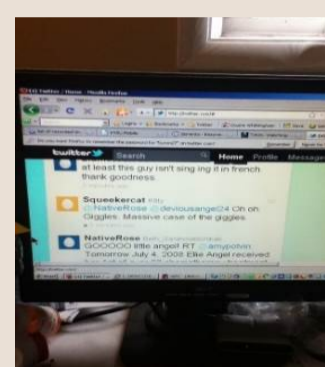


Why Does a Visual Question Have Different Answers?

Nilavra Bhattacharya ¹, Qing Li ², Danna Gurari ¹

¹ The University of Texas at Austin, ² University of California, Los Angeles

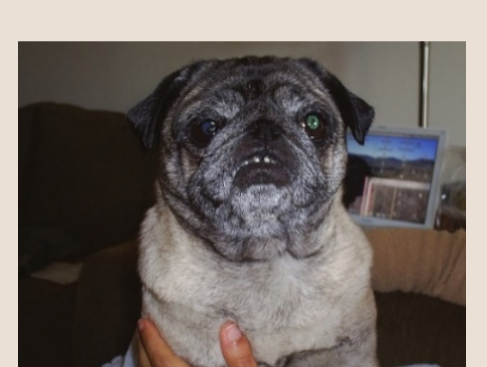
Motivation: Different People Can Offer Different Answers to a Visual Question



Is my monitor on?

Answers

- (1) yes (6) yes
- (2) yes (7) yes
- (3) yes (8) yes
- (4) yes (9) yes
- (5) yes (10) yes



Does this picture look scary?

Answers

- (1) yes (6) yes
- (2) no (7) yes
- (3) no (8) no
- (4) yes (9) no
- (5) no (10) no

- Approximately a 50/50 split for when a single versus multiple answers are observed for ~500,000 visual questions in three datasets [1]

Goal: Understand Why Answers Differ

- Benefits include empowering system designers and users with knowledge in how to interpret, prevent, and resolve answer differences

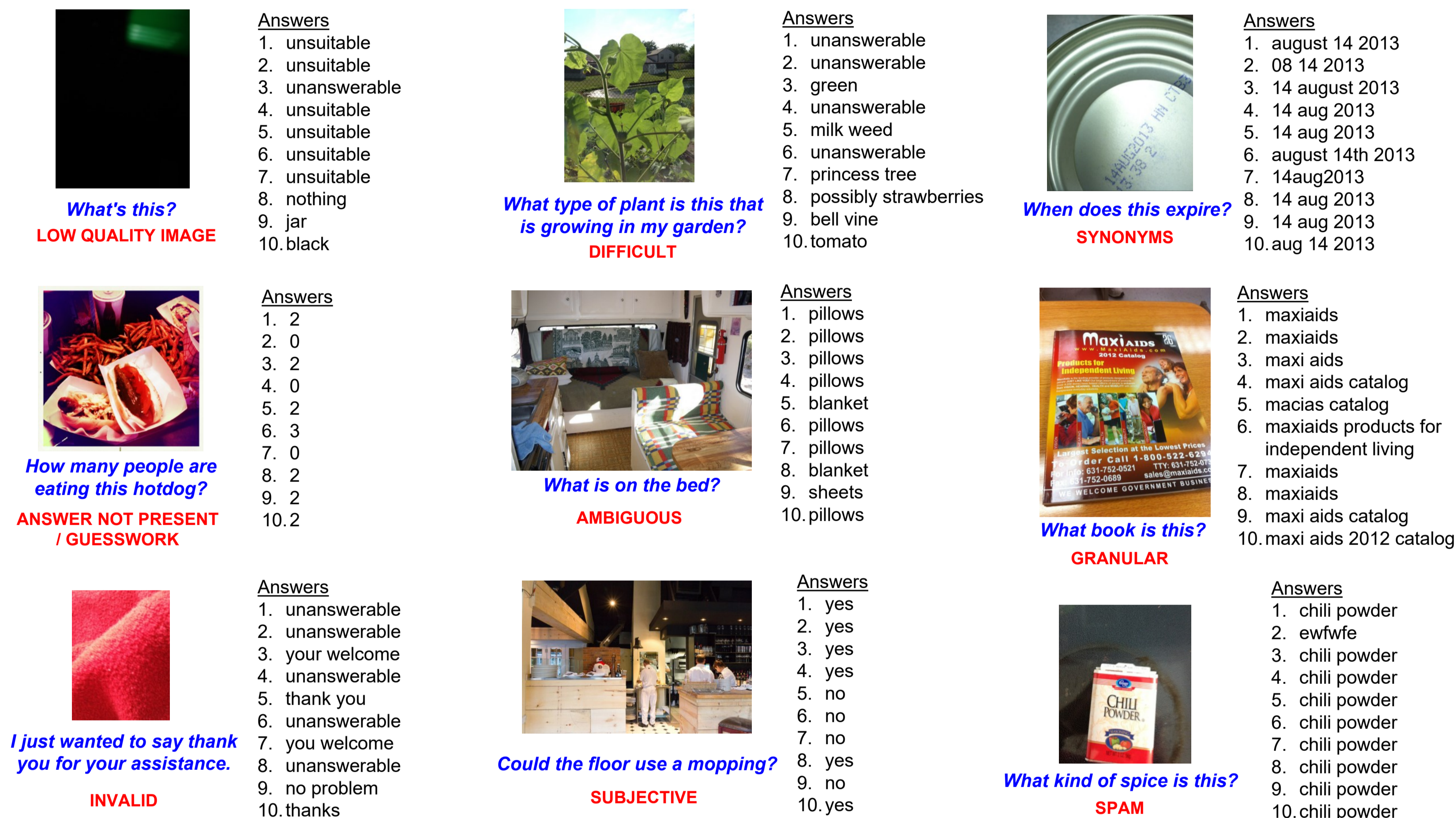
- More generally, support developing algorithmic frameworks that can account for the diversity of perspectives in a crowd

