

Exploring Models and Data for Image Question Answering

Mengye Ren, Ryan Kiros, Richard S. Zemel

University of Toronto



Computer Science
UNIVERSITY OF TORONTO

Problem

- Image Question Answering (QA): given an image and a free-form question, find an answer.
- We assume that answers are one-word, thus we can treat it as a classification problem.

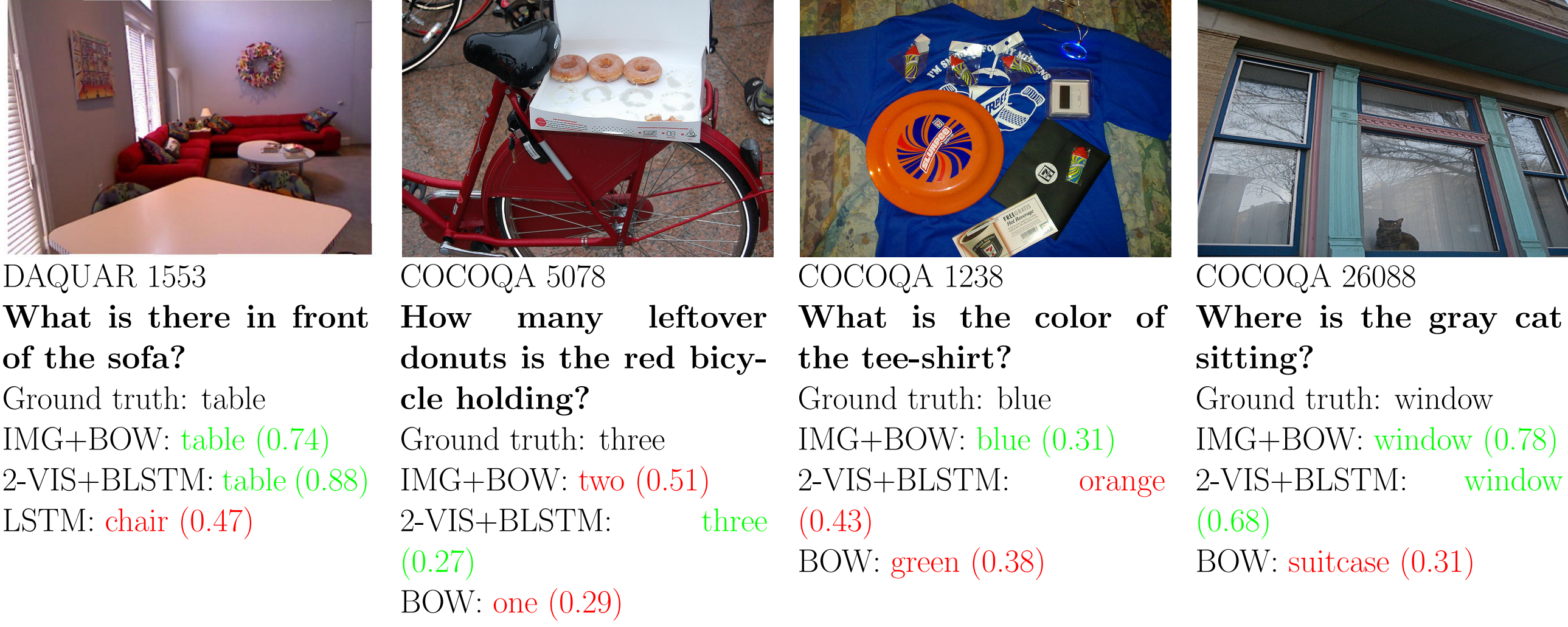


Figure: Sample questions and responses of a variety of models. Correct answers are in green and incorrect in red. The numbers in parentheses are the probabilities assigned to the top-ranked answer by the given model. The leftmost example is from the DAQUAR dataset, and the others are from our new COCO-QA dataset.

