## **Question Answering about Images using Visual Semantic Embeddings**

Q212: what is the object close to the wall right and left of the cabinet? Ground truth: television CNN-LSTM: television (0.3950) BLIND-LSTM: window (0.5004) Complex question and recognized correct object



Q464: what is the object close to the floor left of the wall?
Ground truth: toilet
CNN-LSTM: toilet (0.2292)
BLIND-LSTM: cabinet (0.1588)
Small object



Q585: what is the object on the chair? Ground truth: pillow CNN-LSTM: pillow (0.6475) BLIND-LSTM: clothes (0.3509) Location of object



Q3266: what is to the left of the door way?

Ground truth: lamp CNN-LSTM: lamp (0.3301) BLIND-LSTM: table (0.1434)

Small object



Q2346: what is to left of door? Ground truth: sink CNN-LSTM: sink (0.1992) BLIND-LSTM: door (0.2205) Could have been "toilet" if the model does know left or right



Q232: what is the largest object? Ground truth: sofa CNN-LSTM: sofa (0.6274) BLIND-LSTM: bed (0.3660) Image features are useful mostly when asking about the largest object



Q2136: what is right of table? Ground truth: shelves CNN-LSTM: shelves (0.2780) BLIND-LSTM: window (0.2741) Left/right is good



Ground truth: television
CNN-LSTM: television (0.2888)
BLIND-LSTM: chair (0.1951)
Small object compared to bed and shelf.
But "on the table" gives some hint



Q969: what is on the left side of the table?
Ground truth: refridgerator
CNN-LSTM: refridgerator (0.7176)
BLIND-LSTM: chair (0.1836)
Left/right image feature helps

Table 1. DAQUAR Object Questions



Q705: what is the big object on the big dark brown dresser? Ground truth: television CNN-LSTM: picture (0.2295) BLIND-LSTM: television (0.2216) Fail and the blind model did very good reasoning



Q800: what is biggest black object?

Ground truth: sofa CNN-LSTM: television (0.5182) BLIND-LSTM: sofa (0.2108)

Probably thinks the big picture hanging

on the wall as television



Q876: what is on the table? Ground truth: books CNN-LSTM: lamp (0.3294) BLIND-LSTM: books (0.1915) Books are likely to be on the table



Q31: what is on the wall above the night stand?

Ground truth: picture CNN-LSTM: picture (0.2716) BLIND-LSTM: picture (0.5047) "on the wall" seems to be strong enough



Q1231: what is on the left side of the

Ground truth: night\_stand CNN-LSTM: night\_stand (0.3897) BLIND-LSTM: night\_stand (0.1630) Nightstands are likely to be beside the

bed



O2881: what is on the floor? Ground truth: floor\_mat CNN-LSTM: floor\_mat (0.4828) BLIND-LSTM: floor\_mat (0.4806) Floormats are always on the floor



Q3278: what are around the dining table?

Ground truth: chair CNN-LSTM: chair (0.9625) BLIND-LSTM: chair (0.9878) Chairs are always around the dining ta-

ble

signal



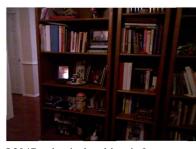
Q1466: what is on the night stand? Ground truth: paper CNN-LSTM: lamp (0.8925) BLIND-LSTM: lamp (0.8738) Very hard question and "lamp" is a more reasonable answer



Q2010: what is to the right of the table

Ground truth: sink CNN-LSTM: chair (0.7257) BLIND-LSTM: chair (0.4172) Very hard question that there are too many cluttered objects

Table 2. DAQUAR Object Questions



Q2047: what is the object is? Ground truth: shelves CNN-LSTM: sofa (0.3215) BLIND-LSTM: sofa (0.0317) Not very clear why CNN fails on this



Q2731: what is found on the floor? Ground truth: bed CNN-LSTM: floor\_mat (0.5103) BLIND-LSTM: floor\_mat (0.7409) CNN doesn't differ too much from blind and image features are not more useful than "common-sense"

Table 3. DAQUAR Object Questions



Q241: what is the largest object? Ground truth: sofa CNN-LSTM: table (0.6090) BLIND-LSTM: bed (0.3660) Again the model fail



Q35: how many drawer knobs are in this picture?

Ground truth: two CNN-LSTM: three (0.2127) BLIND-LSTM: two (0.2479) Very hard questio



Q1908: how many posters are on the

Ground truth: four CNN-LSTM: two (0.4634) BLIND-LSTM: two (0.2821)

Model simply doesn't count even for simple task like this



Q1520: how many shelves are there? Ground truth: three CNN-LSTM: two (0.4801) BLIND-LSTM: two (0.3325) Not counting well



Q1768: how many doors does the cupboard have?

Ground truth: two CNN-LSTM: two (0.5247) BLIND-LSTM: two (0.2680)

Not very obvious question mostly the model gets number question right by

guessing "two"



Q711: how many doors are there? Ground truth: two CNN-LSTM: two (0.3903) BLIND-LSTM: two (0.3339) Not very obvious again



Q1286: how many blinds are there? Ground truth: three CNN-LSTM: three (0.4044) BLIND-LSTM: two (0.3319) Kind of interestin



Q2989: what is the colour of the sofa?

Ground truth: red

CNN-LSTM: red (0.2152) BLIND-LSTM: brown (0.2422) Focus on remote object color



Q2511: what is the colour of the coffee

machine?

Ground truth: black CNN-LSTM: black (0.2621) BLIND-LSTM: white (0.1728)

Small object color



Q746: what color is the sofa?

Ground truth: brown

CNN-LSTM: brown (0.5405) BLIND-LSTM: brown (0.4483) Sofas are mostly brown



Q842: what is the colour dominating?

Ground truth: brown

CNN-LSTM: brown (0.2266) BLIND-LSTM: white (0.3483)

We expect the image features to contain information about dominant color



Q999: what is the colour of pillows on

the sofa?

Ground truth: red

CNN-LSTM: red (0.4382) BLIND-LSTM: brown (0.1990)

Focus on the only red objects in the im-

age



Q1085: what is the colour of the towels

Ground truth: white CNN-LSTM: white (0.1828)

BLIND-LSTM: white (0.1762)

Towels are white



Q2355: what is the colour of roll of tissue paper? Ground truth: white CNN-LSTM: white (0.4944) BLIND-LSTM: white (0.2917) Hard to see non-white toilet tissue pa-



Q1890: what is the colour of the bed

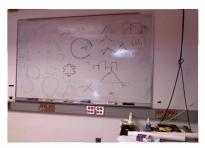
and pillow?

Ground truth: white

CNN-LSTM: brown (0.2249)

BLIND-LSTM: white (0.2566)

Focus to the poster instead of the bed but white is indeed the best blind guess



Q24: what color is the UNK circle on

the whiteboard? Ground truth: green

CNN-LSTM: red (0.2240) BLIND-LSTM: brown (0.1976)

Too hard to answer

Table 5. DAQUAR Color Questions

## Question Answering about Images using Visual Semantic Embeddings



Q681: what color is the chair? Ground truth: black CNN-LSTM: red (0.2508) BLIND-LSTM: brown (0.3813) Sometimes it is oversensitice to red



Q1025: what is the colour of the chair? Ground truth: green CNN-LSTM: red (0.4710) BLIND-LSTM: brown (0.2022) Fails to do a good job on simple one like this

Table 6. DAQUAR Color Questions