References

- [1] D. Gurari and K. Grauman. CrowdVerge: Predicting if People Will Agree on the Answer to a Visual Question. CHI 2017.
- [2] D. Gurari et al. VizWiz Grand Challenge: Answer Visual Questions from Blind People. CVPR 2018.
- [3] Y. Goyal et al. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. CVPR 2017
- [4] A. Mahendru et al. The Promise of Premise: Harnessing the Question Premises in Visual Question Answering. EMNLP 2017.

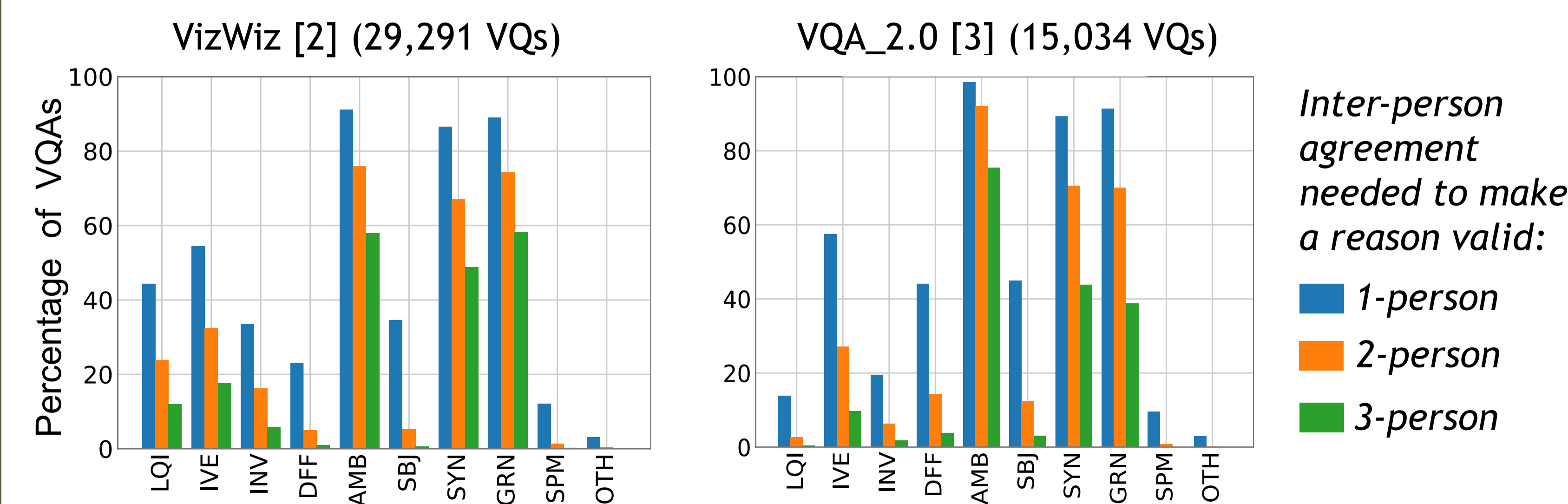
Answer-Differences Dataset

Taxonomy of nine reasons for why answers can differ:



- LOW QUALITY IMAGE**: "What's this?"
- DIFFICULT**: "What type of plant is this that is growing in my garden?"
- SYNONYMS**: "When does this expire?"
- ANSWER NOT PRESENT / GUESSWORK**: "How many people are eating this hotdog?"
- AMBIGUOUS**: "What is on the bed?"
- GRANULAR**: "What book is this?"
- INVALID**: "I just wanted to say thank you for your assistance."
- SUBJECTIVE**: "Could the floor use a mopping?"
- SPAM**: "What kind of spice is this?"

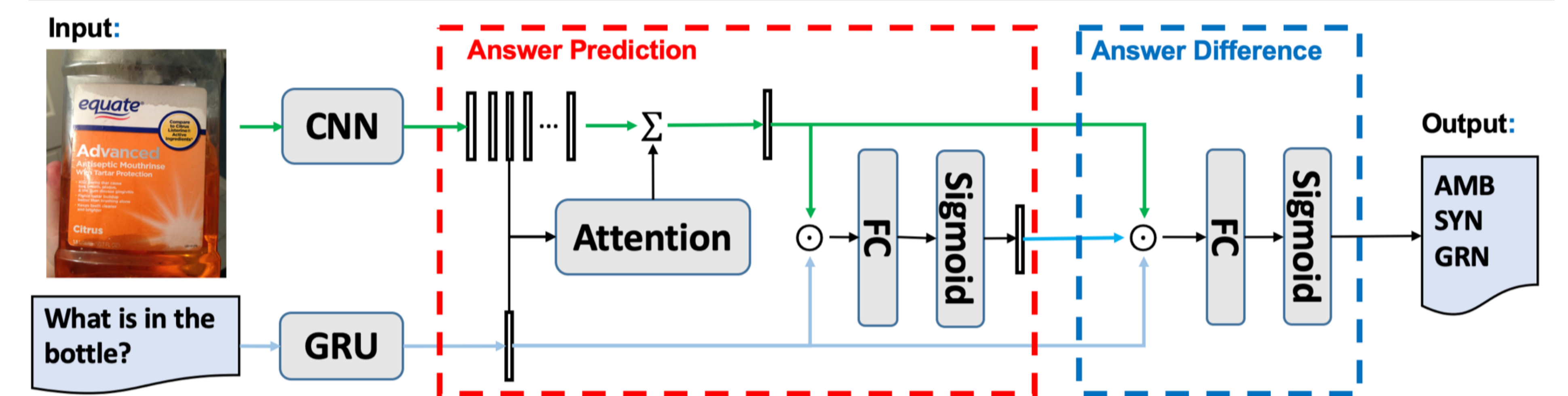
Frequency of each reason leading to answer differences:



Most common reasons: ambiguity, synonymous answers, & varying levels of answer granularity
Least common reason: spam (~1% of VQs)

Novel Task and Algorithm: Given a Visual Question, Predict the Reason(s) Why Answers from a Crowd Will Differ

Algorithm architecture: predicts with **image**, **question**, & **predicted answers**



Evaluation: mean average precision for all reasons

Model	Overall	Model	Overall
Random	30.15	Random	30.24
QI-Relevance [4]	31.71	QI-Relevance [4]	32.23
VQ answerability [2]	35.31	I	31.88
I	40.54	Q	43.47
Q	40.5	Q+I	43.16
Q+I	45.73	Q+I+A	44.55
Q+I+A	50.02	Q+I+A_FT	44.46
Q+I+A_FT	50.01	Q+I+A_GT	44.09
Q+I+A_GT	50.68		

Random: the best a user can achieve today
QI-Relevance: when VQ is irrelevant, then LQI, IVE, & AMB are reasons
VQ answerability: when VQ is unanswerable, then LQI, IVE, & AMB are reasons
FT: fine-tuned network
GT: uses ground truth instead of predicted answers

Algorithms can anticipate why a crowd will offer different answers, outperforming related baselines by >12 percentage points! The models perform best for ambiguity, synonymous answers, & varying answer granularity and perform worst for subjective & difficult VQs