Our Models

- VIS+LSTM Model**

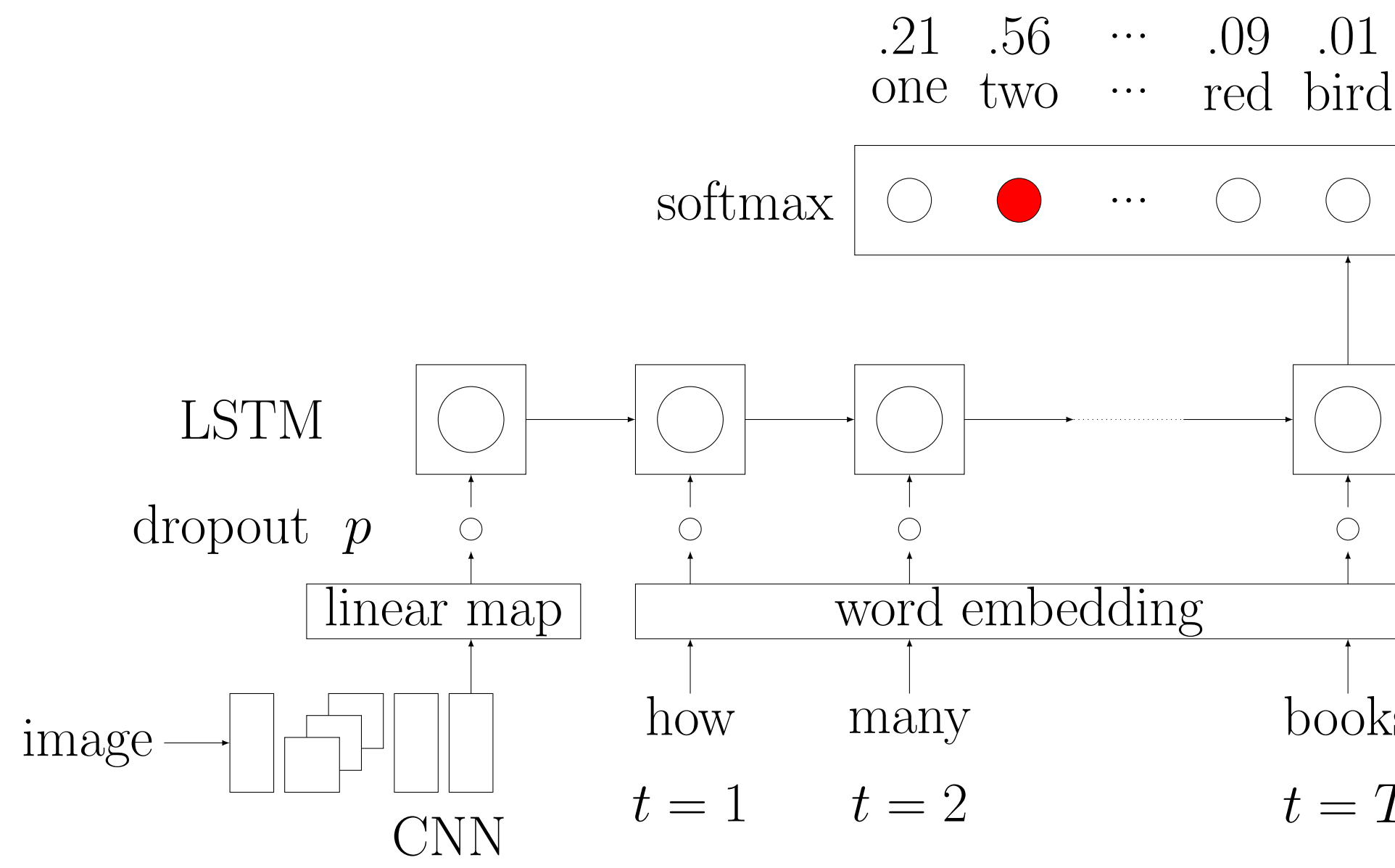


Figure: VIS+LSTM Model

- Borrowed the idea of treating the image as a word from previous caption generation work done by Vinyals et al. [1].
- Image features are passed through a linear transformation to match with word vector dimensions.
- At the last time step, the Long Short Term Memory (LSTM) [2] outputs to the softmax layer to classify answers.
- 2-VIS+BLSTM Model**
 - Two image feature inputs, at the start and the end of the question, with different learned linear transformations.
 - LSTMs going in both the forward and backward directions.
 - Both LSTMs output to the softmax layer.
- IMG+BOW** Multinomial logistic regression based on the CNN image features (4096 dimension), and learned bag-of-word (BOW) vectors.
- FULL** A simple average of the three models above.

Automatic QA Generation

- Motivation** Currently available dataset DAQUAR [3] is very small (1500 images, 7000 QA on 37 classes of objects, 12000 QA on 894 classes of objects). Guessing the modes can yield very good accuracy, and “blind models” can achieve almost equal performance compared the best model.

- We generate 4 types of questions: object, number, color, location, directly from image description.
- All answers are one-word.
- We move the wh-word and verb to the front under certain constraints. For example, “A man is riding a **horse**” => “**What** is **the** man riding?”
- We pruned answers that appear too rarely or too often.

Baselines

- GUESS** Predict the mode based on the question type.
- BOW** Given only the questions without the images and perform logistic regression on the bag-of-words vector to classify answers.
- LSTM** Input the question words into the LSTM alone.
- IMG** Re-train a separate CNN classification layer for each type of question.
- IMG+PRIOR** Combine the prior knowledge of an object and the image understanding from the “deaf model”. Denote c as the color, o as the class of the object of interest, and x as the image. Assuming o and x are conditionally independent given the color. Use the output of the IMG model as $p(c|x)$.
$$p(c|o, x) = \frac{p(o|c)p(c|x)}{\sum_{c \in C} p(o|c)p(c|x)}$$
Empirical estimate: $\hat{p}(o|c) = \frac{count(o,c)}{count(c)}$ and Laplace smoothing.

Experimental Results

Table: COCO-QA question type break-down

Category	Train	%	Test	%
Object	54992	69.84%	27206	69.85%
Number	5885	7.47%	2755	7.07%
Color	13059	16.59%	6509	16.71%
Location	4800	6.10%	2478	6.36%
Total	78736	100.00%	38948	100.00%

Table: DAQUAR and COCO-QA results

	DAQUAR			COCO-QA		
	Acc.	WUPS 0.9	WUPS 0.0	Acc.	WUPS 0.9	WUPS 0.0
MULTI-WORLD [4]	0.1273	0.1810	0.5147	-	-	-
GUESS	0.1824	0.2965	0.7759	0.0665	0.1742	0.7344
BOW	0.3267	0.4319	0.8130	0.3752	0.4854	0.8278
LSTM	0.3273	0.4350	0.8162	0.3676	0.4758	0.8234
IMG	-	-	-	0.4302	0.5864	0.8585
IMG+PRIOR	-	-	-	0.4466	0.6020	0.8624
IMG+BOW	0.3417	0.4499	0.8148	0.5592	0.6678	0.8899
VIS+LSTM	0.3441	0.4605	0.8223	0.5331	0.6391	0.8825
ASK-NEURON [5]	0.3468	0.4076	0.7954	-	-	-
2-VIS+BLSTM	0.3578	0.4683	0.8215	0.5509	0.6534	0.8864
FULL	0.3694	0.4815	0.8268	0.5784	0.6790	0.8952
HUMAN	0.6027	0.6104	0.7896	-	-	-

Table: COCO-QA accuracy per category

	Object	Number	Color	Location
GUESS	0.0211	0.3584	0.1387	0.0893
BOW	0.3727	0.4356	0.3475	0.4084
LSTM	0.3587	0.4534	0.3626	0.3842
IMG	0.4073	0.2926	0.4268	0.4419
IMG+PRIOR	-	0.3739	0.4899	0.4451
IMG+BOW	0.5866	0.4410	0.5196	0.4939
VIS+LSTM	0.5653	0.4610	0.4587	0.4552
2-VIS+BLSTM	0.5817	0.4479	0.4953	0.4734
FULL	0.6108	0.4766	0.5148	0.5028

More Examples

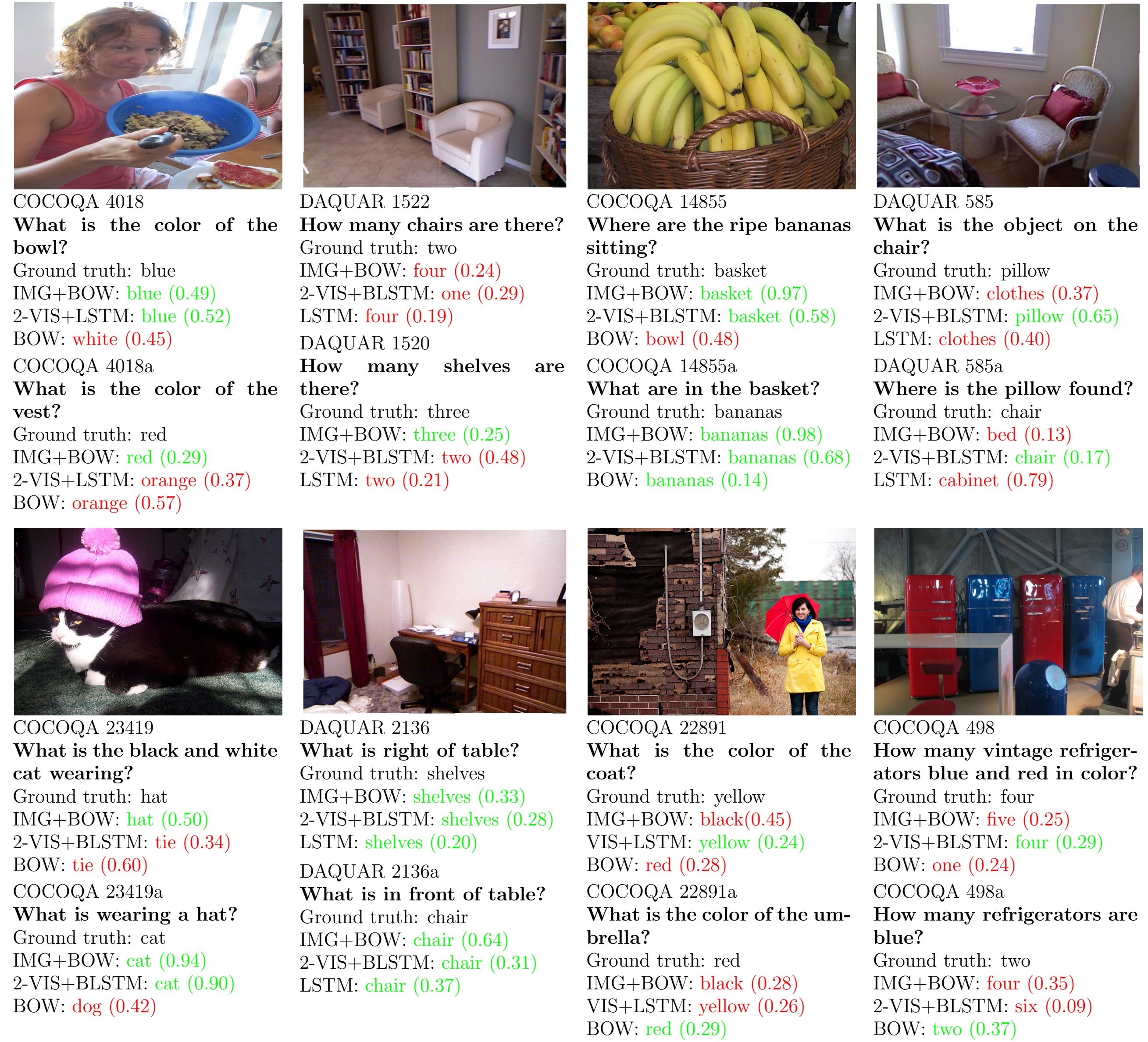


Figure: Sample questions and responses of our system. For some of the examples, we specifically tested extra questions (the ones have an “a” in the question ID).

- Full results: <http://www.cs.toronto.edu/~mren/imageqa/results>

Conclusion

- We present our end-to-end neural network models to image QA.
- Simple bag-of-words can perform equally well compared to recurrent neural network.
- Models have large space for improvement on questions such as color and counting.
- We release an Image QA dataset that is automatically generated from image description. Download: <http://www.cs.toronto.edu/~mren/imageqa/data/cocoqa>

Current Directions

- Free-form text generation model. Similar to image captioning, it will require an automatic free-form answer evaluation metric.
- Extend questions to open domain.
- Use of visual attention to improve results and interpret output (based on recent successes of visual attention in image captioning [6]).

References

- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- Mateusz Malinowski and Mario Fritz. Towards a visual turing challenge. In *NIPS Workshop on Learning Semantics*, 2014.
- Mateusz Malinowski and Mario Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014.
- Mateusz Malinowski, Marcus Rohrbach, and Mario Fritz. Ask Your Neurons: A Neural-based Approach to Answering Questions about Images. *CoRR*, abs/1505.01121, 2015.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML (to appear)*, 2015.