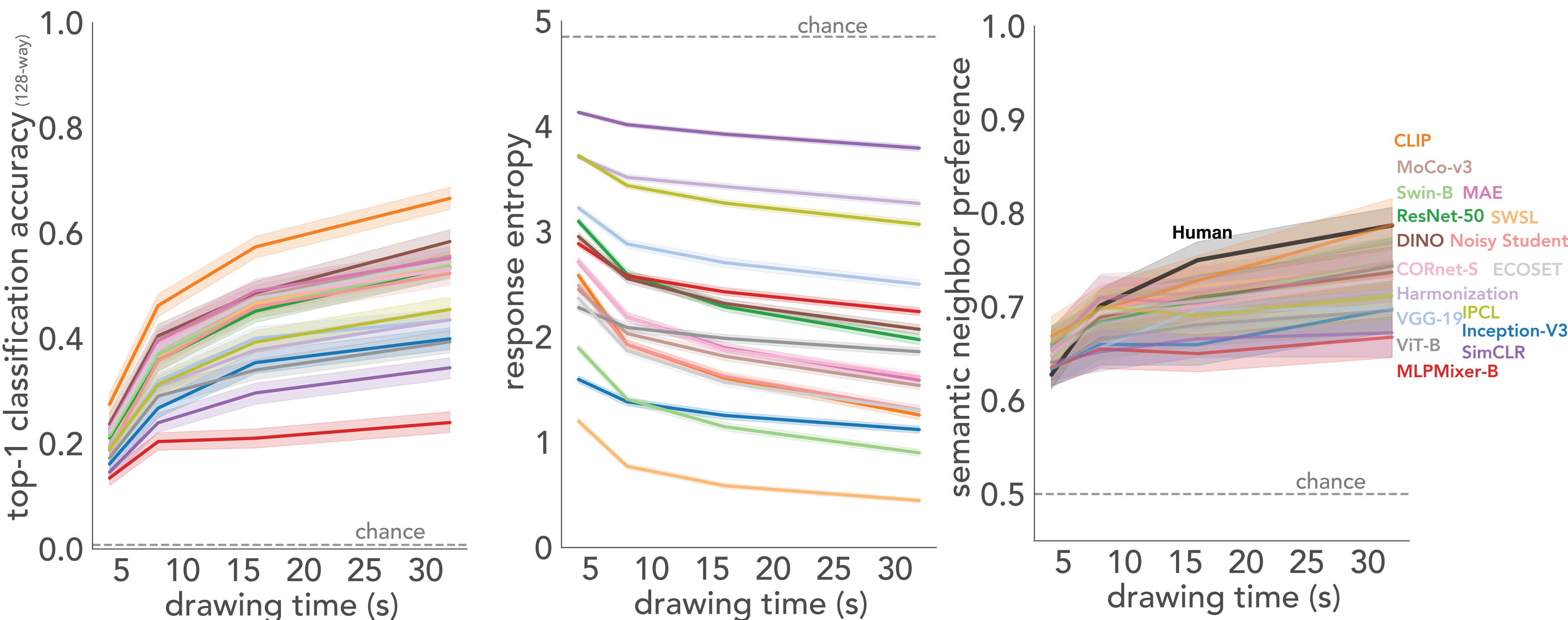


Sparser sketches more semantically ambiguous for models and humans.

