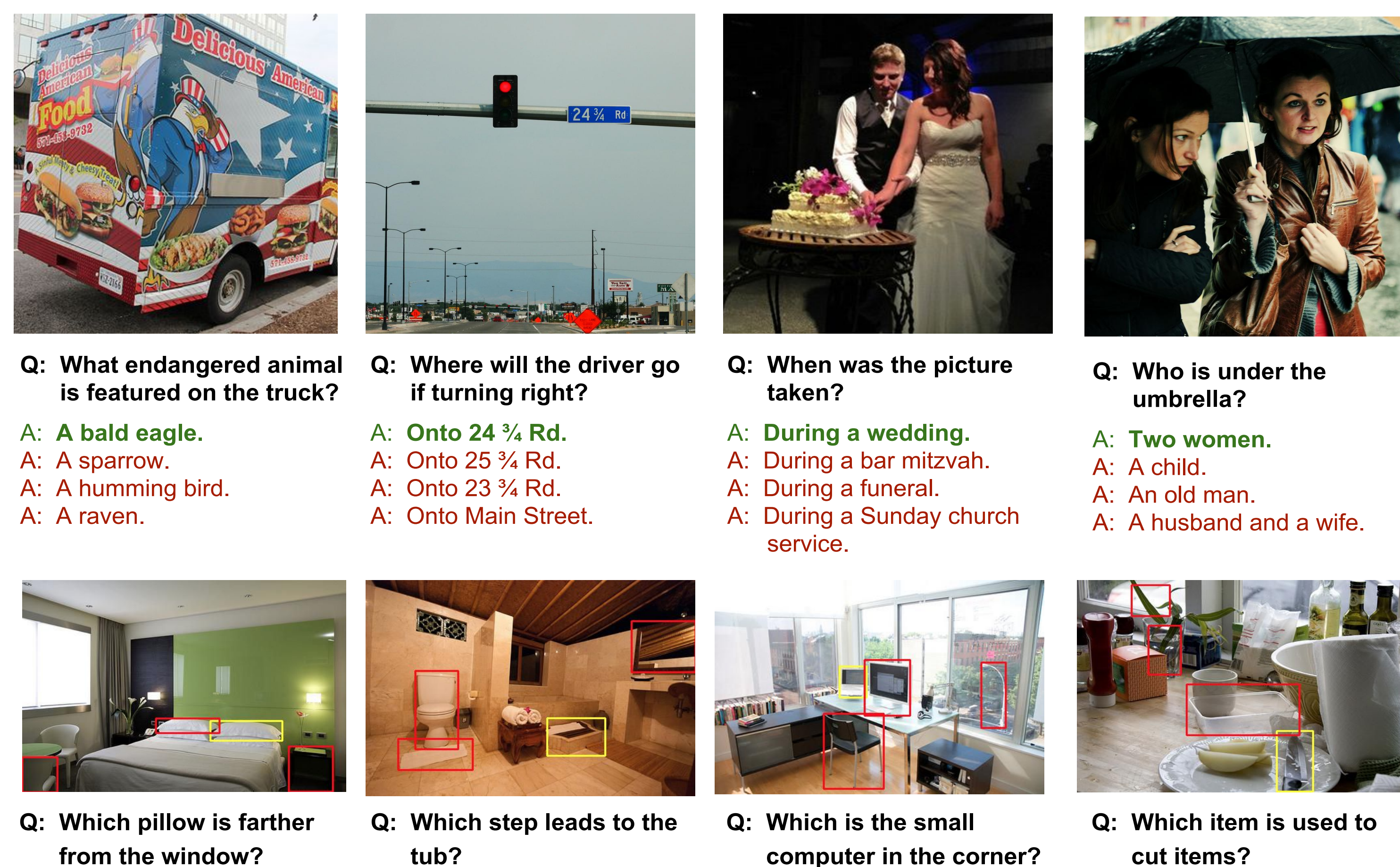


Visual7W Data



Telling QA

Question starts with who, what, where, when, why or how; answer consists of 4 textual multiple choices

Pointing QA

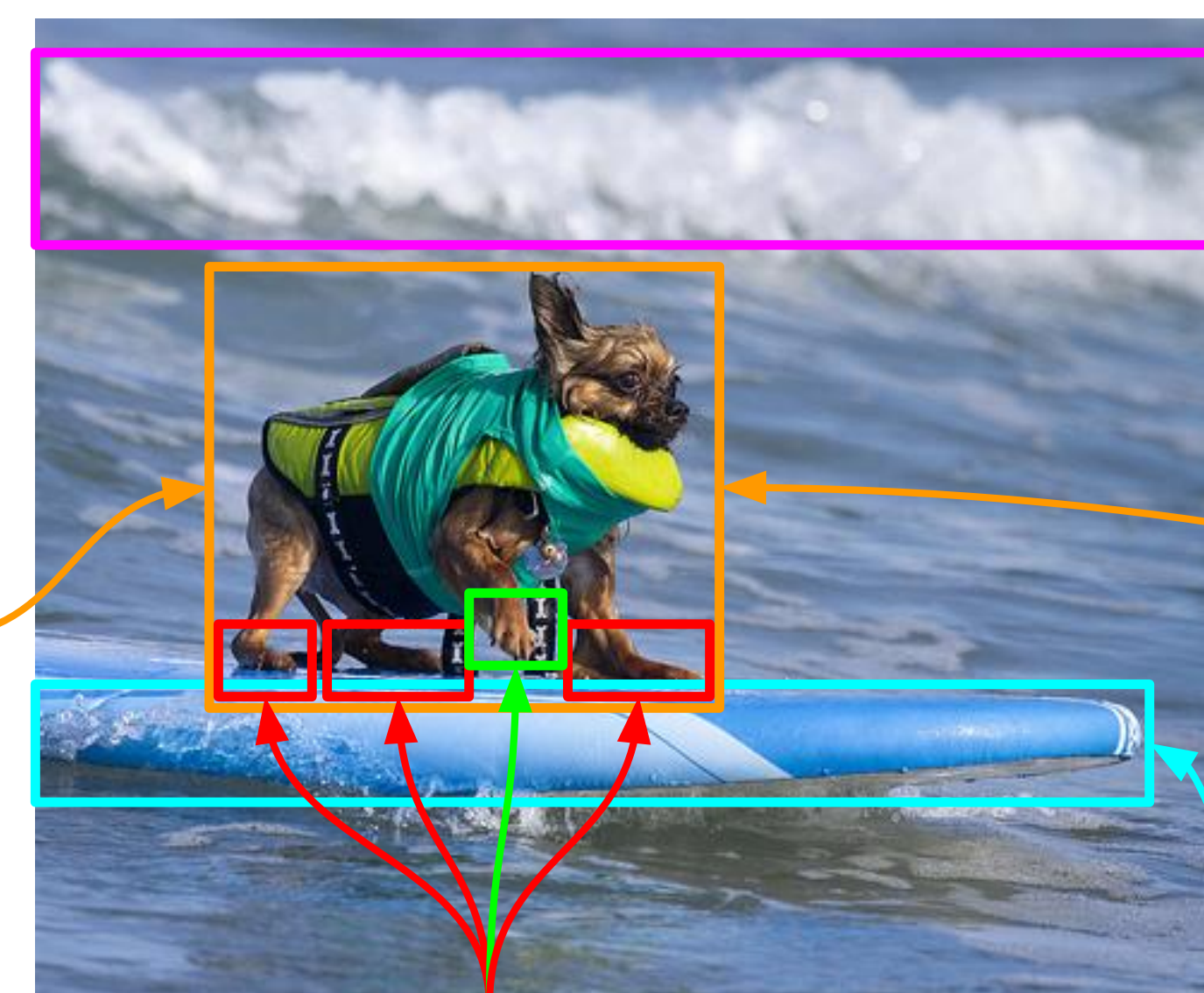
Question starts with which; answer consists of 4 visual bounding boxes in the image

Where does this scene take place?

A) In the sea. ✓
B) In the desert.
C) In the forest.
D) On a lawn.

What is the dog doing?

A) Surfing. ✓
B) Sleeping.
C) Running.
D) Eating.



Why is there foam?

A) Because of a wave. ✓
B) Because of a boat.
C) Because of a fire.
D) Because of a leak.

What is the dog standing on?

A) On a surfboard. ✓
B) On a table.
C) On a garage.
D) On a ball.

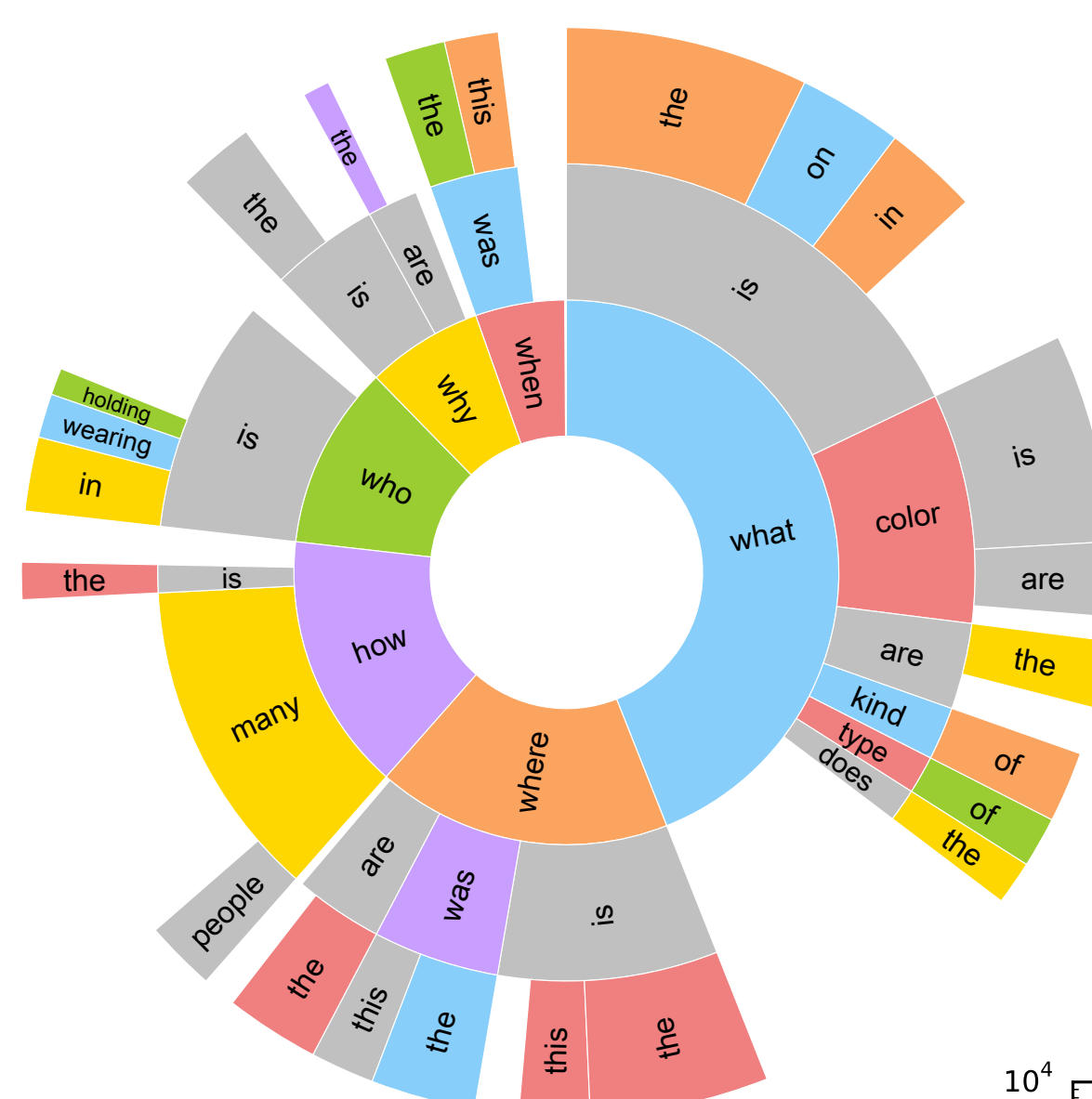
Which paw is lifted?

Statistics

- 47,300 COCO images
- 327,939 QA pairs: kkk, kkk telling QAs, kkk, kkk pointing QAs
- 1,311,756 multiple choices
- 561,459 object groundings
- 36,579 object categories
- avg. question length: 6.9 +/- 2.4
- avg. answer length: 2.0 +/- 1.4
- 27.6% of all answers >2 words
- 1,000 most frequent answers account for 63.5% of all answers

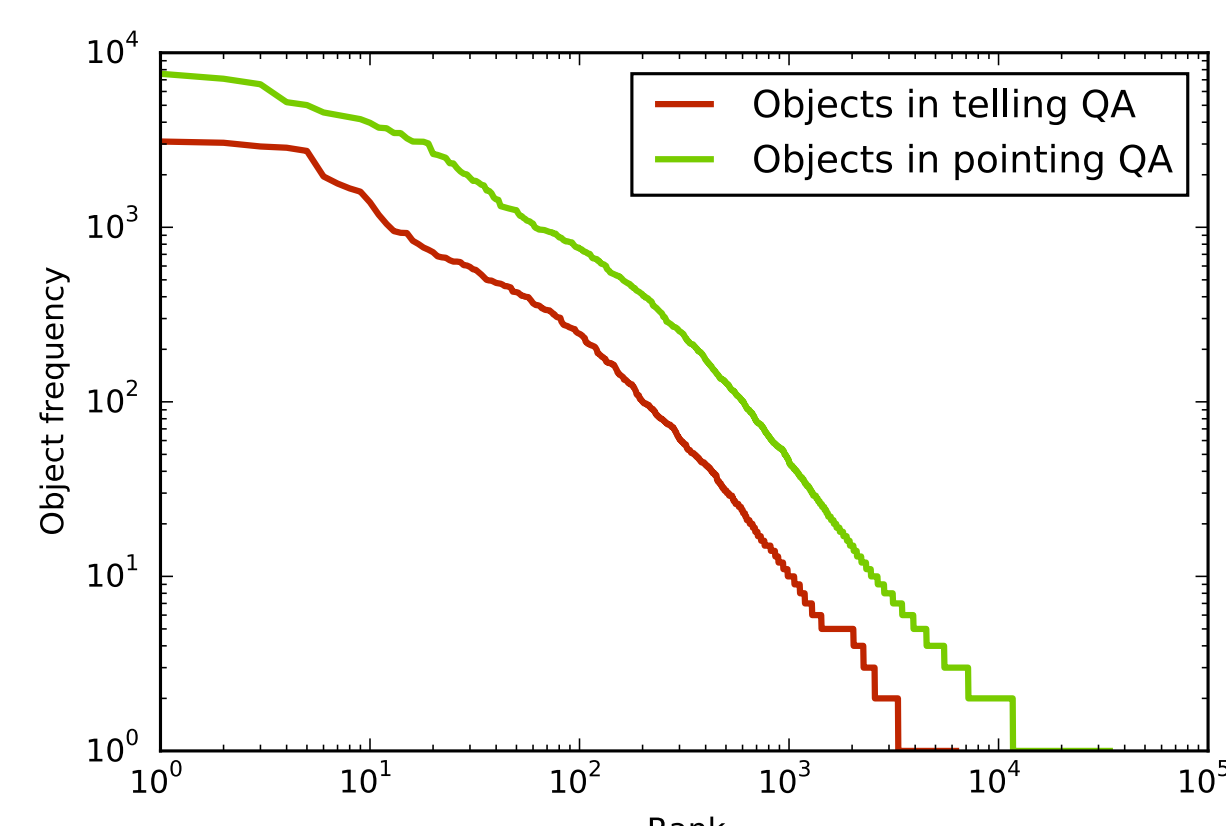
Question Distribution

telling questions grouped by first 3 words (sparse arcs removed from sunburst)

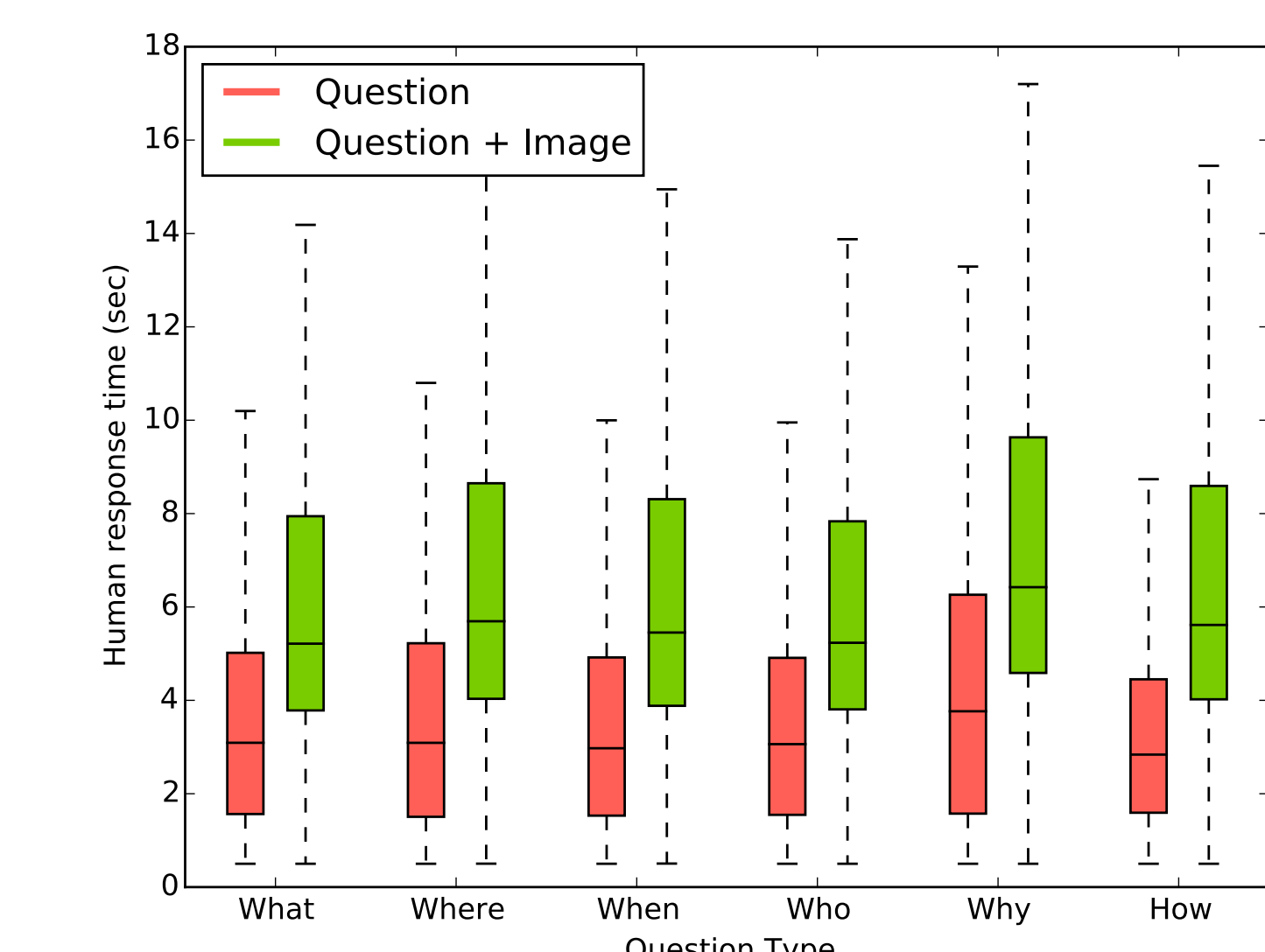


Object Distribution

pointing questions refer to significantly more object categories



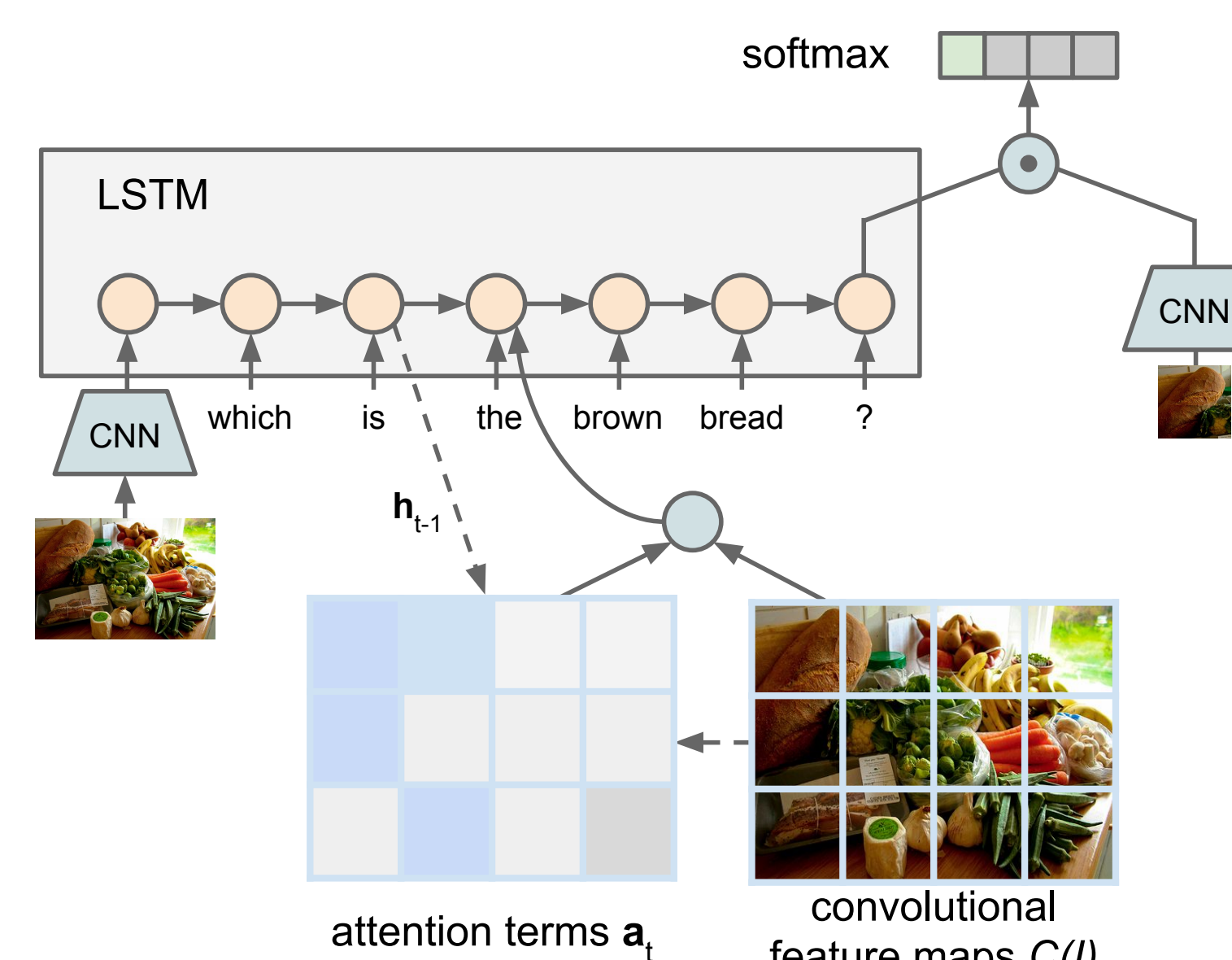
- Top 20 objects**
- man
 - trees
 - shirt
 - wall
 - grass
 - tree
 - woman
 - people
 - sky
 - building
 - person
 - window
 - table
 - head
 - water
 - snow
 - hair
 - sign
 - ground
 - wood



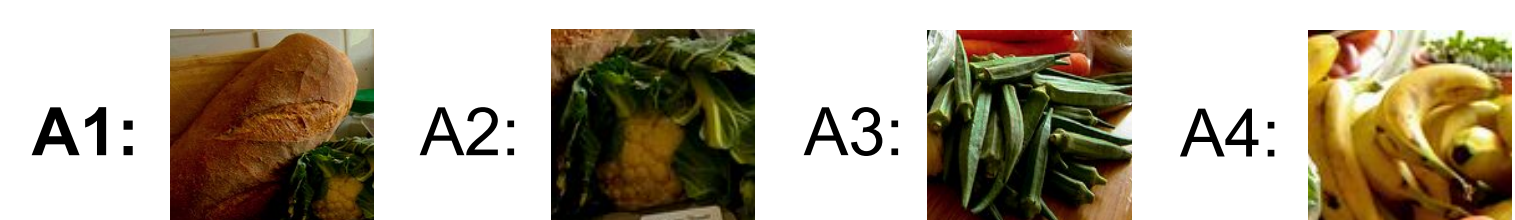
Human Response Times

significantly longer response delays when image absent across all telling QAs; image and QA tightly coupled

Spatial Attention LSTM



Q: Which is the brown bread?



Attention Terms: $w_a^T, W_{he}, W_{ce}, b_a$ trainable
 $a_t = \text{softmax}(e_t)$
 $e_t = w_a^T \tanh(W_{he} h_{t-1} + W_{ce} C(I)) + b_a$
 h_{t-1} hidden state vector
 $C(I)$ convolotional image features

Experiments & Results

Image	Multiple Choices	Q: Who is behind the batter?	Q: What adorns the tops of the post?	Q: What kind of stuffed animal is shown?	Q: What animal is being petted?
		A1: Catcher. A2: Umpire. A3: Fans. A4: Ball girl.	A1: Gulls. A2: An eagle. A3: A crown. A4: A pretty sign.	A1: Teddy Bear. A2: Monkey. A3: Tiger. A4: Bunny rabbit.	A1: A sheep. A2: Goat. A3: Alpaca. A4: Pig.
w/o Image		H: Catcher. ✓ M: Umpire. ✗	H: Gulls. ✓ M: Gulls. ✓	H: Monkey. ✗ M: Teddy Bear. ✓	H: A sheep. ✓ M: A sheep. ✓
w/ Image		H: Catcher. ✓ M: Catcher. ✓	H: Gulls. ✓ M: A crown. ✗	H: Teddy Bear. ✓ M: Teddy Bear. ✓	H: Goat. ✗ M: A sheep. ✓

Top Grid: qualitative results of humans (H) and our model (M) in the V7W multiple choice benchmark with and without showing images while asking

Bottom Table: comparison of model accuracies per question category and overall on the V7W benchmark split

Q indicates question being shown; I indicates image being shown

LogRegr: logistic regression model, classifies concatenation of embeddings

LSTM: model proposed by Malinowski et al., ICCV'15 [2]

LSTM-Att: our spatial attention LSTM, trained on V7W

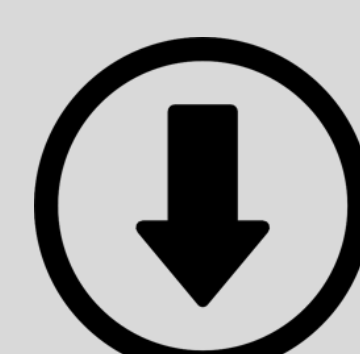
LSTM-Att+: our spatial attention LSTM, trained on VG-QA [4]

Method	What	Where	When	Who	Why	How	Which	Overall
Human (Q)	0.356	0.322	0.393	0.342	0.439	0.337	-	0.353
Human (Q+I)	0.965	0.957	0.944	0.965	0.927	0.942	0.973	0.966
LogRegr (Q)	0.420	0.375	0.666	0.510	0.354	0.458	0.354	0.383
LogRegr (I)	0.408	0.426	0.438	0.415	0.337	0.303	0.256	0.305
LogRegr (Q+I)	0.429	0.454	0.621	0.501	0.343	0.356	0.307	0.352
LSTM (Q)	0.430	0.414	0.693	0.538	0.491	0.484	-	0.462
LSTM (I)	0.422	0.497	0.660	0.523	0.475	0.468	0.299	0.359
LSTM (Q+I)	0.489	0.544	0.713	0.581	0.513	0.503	0.521	0.521
LSTM-Att (Q+I)	0.529	0.560	0.743	0.602	0.522	0.466	0.561	0.541
LSTM-Att+ (Q+I)	0.572	0.613	0.760	0.624	0.590	0.531	-	0.587

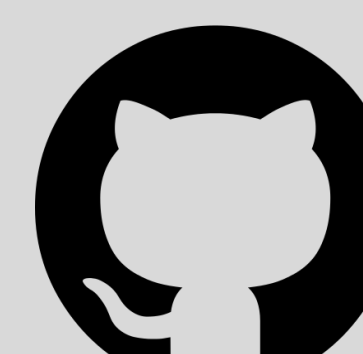


Project Page

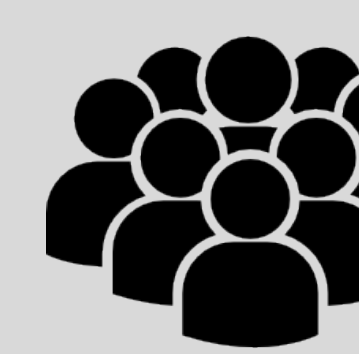
<http://web.stanford.edu/~yukez/visual7w/>



Data
Images, QAs
& Grounding



Code
Evaluation Toolkit
& Torch Model



Crowdsourcing
AMT Interfaces
available

All collected QA data contributed to visualgenome.org
VG-QA (total) [4]: 101,174 Images; 1,773,258 QA Pairs