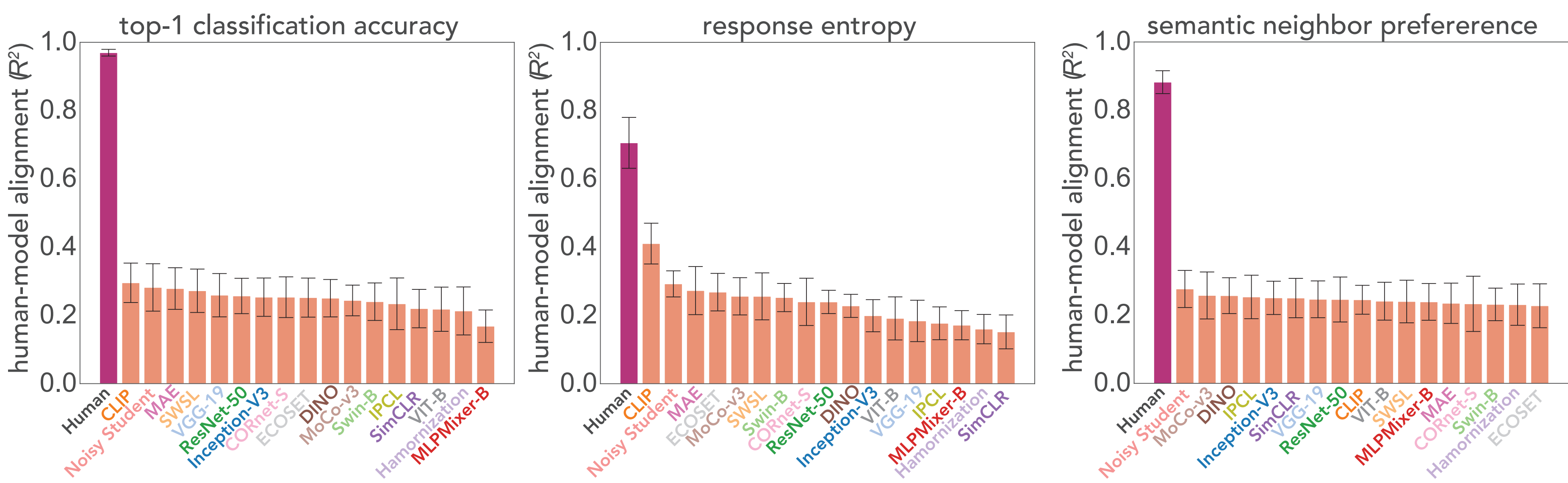
Sketches that are more detailed elicit ...

...higher classification accuracy. ...less variable response labels.

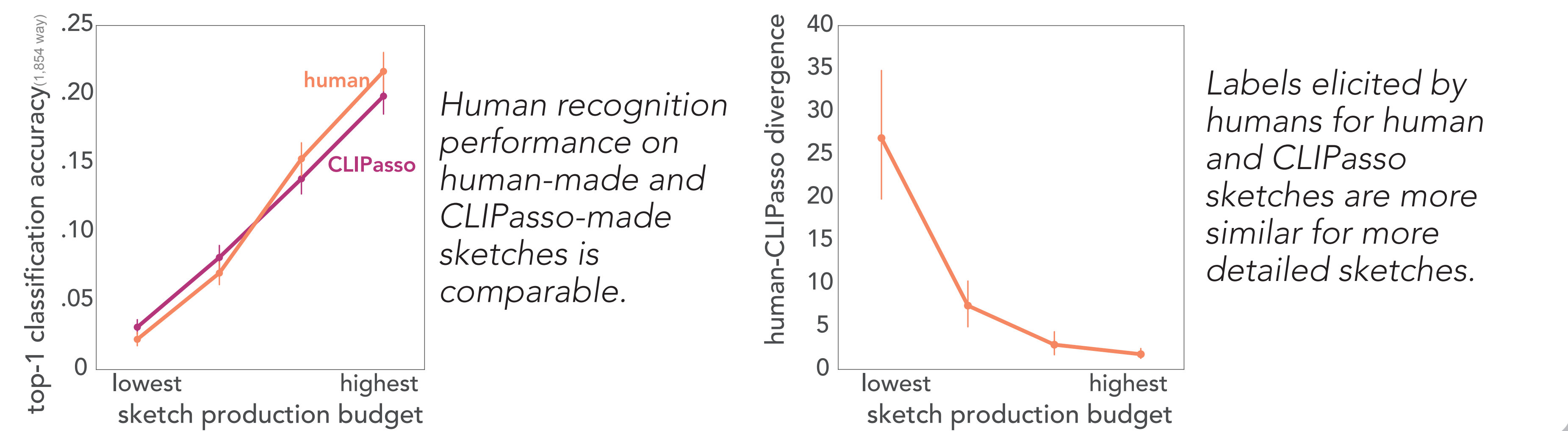
...guesses that are semantically close to the true label when incorrect.



Models vary in their degree of alignment to human behavior but a large gap remains between human and model sketch understanding.



A CLIP-based sketch generation algorithm emulates human sketches at greater levels of detail



TAKEAWAYS

We introduce a new dataset of >90K sketches at varied abstraction levels made by humans & CLIPasso, an AI-sketch generation model.

State-of-the-art vision models are sensitive to variation in the semantic information conveyed by sketches under different production budgets.

A large alignment gap still remains between the most performant vision models and humans.

CLIPasso-generated and human-made sketches elicit similar responses at greater levels of detail.

correspondence to:
kmukherjee2@wisc.edu
data & materials available at:
https://github.com/cogtoolslab/visual-abstractions-benchmarking_public2023
paper: